

# CS 343 Operating Systems, Fall 2022

## Paging Lab

### Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Setup</b>	<b>2</b>
<b>3</b>	<b>Files</b>	<b>3</b>
<b>4</b>	<b>Address spaces in NK</b>	<b>4</b>
<b>5</b>	<b>x64 Paging</b>	<b>5</b>
<b>6</b>	<b>Task 0: Understand the boot loop</b>	<b>8</b>
<b>7</b>	<b>Task 1: Eager page table construction</b>	<b>8</b>
<b>8</b>	<b>Task 2: Memory map data structure</b>	<b>8</b>
<b>9</b>	<b>Task 3: Lazy page table construction</b>	<b>9</b>
<b>10</b>	<b>Task 4: Fleshing out your implementation</b>	<b>9</b>
<b>11</b>	<b>Task 5: Reflection on your implementation</b>	<b>10</b>
<b>12</b>	<b>Grading</b>	<b>10</b>

## 1 Introduction

The purpose of this lab is to introduce you to virtual memory by implementing virtual address spaces using paging. Paging requires you to think at a deep level about indirection and its management via joint hardware/software mechanisms. In this lab, you will build an implementation of virtual memory using x64 paging within NK's address space abstraction.

You may work in a group of up to three people in this lab. Clarifications and revisions will be posted to the course discussion group.

## 2 Setup

You can work on this lab on any modern Linux system, but our class server, Moore, has been specifically set up for it. This lab will work best on our class server, Moore. Although you can set up your own machine with a special build of QEMU by following a guide on Campuswire. We will describe the details of how to access the lab repo via Github Classroom on Campuswire. You will use this information to clone the assignment repo using a command like this:

```
server> git clone [url]
```

This will give you the entire codebase and history of the Nautilus kernel framework ("NK"), just as in the Getting Started Lab. As before, you may want to use `chmod` to control access to your directory.

**Important!** You will need to make sure you have a valid display for NK to run. You can get that through FastX or with `ssh -Y`. See the Getting Started Lab for more details.

Now build it (you may need to run `bash` first if it is not your default terminal):

```
server> cd [assignment-directory]
server> make clean
server> make -j 8 isoimage
```

You can now boot your kernel:

```
server> ./run
```

The `run` command will execute the emulator (QEMU) with a set of options that are appropriate for the lab.

**Boot failure!** The emulated machine will boot NK. *You will see that the kernel is in a boot loop!* It will try to boot, get to a certain point, then the machine will spontaneously reboot. This is because the shell is trying to create and place itself into a new address space. Unfortunately, paging is mostly unimplemented, so what happens is a switch to a bad page table. Fetching the very next instruction immediately causes a *page fault*, which, to be handled, requires paging, and, with a bad page table, this faults again, this time with a *double fault*. To handle a double fault, the machine once again needs paging, so it faults again. This *triple fault* is handled directly by the hardware, by resetting itself. Hence the boot loop. Three strikes and you are out.

### 3 Files

While NK is tiny compared to the Linux, Darwin, or Windows kernels, it does have several hundred thousand lines of code spread over over more than a thousand files. Therefore it is important to focus on what is important for your goals. As with any significant codebase, trying to grok the whole thing is either impossible or will take far too long. Your strategy for approaching the code has to be adaptive. In part, we are throwing you into this codebase to help you learn how to do this.

Here are some important files for this lab. Your edits will be in `src/nautilus/shell.c` and `src/aspace/paging/paging.c` (bolded below):

- `include/nautilus/aspace.h`: This is NK's address space abstraction. You will be creating an address space implementation that conforms to it. You will not change this.
- `src/nautilus/aspace.c`: This is the implementation of NK's address space abstraction. You will not change this.
- **`src/nautilus/shell.c`**: This is the shell implementation, which uses the address space abstraction. It's where you can test things out. The boot loop is occurring as a result of the call to `nk_aspace_move_thread()` called in the `shell()` function. The surrounding code shows how the address space abstraction is used. This is your first test!
- `src/asm/thread_lowlevel.S`: The call instruction to `nk_aspace_switch` in this code is what does a possible address space switch when a different thread is scheduled. You do not need to change this file.
- `include/nautilus/thread.h`: The field `aspace` within `struct nk_thread` points to the address space the thread is in. If this is null, it means the thread is in the default (or boot) address space. You do not need to change this file.
- `include/nautilus/smp.h`: The field `cur_aspace` within `struct cpu` points to the currently active address space for the CPU (the hardware thread). If this is null, it means the CPU is in the default (or boot) address space. You do not need to change this file.
- `include/nautilus/paging.h`: This includes helpful definitions that can be used in paging implementations. For example, `PAGE_SIZE_4KB` can be used in your code instead of remembering how many bytes a standard 4 kB page is. You do not need to change this file.
- **`src/aspace/paging/paging.c`** This is the stub source code for your paging address space implementation. It is heavily commented. You will add to this.
- `src/aspace/paging/paging_helpers.[ch]` These files contain heavily commented helper code for building 4-level x64 page tables. You can leverage this helper code or write your own. You might find the `paging_helper_walk()` and `paging_helper_drill()` utility functions helpful. You are welcome to add to these files if desired.
- **`src/aspace/paging/paging_test.c`** This file contains additional test code, which you can run using the shell command `paginptest`. You are welcome to add additional tests. Also included are several commented-out tests which *should* fail if your implementation is correct. When you believe your implementation is complete, you should bring them back, one at a time, and check that they fail as expected.

