

# CSE 506: Lab 1: x86 Assembly and Bootloader

Home	Schedule	Syllabus	Labs	Tools	Reference	Announcements	Mailing List	
------	----------	----------	------	-------	-----------	---------------	--------------	--

Due 11:59 PM, Wednesday, September 10, 2014

## Introduction

This lab is split into three parts. The first part concentrates on getting familiarized with x86 assembly language, the QEMU x86 emulator, and the PC's power-on bootstrap procedure. The second part examines the boot loader for our CSE 506 kernel, which resides in the boot directory of the lab tree. Finally, the third part delves into the initial template for our CSE 506 kernel itself, named JOS, which resides in the kernel directory.

## Software Setup

The files you will need for this and subsequent lab assignments in this course are distributed using the [Git](#) version control system.

Git is a powerful, but tricky, version control system. We highly recommend taking time to understand git so that you will be comfortable using it during the labs. We recommend the following resources to learn more about git:

1. [Understanding git conceptually](#) This is a MUST READ if you want to work on git smoothly. (You may skip the last part: Rebasing, for now)
2. [Quick 15-20 mins online exercise to get to know git.](#)
3. [Git user's manual](#)
4. If you are already familiar with other version control systems, you may find this [CS-oriented overview of Git](#) useful.

Each student (enrolled and on the waiting list) will be given a virtual machine on the OS teaching cluster with the basic required software installed. You will have root access to this machine, and be able to install additional tools (editors, debuggers, etc) as you see fit.

You are also welcome to install the needed software on your own laptop. The course staff is not available to help you debug your personal laptop configuration. The [tools page](#) has directions on how to set up qemu and gcc for use with JOS.

## VM Setup

Each student currently enrolled or on the waiting list will be assigned a VM. Newly enrolled students will be assigned a VM within a few days of enrolling.

To get started with your VM, you must do some system setup. You must initially setup your VM via the vSphere client; after setup is complete you can use ssh to access your VM.

Use a remote desktop client to connect to connect to ts1.cs.stonybrook.edu:22. Use **CS\userid** to log in, using the same username and password you use for department email.

Start the vSphere client. Select the vm server assigned to you, eg. **esx6sc.cs.stonybrook.edu** and again use your CS department credentials (student VMs are distributed across esx6--9). Ignore the SSL warning. You will need to try each server to identify your VM, which should be named: cse506-{USERID}, where USERID is your netid.

Once you have logged in, click on "Inventory", then select your VM from the list. Right click on it, and select "Open Console." Click the green "Play" button to start the VM. You can log in with the provided account and initial password. **Immediately change your password!**

After booting, the VM will be available via an ssh client on port 130 (use the -p option to specify a port, or add an entry to your .ssh/config file.).

**Warning:** It is highly advisable that you take a checkpoint of your VM at this point, especially if you intend to modify packages on the VM. This allows you to rollback the VM to a working state without CS department administrator help (which is not available on the weekends or late at night) if something goes horribly wrong. The button to take a checkpoint in vSphere is next to the button you used to power on the VM.

**Good citizenship.** You have administrator access to this VM and can install anything you like. That said, to keep load down, DO NOT INSTALL A GRAPHICAL DESKTOP on the VM. You may tunnel the X protocol (ssh -X) to your local machine and display your editor in a window.

## Picking your group

You may do the lab alone, or in pairs. If you work as a pair, you will continue working together all semester.

**Please email the instructor your group preference as soon as possible.** Once we have your group membership, we will create a git repository, which you will use to hand in the assignment. If you choose to work alone, please email this to the instructor.

## Getting started with git

Once you have turned in your group selection, your group will be assigned a private git repository initially populated with the lab 1 JOS skeleton code. You will use this repository to turn in your code (see below). To download the files into your development environment, you need to *clone* the course repository, by running the commands below. Substitute your group id GROUP and your CS netid for USER below.

```
kernit% mkdir ~/cse506
```

```

kermit% cd ~/cse506
kermit% chmod 0700 . # (sets appropriate permissions)
kermit% git clone ssh://USER@scm.cs.stonybrook.edu:130/scm/cse506git-f14/hw-GR0UP lab

Initialized empty Git repository in ...../cse506/lab/.git/
got f6ec6e08634de9b9c4d73ab5af92da16cc610f44
walk f6ec6e08634de9b9c4d73ab5af92da16cc610f44
got a8d9dd484df67d928a51127ce4c6d9f6d01c5a6a
...
got c9dab101498914dbdce377b89a6eb0f6a421d018
Checking out files: 100% (44/44), done.
kermit% cd lab
kermit%

```

Notice that each group will have their own git repository hosted on scm; your group ID should substituted in the URL above, and you will use your CS department account to log into scm. If you are in the course (not on the waitlist) and do not have a repository, please email the course staff with your group selection.

