

Introduction

Your OS will be managing all of the memory. As you have already probably figured out, we do not have the `sbrk` or `mmap` functions to ask for more memory. Your OS will instead be asked by the OS itself as well as user applications for dynamic memory.

Since your OS will be running in virtual memory, it is important to create a **page-grained memory allocator**. This allocator hands out individual **pages** of memory using virtual memory addresses. Recall that RISC-V has three page sizes: (1) 1GB, (2) 2MB, and (3) 4KB.

Your page-grained allocator will hand out individual or contiguous 4KB pages.

Page Allocator

Edit a file called `page.h`. This file will export the page-grained allocator functions, listed below.

page.h

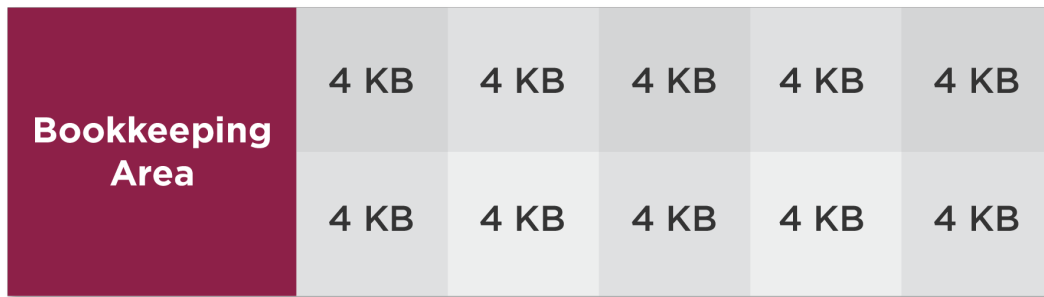
```
1 void page_init(void);
2 void *page_nalloc(unsigned int n);
3 void *page_znalloc(unsigned int n);
4 void page_free(void *p);
5
6 #define page_alloc() page_nalloc(1)
7 #define page_zalloc() page_znalloc(1)
```

Function	Description
<code>page_init</code>	Initializes the page allocator system. Usually, this means setting up the bookkeeping area.
<code>page_nalloc</code>	Allocate a contiguous set of pages given by the parameter <code>n</code> . This function does NOT clear the memory returned.
<code>page_znalloc</code>	Allocate a contiguous set of pages given by the parameter <code>n</code> . This function clears the memory returned to 0.
<code>page_free</code>	Returns the page(s) allocated and returned by <code>p</code> . Recall that the page address in <code>p</code> may be only a single page or multiple pages.

Write your `page_*` functions in a file called `page.c`. The Makefile will automatically compile all C files in `src/` and expect all headers be in `src/include/`

Your page grained allocator will use the top portion of the memory pool (`_heap_start`) to store bookkeeping information, which is two bits: (1) page tak and (2) last page. Therefore, we bookkeep four pages per byte.

Memory Layout for Page-Grained Allocator



- $_heap_start + n$ will be for pages $(4 \times n)$, $(4 \times n + 1)$, $(4 \times n + 2)$, and $(4 \times n + 3)$ $(4 \times \text{? ? ? ?})$, $(4 \times \text{? ? ? ?} + 1)$, $(4 \times \text{? ? ? ?} + 2)$, and $(4 \times \text{? ? ? ?} + 3)$
- $_heap_start + 0$ will be for pages 0, 1, 2, and 3.
- $_heap_start + 1$ will be for pages 4, 5, 6, and 7.

Finding Things in the Heap

1. $_heap_size = _heap_end - _heap_start$ $_heap_size = _heap_end - _heap_start$.

a. This gives us the total number of bytes in the heap.

2. $_num_pages = \frac{_heap_size}{4096}$ $_num_pages = \frac{_heap_size}{4096}$.

a. This is the total number of pages in the heap.

3. $_num_pages \div 4$ $_num_pages \div 4$

a. Since we bookkeep four pages per byte, we divide by four.