Note that we are pointing out a lot of different files above. This is to show you how deeply embedded the notion of a virtual memory tends to be in a kernel, and to be complete. In this lab, you will mostly be working in `src/aspace/paging/paging.c`, which is heavily commented to help you.

## 4 Address spaces in NK

NK is somewhat different than the general purpose kernels, such as Linux, Windows, MacOS, BSD, etc, described in the book, as well as from the microkernels your author likes. In particular:

- The use of virtual memory is *optional* in NK. Using physical memory directly (or as close as possible) is the common case.
- There is no kernel/user distinction in NK by default.
- There is no process abstraction in NK by default.
- There is an address space abstraction designed to allow the use of models of virtual memory that do *not* involve paging as well as those that do. The address space abstraction is optional.

Even if you're not using them, on x64 hardware, page tables *must* be installed. When NK boots, it builds a page table hierarchy that does an *identity map*, meaning that every virtual address maps to exactly the same physical address, with full kernel privileges. NK implements this page table hierarchy using the largest possible pages. Within the address space abstraction, this forms the *default address space*. Everything lives within this single address space unless it chooses otherwise.

The *address space abstraction* (`include/nautilus/aspace.h`) allows the creation and management of additional address spaces. A thread can choose to join an address space (`nk_aspace_move_thread()`). When a new thread is spawned, it joins the address space of its parent. Each CPU has a current address space, that of the currently running thread. Interrupt handlers and the rest of the kernel run in the context of the current address space of the CPU. You can find the current list of all address spaces on the system using the `ases` shell command.

An address space is implemented by a *address space implementation*, and a design goal here is to allow very unusual implementations that manage address spaces at arbitrary granularities, and to allow address spaces from multiple implementations to coexist at runtime. You are writing an address space implementation based on x64 4-level paging. You can get a list of all the available address space implementations with the `asis` shell command.

The address space abstraction centers around a *region* (`nk_aspace_region_t`), which is a mapping from a virtual address to a physical address for some number of bytes (not pages), combined with a set of protection flags. An address space contains a set of regions, and the set constitutes the address space's *memory map*. The memory map is an implementation-independent representation of the virtual address to physical address mapping.<sup>1</sup> In this lab, you will implement this mapping using paging.

Once an address space is created, it can be manipulated using regions:

- *Add region*: This expands the memory map. If the region is *eager*, then this part of the memory map must be immediately implemented by the underlying mechanism. For example, you would need to build matching page table entries in this lab. On the other hand, the page table entries for a lazy region could be built on demand.

---

<sup>1</sup>Note that this is different from the concept of memory map in Linux, which instead provides a mapping of virtual memory regions to file sections.

- *Remove region*: This shrinks the memory map. For paging you need to be sure that any page table entries you created for the region are also deleted. Page table entries may be cached in the TLB, so you also need to flush them from the TLB.
- *Move region*: The idea here is that we are changing the virtual to physical mapping of a region in the memory map. The virtual address stays the same, but the physical address changes. Similar to removing a region, you need to assure that old page table entries are edited, and that old entries are flushed from the TLB.
- *Protect region*: Here, we are changing the protections of an existing region. The virtual and physical addresses stay the same, but the protections change. Similar to moving or removing a region, you need to edit page table entries and flush the old ones from the TLB.<sup>2</sup>

Your address space also needs to handle the following requests:

- *Switch from*: This is invoked when your address space is about to stop being the current address space for the CPU. For paging, there is little you probably need to do.
- *Switch to*: This is invoked when your address space is about to become the current address space for the CPU. For paging, you need to install your page tables at this point.
- *Exception*: This is invoked in interrupt context whenever a page fault or general protection fault is encountered. For paging, you might build a page table entry for a lazy region at this point.<sup>3</sup>
- *Add thread*: This is invoked when a thread is joining the address space. For paging, you probably don't care.
- *Remove thread*: This is invoked when a thread is leaving the address space. For paging, you probably don't care.
- *Print*: Display the details of the address space. This supports the `ases` shell command.
- *Destroy*: The address space will no longer be used, and you should free all of its state.