For students that do not have a private git repository yet, you may use the read-only git repository on the course webpage, using the alternate command below:

```

kermit% git clone http://www.cs.stonybrook.edu/~porter/courses/cse506/f14/jos.git lab

```

Git allows you to keep track of the changes you make to the code. For example, if you are finished with one of the exercises, and want to checkpoint your progress, you can *commit* your changes by running:

```

kermit% git commit -am 'my solution for lab1 exercise9'
Created commit 60d2135: my solution for lab1 exercise9
1 files changed, 1 insertions(+), 0 deletions(-)
kermit%

```

You can keep track of your changes by using the **git diff** command. Running **git diff** will display the changes to your code since your last commit, and **git diff origin/lab1** will display the changes relative to the initial code supplied for this lab. Here, **origin/lab1** is the name of the git branch with the initial code you downloaded from our server for this assignment.

We have set up the appropriate compilers and simulators for you on the CS lab machines.

## Hand-In Procedure

Labs will be handed in using the **make handin** command. This creates a tag in git and pushes the tag and changes to the source repository. You must commit all changes you want included in the handin.

When you are ready to hand in your lab, add an entry to `slack.txt` noting how many late hours you have used both for this assignment. (This is to help us agree on the number that you have used.) Then run **make handin** in the `lab` directory. *If you submit multiple times, we will take the latest submission and count late hours accordingly.*

In this and all other labs, you may complete challenge problems for extra credit. If you do this, please add entries to the file called `challenge1.txt`, which includes a short (e.g., one or two paragraph) description of what you did to solve your chosen challenge problem and how to test it. If you implement more than one challenge problem, you must describe each one. Be sure to list the challenge problem number.

*If you submit multiple times, we will take the latest submission and count late hours accordingly.*

You do not need to turn in answers to any of the questions in the text of the lab. (Do answer them for yourself though! They will help with the rest of the lab.) **Note, lab questions may reappear as exam questions!**

We will be grading your solutions with a grading program. You can run **make grade** to test your solutions with the grading program.

## Migrating to your private repository

If you initially cloned the read-only (http) repository, you must update your git configuration. First, do a **git pull** to pull any changes from the read-only source. Once you have a private repository for handing in the labs, you should issue this command to add the handin repository (substituting your user and group IDs appropriately):

```

git remote add handin ssh://USER@scm.cs.stonybrook.edu:130/scm/cse506git-f14/hw-GR0UP

```

You can double check the **.git/config** file to verify the update was successful. Initially, the config will look something like this:

```

[core]
    repositoryformatversion = 0
    filemode = true
    bare = false
    logallrefupdates = true
[remote "origin"]

```

```

fetch = +refs/heads/*:refs/remotes/origin/*
url = http://www.cs.stonybrook.edu/~porter/courses/cse506/f14/jos.git
[branch "lab1"]
    remote = origin
    merge = refs/heads/lab1

```

There should now be a new remote entry like this:

```

[core]
    repositoryformatversion = 0
    filemode = true
    bare = false
    logallrefupdates = true
[remote "origin"]
    fetch = +refs/heads/*:refs/remotes/origin/*
    url = http://www.cs.stonybrook.edu/~porter/courses/cse506/f14/jos.git
[branch "lab1"]
    remote = origin
    merge = refs/heads/lab1
[remote "handin"]
    fetch = +refs/heads/*:refs/remotes/handin/*
    url = ssh://USER@scm.cs.stonybrook.edu:130/scm/cse506git-f14/hw-GROUP

```

## Set your umask

If you are working in a group, you also need to log in to scm and adjust your umask. The default settings on scm (0022) will cause files in your repository to be created without group write permission. You need to change your umask to (0002) so that files on scm are writable by your partner.