4. $_bk_size = \frac{(_heap_end - _heap_start)}{4096 \times 4}$ $_bk_size = \frac{(_heap_end - _heap_start)}{4096 \times 4}$.

a. This calculates the number of bookkeeping bytes needed to manage the heap.

The first page needs to start at a page boundary, so the first page starts at:

$$\text{page}_0 = (_heap_start + \text{ALIGN_UP}(_bk_size, 4096)) \quad \text{page}_0 = (_heap_start + \text{ALIGN_UP}(_bk_size, 4096))$$

Do not forget to set the taken bits for ALL of the bookkeeping bytes. Otherwise, your allocator **may** allocate pages already taken for the book keeping bits.

Bookkeeping Bits

There are two bits per page: (1) taken and (2) last. Since we can allocate multiple, contiguous pages, we need to mark the last page we handed out in a sequence of pages.

TAKEN

LAST

Example

We can visualize how to dole out pages.



You can see that the last bits are 0 for every page in a sequence that is not the last in that sequence.

Recall that `page_nalloc` and `page_free` deal with memory addresses, so you will need to take apart the memory address to find the index in the heap pool.

Memory Example

Let's take, for example, that `_heap_start` is `0xbeef000`, and `_heap_end` is `0xdead000`.

The linker script will ensure that `_heap_start` and `_heap_end` are page aligned. It sets the kernel stack above the heap since the kernel stack is a fixed size.

The full heap is therefore $\text{heap_size} = \text{_heap_end} - \text{_heap_start} = 33284096$ $\text{heap_size} = \text{_heap_end} - \text{_heap_start} = 33284096$ byte:

This means we have $\frac{\text{heap_size}}{4096} = \frac{33284096}{4096} = 8126$ $\frac{\text{heap_size}}{4096} = \frac{33284096}{4096} = 8126$ pages.

To manage 8126 pages, we need $\frac{8126}{4} = 2031.5 = 2032$ $\frac{8126}{4} = 2031.5 = 2032$ bytes.

Recall we need to align the bookkeeping bytes to the next page, so we actually need 4096 bytes (one page) to manage 8126 pages.

Therefore, the first page we can allocate is one page after `_heap_start`, which is $0xbeef000 + 4096 = 0xbef0000$ $0xbeef000 + 4096 = 0xbef0000$

Write short, static functions to test/set/clear the taken and last bits. Also, write short, static functions to calculate the math to map page addresses to the bookkeeping bits.

Virtual Memory

`_heap_start` and `_heap_end` are **physical memory addresses**, since they come from the linker script. After you get the MMU working properly, you need to ensure that your heap is managed **virtually**.

MMU Functions

Your MMU functions need to be aware of the MMU design. For the RISC-V architecture, we are going to use the Sv39 (supervisor, 39-bit virtual addresses) system. Therefore, the following defines/macros will be important.

mmu.h

```

1  #define MMU_LEVEL_1G      2
2  #define MMU_LEVEL_2M      1
3  #define MMU_LEVEL_4K      0
4
5  #define PAGE_SIZE_4K      PAGE_SIZE_AT_LVL(MMU_LEVEL_4K)
6  #define PAGE_SIZE_2M      PAGE_SIZE_AT_LVL(MMU_LEVEL_2M)
7  #define PAGE_SIZE_1G      PAGE_SIZE_AT_LVL(MMU_LEVEL_1G)
8
9  #define PAGE_SIZE_AT_LVL(x) (1 << (x * 9 + 12))
10
11 #define PAGE_SIZE          PAGE_SIZE_4K
12
13 // PB_* - page bits
14 #define PB_NONE            0
15 #define PB_VALID          (1UL << 0)
16 #define PB_READ           (1UL << 1)
17 #define PB_WRITE          (1UL << 2)
18 #define PB_EXECUTE        (1UL << 3)
19 #define PB_USER           (1UL << 4)
20 #define PB_GLOBAL         (1UL << 5)
21 #define PB_ACCESS         (1UL << 6)
22 #define PB_DIRTY          (1UL << 7)
23
24 // MODE[63:60] (4 bits) in the SATP register.
25 #define SATP_MODE_BIT      60
26
27 // SV39 is MODE=8
28 #define MODE_SV39          8UL
29 #define SATP_MODE_SV39     (MODE_SV39 << SATP_MODE_BIT)
30
31 // ASID[59:44] (16 bits) in the SATP register.
32 #define SATP_ASID_BIT      44
33
34 // PPN[43:0] (44 bits) in the SATP register
35 #define SATP_PPN_BIT       0
36 #define SATP_SET_PPN(x)    (((uint64_t)(x)) >> 12) & 0xFFFFFFFFFUL)
37 #define SATP_SET_ASID(x)   (((uint64_t)(x)) & 0xFFFF) << SATP_ASID_BIT)
38
39 #define SATP(table, asid) (SATP_MODE_SV39 | SATP_SET_PPN(table) |
40 SATP_SET_ASID(asid))
41 #define SATP_KERNEL        SATP(kernel_mmu_table, KERNEL_ASID)