## 5 x64 Paging

You will implement 4-level x64 paging. This is the most basic form of paging used on x86 (Intel/AMD/etc) processors when running in 64-bit mode (“long mode”). Please note that many teaching OS examples you might find (i.e. xv6, Pebbles, GeekOS, etc) use 32-bit mode, where paging is substantially different.

Three references for 64 bit paging that you should be aware of are the following:

- *CS 213 Textbook*: R. Bryant, D. O’Hallaron, Computer Systems: A Programmer’s Perspective, 3rd edition, Section 9.7, shows the big picture of this kind of machine.

---

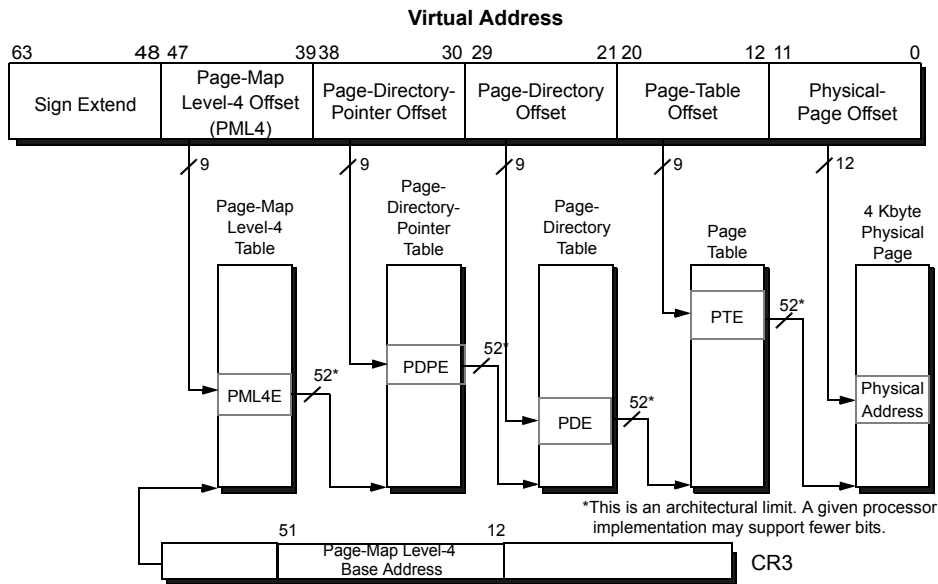
<sup>2</sup>Technically, whenever you flush entries from your CPU’s TLB, you also need to make sure they are flushed from the other CPUs’ TLBs. Typically, this is done using an inter-processor interrupt (IPI) called a TLB shutdown.

<sup>3</sup>If the reason for the page fault is unfixable by the kernel, then you should panic. In a kernel with a user space, if the page fault originated in user space, you would instead inject a signal (SIGSEGV (i.e., segfault) on Unix-like systems) into the user space program. If it originated in the kernel, you would panic.

- *AMD Documentation:* AMD64 Architecture Programmer’s Manual Volume 2: System Programming, Chapter 5, and in particular 5.3, describe paging on this form of architecture.
- *Intel Documentation:* Intel 64 and IA-32 Architectures Software Developer’s Manual: Volume 3, Chapter 4, and, in particular 4.5 (“IA32e” is Intel’s name for 64 bit mode)

The material, particularly the Intel documentation, can be quite daunting. In part this is because it is explaining all of the various modes and aspects of the machine tied to paging all at once. We are looking for you to build one thin slice through this. Also keep in mind that we have implemented a lot of support for you.

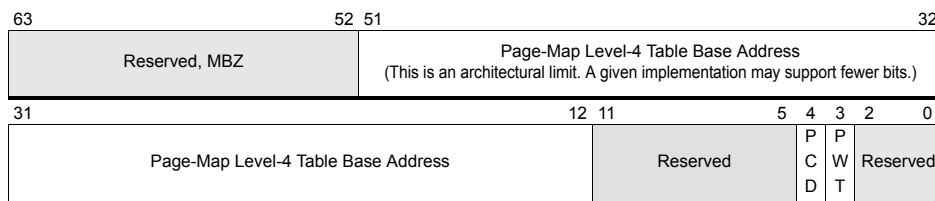
The following figure<sup>4</sup> shows how a virtual address is translated into a physical address by the hardware:



**Figure 5-17. 4-Kbyte Page Translation—Long Mode**

Here, the page table hierarchy we will use is selected by a pointer stored in the CR3 register—this pointer points to the root page table. The first group of 9 bits in the address are used to select one of the 512 entries on this table. The entry contains a pointer to the next level page table. The next group of 9 bits in the address are used to select an entry within it. And so on, all the way down to the last level page table, where the pointer indicates the physical page that corresponds to the virtual address.