Add an entry to the `.bashrc` file in your home directory so that it looks like this:

```

# .bashrc

# Source global definitions
if [ -f /etc/bashrc ]; then
    . /etc/bashrc
fi

# User specific aliases and functions
umask 0002

```

## Part 1: PC Bootstrap

The purpose of the first exercise is to introduce you to x86 assembly language and the PC bootstrap process, and to get you started with QEMU and QEMU/GDB debugging. You will not have to write any code for this part of the lab, but you should go through it anyway for your own understanding and be prepared to answer the questions posed below.

### Getting Started with x86 assembly

If you are not already familiar with x86 assembly language, you will quickly become familiar with it during this course! The [PC Assembly Language Book](#) is an excellent place to start. Hopefully, the book contains mixture of new and old material for you.

**Warning:** Unfortunately the examples in the book are written for the NASM assembler, whereas we will be using the GNU assembler. NASM uses the so-called *Intel* syntax while GNU uses the *AT&T* syntax. While semantically equivalent, an assembly file will differ quite a lot, at least superficially, depending on which syntax is used. Luckily the conversion between the two is pretty simple, and is covered in [Brennan's Guide to Inline Assembly](#).

**Exercise 1.** Read or at least carefully scan the entire [PC Assembly Language](#) book, except that you should skip all sections after 1.3.5 in chapter 1, which talk about features of the NASM assembler that do not apply directly to the GNU assembler. You may also skip chapters 5 and 6, and all sections under 7.2, which deal with processor and language features we won't use. This reading is useful when trying to understand assembly in JOS, and writing your own assembly. If you have never seen assembly before, read this book carefully.

Also read the section "The Syntax" in [Brennan's Guide to Inline Assembly](#) to familiarize yourself with the most important features of GNU assembler syntax. JOS uses the GNU assembler.

We will be developing JOS for the 64-bit version of the x86 architecture (also known as amd64). The assembly is very similar to 32-bit, with a few key differences. Read [this guide](#), which explains the key differences between the assembly.

Become familiar with inline assembly by writing a simple program. Modify the program [ex1.c](#) to include inline assembly that increments the value of `x` by 1. **Add this file to your lab directory so that it is turned in for grading with the rest of your code.**

Certainly the definitive reference for x86 assembly language programming is Intel's instruction set architecture reference, which you can find on [the CSE 506 reference page](#) in two flavors: an HTML edition of the old [80386 Programmer's Reference Manual](#), which is much shorter and easier to

navigate than more recent manuals but describes all of the x86 processor features that we will make use of in CSE 506; and the full, latest and greatest [Intel 64 and IA-32 Combined Software Developer's Manuals](#) from Intel, covering all the features of the most recent processors that we won't need in class but you may be interested in learning about. An equivalent (but even longer) set of manuals is [available from AMD](#), which also covers the new 64-bit extensions now appearing in both AMD and Intel processors.

You should read the recommended chapters of the PC Assembly book, "The Syntax" section in Brennan's Guide, and the gentle introduction to AMD 64 now. Save the Intel/AMD architecture manuals for later or use them for reference when you want to look up the definitive explanation of a particular processor feature or instruction.

## Simulating the x86

Instead of developing the operating system on a real, physical personal computer (PC), we use a program that faithfully emulates a complete PC: the code you write for the emulator will boot on a real PC too. Using an emulator simplifies debugging; you can, for example, set break points inside of the emulated x86, which is difficult to do with the silicon-version of an x86.

In CSE 506 we will use the [QEMU Emulator](#), a modern and relatively fast emulator. While QEMU's built-in monitor provides only limited debugging support, QEMU can act as a remote debugging target for the [GNU debugger](#) (GDB), which we'll use in this lab to step through the early boot process.

To get started, extract the Lab 1 files into your own directory as described above in "Software Setup", then type **make** in the lab directory to build the minimal CSE 506 boot loader and kernel you will start with. (It's a little generous to call the code we're running here a "kernel," but we'll flesh it out throughout the semester.)

```

kermit% cd lab
kermit% make
+ as kern/entry.S
+ cc kern/init.c
+ cc kern/console.c
+ cc kern/monitor.c
+ cc kern/printf.c
+ cc lib/printfmt.c
+ cc lib/readline.c
+ cc lib/string.c
+ ld obj/kern/kernel
+ as boot/boot.S
+ cc -Os boot/main.c
+ ld boot/boot
boot block is 414 bytes (max 510)
+ mk obj/kern/kernel.img

```

(If you get errors like "undefined reference to `\_\_udivdi3'", you probably don't have the 32-bit gcc multilib. If you're running Debian or Ubuntu, try installing the gcc-multilib package.)