```

You will also need to write four functions.

mmu.h

```

1  struct page_table {
2      uint64_t entries[PAGE_SIZE / 8];
3  };
4
5  bool      mmu_map(struct page_table *tab,
6                  uint64_t vaddr,
7                  uint64_t paddr,
8                  uint8_t lvl,
9                  uint64_t bits);
10
11 void      mmu_free(struct page_table *tab);
12 uint64_t mmu_translate(const struct page_table *tab, uint64_t vaddr);
13
14 #define mmu_translate_ptr(tab, ptr) (void *)mmu_translate(tab, (uint64_t)ptr)
15 #define mmu_translate_ptr_to_u64(tab, ptr) mmu_translate(tab, (uint64_t)ptr)
16 #define MMU_TRANSLATE_PAGE_FAULT -1UL
17
18 uint64_t mmu_map_range(struct page_table *tab,
19                       uint64_t start_virt,
20                       uint64_t end_virt,
21                       uint64_t start_phys,
22                       uint8_t lvl,
23                       uint64_t bits);

```

mmu_map

The `mmu_map` function will map the given virtual address to the physical address and create a leaf at the given level. The `bits` needs to be OR'd with the `PB_VALID` bit. This allows the programmer to specify the permission bits, such as `PB_USER` and/or `PB_READ` and so forth. The page table passed will be the root table. The reason it is passed in is because you will need to create mappings for the kernel and any user space applications, so this same function will be used for all page tables.

This function returns `true` if the mapping was made, or `false` otherwise. A `false` can be returned if the parameters don't make sense (e.g., `lvl` is not 0, 1, or 2) or if there is not enough memory to create branch tables.

This function will **overwrite** previous mappings if the same virtual address is provided on the same table.

mmu.c

```

1  bool mmu_map(struct page_table *tab,
2              uint64_t vaddr,
3              uint64_t paddr,
4              uint8_t lvl,
5              uint64_t bits)
6  {
7      int i;
8      // Error check tab, lvl, and bits
9      if (tab == NULL || lvl > MMU_LVL_1GB || (bits & 0xE) == 0) {
10         return false;
11     }
12
13     // Get vpn[0], vpn[1], and vpn[2].
14     const uint64_t vpn[] = {(vaddr >> ADDR_0_BIT) & 0x1FF, (vaddr >> ADDR_1_BIT) &
15 0x1FF,
16                             (vaddr >> ADDR_2_BIT) & 0x1FF};
17
18     // Get ppn[0], ppn[1], and ppn[2].
19     const uint64_t ppn[] = {(paddr >> ADDR_0_BIT) & 0x1FF, (paddr >> ADDR_1_BIT) &
20 0x1FF,
21                             (paddr >> ADDR_2_BIT) & 0x3FFFFFF};
22
23
24     for (i = MMU_LEVEL_1G; i > lvl; i -= 1) {
25         // Go through the branches.
26         // NOTE: you may need to create additional tables.
27     }
28
29     // After the loop, you're looking at the leaf @ i.
30
31     return true;
32 }

```

Recall that C does not have a **bool** data type unless you include `<stdbool.h>`.