The CR3 register has a special format:



**Figure 5-16. Control Register 3 (CR3)—Long Mode**

as do the page table entries at the different levels:

<sup>4</sup>Figures are from the AMD documentation unless otherwise noted.



- `excp_entry_t`, field `error_code`: On a page fault, this value is the reason for the page fault (why it occurred). See the stub code in for more.
- `invlpg()` Given a virtual address, this function flushes the corresponding entry out of the TLB, if it's in the TLB.

## 6 Task 0: Understand the boot loop

Start by walking through the code, starting with the `aspace`-related calls that occur in `src/nautilus/shell.c`. Your goal is to understand for yourselves why the boot loop is occurring. You might want to try commenting out the `nk_aspace_move_thread()` call to see what happens. Then bring it back in, but add incremental printouts so that you can see how far it gets. Possibly also attach `gdb` and use it to single-step through this processing.

Note that the abstraction creates an indirection that can make it hard to follow. In this case, we can assure you that a call to a function like `nk_aspace_add_region()` will result in call to the `add_region()` function within the paging implementation. The other high-level calls will also route to the similarly named functions in the paging implementation.

## 7 Task 1: Eager page table construction

Your next step is to stop the boot loop by building a page table hierarchy that has the necessary parts of the virtual address space mapped. Note that the test code in the shell adds two regions. The first of these is an “eager” region. This means you should build page tables for it immediately, right in your `add_region()` function. You will find the `paging_helper_drill()` function useful for this.

Drilling a page builds the page table hierarchy that represents the translation for that page. Multiple pages will end up sharing parts of their hierarchy if they are adjacent. Drilling requires a virtual and physical address for the page as well as a `ph_pf_access_t` union that specifies permissions for the page. The union is defined in `src/aspace/paging/paging_helpers.c` and is identical to `ph_pf_error_t`. In an `access_type` the bits are used to represent the permissions granted to the page.

Once you correctly construct the page table entries corresponding to this eager region, your boot loop should go away. However, the second region we ask you to add in the shell code is not an eager region. You should not build it eagerly. As a consequence, you will now likely get a repeated page fault. This is due to the `memcmp()` test in the shell code. If you comment this out, you should boot all the way to the shell prompt. At this point, you should be able to run the `ases` command or `threads` command and see that your shell is running in a new address space.

## 8 Task 2: Memory map data structure

To continue the lab, you will need a way to manage the regions that comprise your memory map. Recall that the memory map is a set of regions. You need to design and implement a data structure that contains an set of regions. You need to be able to add and remove regions from the data. Regions may not overlap by virtual addresses. That is, a virtual address must map to at most one region. On the other hand, a physical address can map to multiple regions. The data structure also needs to be searchable by virtual address. That



is, given any virtual address (not just the start of a region), you need to find the region that contains it, if any exists.

This does not have to be fancy since we will not grade you on performance. A simple linked list can work fine. Most students create their own data structure, which is totally fine. You could even add new `.c` and `.h` files inside the paging directory and add them to the `Makefile`. For now, keep the implementation simple, and you can add features as needed when working on later tasks.

Note that your data structure should hold region structs that are copies of the passed in region data, rather than storing the region pointer itself. The region pointer passed into the `add_region()` function may not persist past the return of the function or may be modified in the future.

We also provide a Linux-style intrusive linked list implementation (`include/nautilus/list.h`) if you'd like to leverage that. Here is a guide with more details on intrusive linked lists: <https://www.data-structures-in-practice.com/intrusive-linked-lists/>

## 9 Task 3: Lazy page table construction

Now that you can keep track of your memory map, you can start handling page faults, and possibly construct new page table entries based on them.

Enable the second region in the shell test code (the lazy one), and enable the `memcmp()` test. The region should now be included in your memory map, but, at least initially, you won't have any page table entries to support it. As a consequence, the `memcmp()` will cause a page fault when it reads from the high address.