Now you're ready to run QEMU, supplying the file `obj/kern/kernel.img`, created above, as the contents of the emulated PC's "virtual hard disk." This hard disk image contains both our boot loader (`obj/boot/boot`) and our kernel (`obj/kern/kernel`).

```

kermit% make qemu

```

This executes QEMU with the options required to set the hard disk and direct serial port output to the terminal. (You could also use **make qemu-nox** to run QEMU in the current terminal instead of opening a new one.)

Some text should appear in the QEMU window:

```

Booting from Hard Disk...
6828 decimal is XXX octal!
entering test_backtrace 5
entering test_backtrace 4
entering test_backtrace 3
entering test_backtrace 2
entering test_backtrace 1
entering test_backtrace 0
leaving test_backtrace 0
leaving test_backtrace 1
leaving test_backtrace 2
leaving test_backtrace 3
leaving test_backtrace 4
leaving test_backtrace 5
Welcome to the JOS kernel monitor!
Type 'help' for a list of commands.
K>

```

Everything after 'Booting from Hard Disk...' was printed by our skeletal JOS kernel; the `K>` is the prompt printed by the small *monitor*, or interactive control program, that we've included in the kernel. These lines printed by the kernel will also appear in the regular shell window from which you ran QEMU. This is because for testing and lab grading purposes we have set up the JOS kernel to write its console output not only to the virtual VGA display (as seen in the QEMU window), but also to the simulated PC's virtual serial port, which QEMU outputs to its own standard output because of the `-serial` argument. Likewise, the JOS kernel will take input from both the keyboard and the serial port, so you can give it commands in either the VGA display window or the terminal running QEMU.

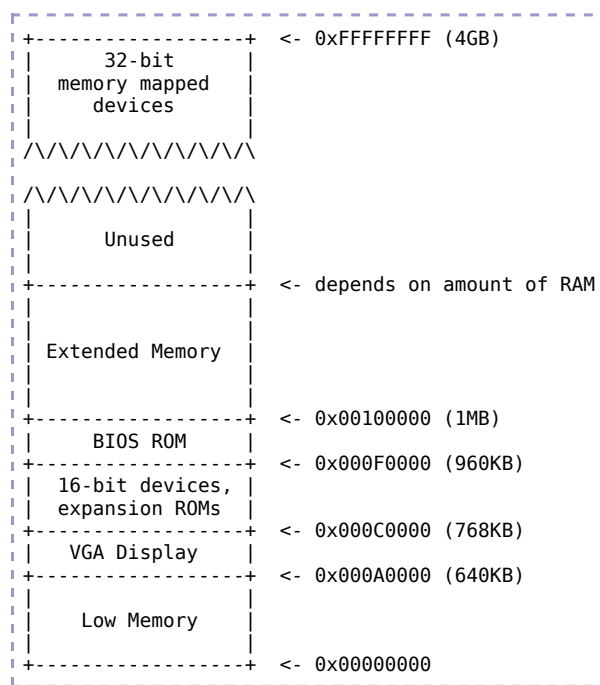
There are only two commands you can give to the kernel monitor, `help` and `kerninfo`. Make sure you type them (and all other input to JOS) into the VGA display window, not into the xterm running QEMU.

```
K> help
help - display this list of commands
kerninfo - display information about the kernel
K> kerninfo
Special kernel symbols:
_start f010000c (virt) 0010000c (phys)
etext f0101a75 (virt) 00101a75 (phys)
edata f010f320 (virt) 0010f320 (phys)
end f010f980 (virt) 0010f980 (phys)
Kernel executable memory footprint: 63KB
K>
```

The `help` command is obvious, and we will shortly discuss the meaning of what the `kerninfo` command prints. Although simple, it's important to note that this kernel monitor is running "directly" on the "raw (virtual) hardware" of the simulated PC. This means that you should be able to copy the contents of `obj/kern/kernel.img` onto the first few sectors of a *real* hard disk, insert that hard disk into a real PC, turn it on, and see exactly the same thing on the PC's real screen as you did above in the QEMU window. (We don't recommend you do this on a real machine with useful information on its hard disk, though, because copying `kernel.img` onto the beginning of its hard disk will trash the master boot record and the beginning of the first partition, effectively causing everything previously on the hard disk to be lost!)