The `mmu_map_range` function needs to map a range of addresses, mapping each page to the corresponding physical address. The following code represents an implementation.

mmu.c

```

1  uint64_t mmu_map_range(struct page_table *tab,
2                          uint64_t start_virt,
3                          uint64_t end_virt,
4                          uint64_t start_phys,
5                          uint8_t lvl,
6                          uint64_t bits)
7  {
8      start_virt          = ALIGN_DOWN_POT(start_virt, PAGE_SIZE_AT_LVL(lvl));
9      end_virt            = ALIGN_UP_POT(end_virt, PAGE_SIZE_AT_LVL(lvl));
10     uint64_t num_bytes   = end_virt - start_virt;
11     uint64_t pages_mapped = 0;
12
13     uint64_t i;
14     for (i = 0; i < num_bytes; i += PAGE_SIZE_AT_LVL(lvl)) {
15         if (!mmu_map(tab, start_virt + i, start_phys + i, lvl, bits)) {
16             break;
17         }
18         pages_mapped += 1;
19     }
20     return pages_mapped;
21 }

```

The `mmu_map_range` function returns the number of pages that were properly mapped. This is because any individual `mmu_map` may fail, but we don't want to unwind it. Instead, we will let the programmer decide what to do if all pages aren't mapped.

mmu_free

The `mmu_free` function needs to **recursively** free all of the entries of a given table. Recall that each table could be a branch, which means that the memory address stored in the entry is a page you allocated. All of these pages need to be freed.

mmu.c

```

1  void mmu_free(struct page_table *tab)
2  {
3      uint64_t entry;
4      int i;
5      if (tab == NULL) {
6          return;
7      }
8      // Each entry is 8 bytes, so there are PAGE_SIZE / 8 entries.
9      for (i = 0; i < (PAGE_SIZE / 8); i += 1) {
10         entry = tab->entries[i];
11         // Check if this is a branch, if it is, recurse
12         // to the branch.
13
14         // ALL entries should be cleared to 0 after branches
15         // return back from the recursion.
16         tab->entries[i] = 0;
17     }
18     page_free(tab);
19 }

```


mmu_translate

The `mmu_translate` function will translate a virtual address to a physical address given a table. This will be helpful getting the physical address when worrying about hardware drivers, etc.s

mmu.c

```
1  uint64_t mmu_translate(const struct page_table *tab, uint64_t vaddr)
2  {
3      int i;
4      // Can't translate without a table.
5      if (tab == NULL) {
6          return MMU_TRANSLATE_PAGE_FAULT;
7      }
8      uint64_t vpn[] = {(vaddr >> ADDR_0_BIT) & 0x1FF, (vaddr >> ADDR_1_BIT) &
9      0x1FF,
10                      (vaddr >> ADDR_2_BIT) & 0x1FF};
11
12      // Translate and return the physical address.
13  }
```

Copy To/From

We have made things easier by identity mapping some physical pages with the same virtual address. However, we will eventually need to copy data from/to process which is not identity mapped *and* perhaps the memory addresses are virtually contiguous but not necessarily physically contiguous.

Edit the files called `uaccess.h` (user access) and `uaccess.c` to support the following two functions.

uaccess.h

```
1  #pragma once
2
3  #include <stdint.h>
4
5  struct page_table;
6
7  uint64_t copy_from(void *dst,
8                    const struct page_table *from_table,
9                    const void *from,
10                   uint64_t size);
11
12  uint64_t copy_to(void *to,
13                  const struct page_table *to_table,
14                  const void *src,
15                  uint64_t size);
```

uaccess.c

```

1  #include <uaccess.h>
2  #include <util.h>    // for memcpy
3  #include <mmu.h>     // for struct page_table
4
5  uint64_t copy_from(void *dst,
6                     const struct page_table *from_table,
7                     const void *from,
8                     uint64_t size)
9  {
10     uint64_t bytes_copied = 0;
11
12     // ...
13
14     return bytes_copied;
15 }
16
17 uint64_t copy_to(void *to,
18                 const struct page_table *to_table,
19                 const void *src,
20                 uint64_t size)
21 {
22     uint64_t bytes_copied = 0;
23
24     // ...
25
26     return bytes_copied;
27 }

```

Copy From

```

uint64_t copy_from(void *dst,
                  const struct page_table *from_table,
                  const void *from,
                  uint64_t size);

```

The `copy_from` function copies data from the virtual address in `from` to the virtual address in `dst`. The destination is translated by the MMU, but the source will be translated with `mmu_translate` using `from_table` as the page table. The number of bytes to copy is provided in `size`. This function returns the number of bytes copied.