Your page fault handler implemented in `exception()` can now do something about this. It should get the faulting virtual address and error, and then look up the virtual address in the memory map. If a region exists, and the permissions of the region are appropriate given the error, then drill a page table entry for the address, with the physical address corresponding to the region.

To determine if the permissions of the region are valid, you should look at the `ph_pf_error_t` created from the value passed into `exception()`. The union is defined in `src/aspace/paging/paging_helpers.c` and each field provides characteristics of the access that faulted. You can compare those with the permissions of the region to determine if they may be problematic.

Once you have this right, you will survive the boot process all the way to the shell prompt. You will also see numerous page faults during the boot process, all of which are satisfied by your lazy page table construction logic.

## 10 Task 4: Fleshing out your implementation

Finally, you will expand your implementation to support deleting and moving regions, as well as changing their protections. Note that applying these changes may be different depending on whether the region has actually affected the page tables due to being an eager region or at least one page fault having already occurred.

Regions that are removed, but have already been drilled in the page table hierarchy, should be marked as not present. A full implementation would scan the page table hierarchy and free branches that are entirely not present, but you do not need to implement that logic for this lab.

Regions are only moved from one physical address to another. This could represent memory being placed in a new operation in RAM after a swap to disk and back. Region virtual addresses will never change.

Regions can be created with or modified to have arbitrary permissions. Any combination of writeable and/or executable is possible and should be enforced, but regions will always be readable. Regions can also be “pinned”. Pinned regions cannot be moved or removed unless they are first unpinned. This prevents important memory regions from being swapped to disk.

Before testing write permissions on an address space, you will need to enable write protections in the kernel. By default, write permissions are normally ignored when in kernel mode. To enable them add the following line to your testing code once, sometime before calling `nk_aspace_protect_region()`.

```
write_cr0(read_cr0() | (1<<16));
```

We already do this for you as an example in the `pagingtest` code (`src/aspace/paging/paging_test.c`).

Generally, the rest of the paging logic is straightforward once you can handle lazy page table construction because it depends similarly on the memory map data structure. However, read carefully the earlier comments about TLB invalidation. When a page table entry can be cached in a TLB, it is important to make sure it is removed from the TLB after you edit the in-memory copy.

## 11 Task 5: Reflection on your implementation

Answer the following questions about your implementation and how it functions. Put the answers in your `STATUS` file. Don't feel like you need to spend pages answering these. Simple, concise answers are greatly preferred. For many of these questions there is no correct answer, only an answer that is how it would apply to your implementation.

1. Explain what your data structure for the memory map / region set is. What are the costs to insert/remove/change/search for regions?
2. Explain how your implementation handles the following situations/questions:
  - (a) `add_region` that has a virtual memory overlap with an existing region in the memory map
  - (b) `add_region` that has a physical memory overlap with an existing region in the memory map
  - (c) `move_region` on the current thread's address space where the move would end up causing a physical memory overlap. This is for a move involving multiple pages.
  - (d) `protect_region` on a non-existent region. How do you find out it is non-existent?
  - (e) when is it necessary to flush the whole TLB (move to `cr3`)?
  - (f) when is it necessary to flush a single page from the TLB (`invlpg`)?
  - (g) what happens if a valid `delete_region` request is for a region that contains `%rip`? What will happen after the paging library code completes?

## 12 Grading

Your group should regularly push commits to Github. You also should create a file named `STATUS` in which you regularly document (and push) what is going on, todos, what is working, etc. Your commits are visible to us, but not to anyone else outside of your group. The commits that we see up to deadline will constitute

your hand-in of the code. The `STATUS` file should, at that point, clearly document that state of your lab (what works, what doesn't, etc).

In addition to your `STATUS` file, you should regularly push your work within `src/ospace/paging/*`, `src/nautilus/shell.c`, and all other files you are changing.

The breakdown in score will be as follows:

- 20% Task 1—Functional and sensible implementation of eager page table construction. With `memcmp()` disabled, it should boot to the kernel prompt. Furthermore, creating additional shells should create additional address spaces that work correctly.
- 20% Task 2—Sensible implementation of a memory map data structure. It needs to provide `add`, `remove`, `change`, and `lookup` functionality. The `lookup` needs to work for arbitrary addresses.
- 20% Task 3—Functional and sensible implementation of lazy page table construction. With `memcmp()` enabled, it should still boot to the kernel prompt. Creating additional shells should work as before.
- 30% Task 4—Fleshed out implementation. There should be support for removing, moving, and changing the protection of regions. Matching testcases implemented in the shell code or `pagingtest`, are expected.
- 10% Task 5—Reflection on your implementation. Put the answers in your `STATUS` file please.