## The PC's Physical Address Space

We will now dive into a bit more detail about how a PC starts up. A PC's physical address space is hard-wired to have the following general layout:



The first PCs, which were based on the 16-bit Intel 8088 processor, were only capable of addressing 1MB of physical memory. The physical address space of an early PC would therefore start at `0x00000000` but end at `0x000FFFFF` instead of `0xFFFFFFFF`. The 640KB area marked "Low Memory" was the *only* random-access memory (RAM) that an early PC could use; in fact the very earliest PCs only could be configured with 16KB, 32KB, or 64KB of RAM!

The 384KB area from `0x000A0000` through `0x000FFFFF` was reserved by the hardware for special uses such as video display buffers and firmware held in non-volatile memory. The most important part of this reserved area is the Basic Input/Output System (BIOS), which occupies the 64KB region from `0x000F0000` through `0x000FFFFF`. In early PCs the BIOS was held in true read-only memory (ROM), but current PCs store the BIOS in updateable flash memory. The BIOS is responsible for performing basic system initialization such as activating the video card and checking the amount of memory installed. After performing this initialization, the BIOS loads the operating system from some appropriate location such as floppy disk, hard disk, CD-ROM, or the network, and passes control of the machine to the operating system.

When Intel finally "broke the one megabyte barrier" with the 80286 and 80386 processors, which supported 16MB and 4GB physical address spaces respectively, the PC architects nevertheless preserved the original layout for the low 1MB of physical address space in order to ensure backward compatibility with existing software. Modern PCs therefore have a "hole" in physical memory from `0x000A0000` to `0x00100000`, dividing RAM into "low" or "conventional memory" (the first 640KB) and "extended memory" (everything else). In addition, some space at the very top of the PC's 32-bit physical address space, above all physical RAM, is now commonly reserved by the BIOS for use by 32-bit PCI devices.

Recent x86 processors can support *more* than 4GB of physical RAM, so RAM can extend further above `0xFFFFFFFF`. In this case the BIOS must arrange to leave a *second* hole in the system's RAM at the top of the 32-bit addressable region, to leave room for these 32-bit devices to be mapped.

Because of design limitations JOS will use only the first 256MB of a PC's physical memory anyway, so for now we will pretend that all PCs have "only" a 32-bit physical address space. But dealing with complicated physical address spaces and other aspects of hardware organization that evolved over many years is one of the important practical challenges of OS development.

## The ROM BIOS

In this portion of the lab, you'll use QEMU's debugging facilities to investigate how an IA-32 compatible computer boots.

Open two terminal windows. In one, enter `make qemu-gdb` (or `make qemu-nox-gdb`). This starts up QEMU, but QEMU stops just before the processor executes the first instruction and waits for a debugging connection from GDB. In the second terminal, from the same directory you ran `make`, run `gdb`. You should see something like this,

```

kermit% gdb
GNU gdb (GDB) 6.8-debian
Copyright (C) 2008 Free Software Foundation, Inc.
License GPLv3+: GNU GPL version 3 or later <http://gnu.org/licenses/gpl.html>
This is free software: you are free to change and redistribute it.
There is NO WARRANTY, to the extent permitted by law. Type "show copying"
and "show warranty" for details.
This GDB was configured as "i486-linux-gnu".
+ target remote localhost:1234
The target architecture is assumed to be i8086
[f000:fff0] 0xfffff0:  jmp  $0xf000,$0xe05b
0x0000ffff in ?? ()
(gdb)

```

We provided a `.gdbinit` file that set up GDB to debug the 16-bit code used during early boot and directed it to attach to the listening QEMU.