You need to make sure that the from addresses are properly mapped. The reason you return the number of bytes copied is because you may successfully translate say the first three pages, but the fourth page faults. Therefore, you only copy the bytes from the first three pages.

Do not translate every byte. Instead, copy using `memcpy` until you hit a page boundary, then translate, copy another page, and so forth until you hit `size` or a page fault, whichever comes first.

Recall that `mmu_translate` produces a physical address, but since you identity mapped the physical pool in the kernel page table, you can treat the physical address as a virtual address. Usually, you would have to do a reverse lookup using some data structure like a map.

Copy To

```
uint64_t copy_to(void *to,
                 const struct page_table *to_table,
                 const void *src,
                 uint64_t size);
```

The `copy_to` function is analogous to `copy_from`, except it will copy bytes from the virtual address in `src` to the destination address in `to`. The source address is translated by the MMU, but the `to` address needs to be translated given the passed table in `to_table`. Just like `copy_from`, the `to` memory address may span multiple pages which are not necessarily contiguous.

Also, like `copy_from`, this function will return the total number of bytes copied from the `src` address to the `to` address.

Do not translate every byte. Instead, copy using `memcpy` until you hit a page boundary, then translate, copy another page, and so forth until you hit `size` or a page fault, whichever comes first.

Recall that `mmu_translate` produces a physical address, but since you identity mapped the physical pool in the kernel page table, you can treat the physical address as a virtual address. Usually, you would have to do a reverse lookup using some data structure like a map.

Enabling Code in the Template

In `config.h`, uncomment `USE_MMU` and `USE_HEAP`.

These two defines control certain code in `main.c` that call your `page_init()` function as well as the `heap_init()` function. The `heap_init()` will request a certain number of **continuous** pages from `page_nalloc()` to act as the kernel heap.

You need to get this right. You will be requesting memory for a lot of things, and the heap is the only way to have **persistent** memory. If you don't get this to work properly, many things in your kernel, including the utility library, will not function properly.

You can see in `main.c` that uncommenting `USE_MMU` performs the following actions.

```

1  #ifdef USE_MMU
2      page_init();
3      struct page_table *pt      = page_zalloc();
4      // kernel_mmu_table is global and exported throughout.
5      kernel_mmu_table = pt;
6      // Map memory segments for our kernel
7      mmu_map_range(pt, sym_start(text), sym_end(heap), sym_start(text),
8  MMU_LEVEL_1G,
9                      PB_READ | PB_WRITE | PB_EXECUTE);
10     // PLIC
11     mmu_map_range(pt, 0x0C000000, 0x0C2FFFFFF, 0x0C000000, MMU_LEVEL_2M, PB_READ |
12  PB_WRITE);
13     // PCIe ECAM
14     mmu_map_range(pt, 0x30000000, 0x30FFFFFF, 0x30000000, MMU_LEVEL_2M, PB_READ |
15  PB_WRITE);
16     // PCIe MMIO
17     mmu_map_range(pt, 0x40000000, 0x4FFFFFFF, 0x40000000, MMU_LEVEL_2M, PB_READ |
18  PB_WRITE);
19
20     // TODO: turn on the MMU when you've written the src/mmu.c functions
21     CSR_WRITE("satp", SATP_KERNEL);
22     SFENCE_ALL();
23 #endif

```

Therefore, you need to make sure you have your page grained allocator as well as the mmu mapping functions working properly before uncommenting USE_MMU and USE_HEAP .