The following line:

```
[f000:fff0] 0xfffff0:  jmp  $0xf000,$0xe05b
```

is GDB's disassembly of the first instruction to be executed. From this output you can conclude a few things:

- The IBM PC starts executing at physical address `0x000ffff0`, which is at the very top of the 64KB area reserved for the ROM BIOS.
- The PC starts executing with `CS = 0xf000` and `IP = 0xffff0`.
- The first instruction to be executed is a `jmp` instruction, which jumps to the segmented address `CS = 0xf000` and `IP = 0xe05b`.

Why does QEMU start like this? This is how Intel designed the 8088 processor, which IBM used in their original PC. Because the BIOS in a PC is "hard-wired" to the physical address range `0x000f0000-0x000fffff`, this design ensures that the BIOS always gets control of the machine first after power-up or any system restart - which is crucial because on power-up there is no other software anywhere in the machine's RAM that the processor could execute. The QEMU emulator comes with its own BIOS, which it places at this location in the processor's simulated physical address space. On processor reset, the (simulated) processor enters real mode and sets `CS` to `0xf000` and the `IP` to `0xffff0`, so that execution begins at that (`CS:IP`) segment address. How does the segmented address `0xf000:fff0` turn into a physical address?

To answer that we need to know a bit about real mode addressing. In real mode (the mode that PC starts off in), address translation works according to the formula: *physical address* = *16 \* segment + offset*. So, when the PC sets `CS` to `0xf000` and `IP` to `0xffff0`, the physical address referenced is:

```

16 * 0xf000 + 0xffff0  # in hex multiplication by 16 is
= 0xf0000 + 0xffff0    # easy--just append a 0.
= 0xfffff0

```

`0xfffff0` is 16 bytes before the end of the BIOS (`0x100000`). Therefore we shouldn't be surprised that the first thing that the BIOS does is `jmp` backwards to an earlier location in the BIOS; after all how much could it accomplish in just 16 bytes?

**Exercise 2.** Use GDB's `si` (Step Instruction) command to trace into the ROM BIOS for a few more instructions, and try to guess what it might be doing. You might want to look at [Phil Storrs I/O Ports Description](#), as well as other materials on the [CSE 506 reference materials page](#). No need to figure out all the details - just the general idea of what the BIOS is doing first.

When the BIOS runs, it sets up an interrupt descriptor table and initializes various devices such as the VGA display. This is where the "Starting SeaBIOS" messages you see in the QEMU window come from.

After initializing the PCI bus and all the important devices the BIOS knows about, it searches for a bootable device such as a floppy, hard drive, or CD-ROM. Eventually, when it finds a bootable disk, the BIOS reads the *boot loader* from the disk and transfers control to it.

## Part 2: The Boot Loader

Floppy and hard disks for PCs are divided into 512 byte regions called *sectors*. A sector is the disk's minimum transfer granularity: each read or write operation must be one or more sectors in size and aligned on a sector boundary. If the disk is bootable, the first sector is called the *boot sector*, since this is where the boot loader code resides. When the BIOS finds a bootable floppy or hard disk, it loads the 512-byte boot sector into memory at physical addresses `0x7c00` through `0x7dff`, and then uses a `jmp` instruction to set the `CS:IP` to `0000:7c00`, passing control to the boot loader. Like the BIOS load address, these addresses are fairly arbitrary - but they are fixed and standardized for PCs.



The ability to boot from a CD-ROM came much later during the evolution of the PC, and as a result the PC architects took the opportunity to rethink the boot process slightly. As a result, the way a modern BIOS boots from a CD-ROM is a bit more complicated (and more powerful). CD-ROMs use a sector size of 2048 bytes instead of 512, and the BIOS can load a much larger boot image from the disk into memory (not just one sector) before transferring control to it. For more information, see the ["El Torito" Bootable CD-ROM Format Specification](#).

For CSE 506, however, we will use the conventional hard drive boot mechanism, which means that the BIOS will only load the first 512 bytes. The bootloader (boot/boot.S and boot/main.c) does the bootstrapping work. Look through these source files carefully and make sure you understand what's going on. The boot loader must perform the following main functions:

1. First, the boot loader obtains a map of the physical memory present in the system from the BIOS. This is done using a system call to the BIOS (int 0x15), which returns a structure called an e820 map. This call is possible only while the processor is still in real mode. JOS's bootloader constructs a multiboot structure, which it passes to the kernel. [Multiboot](#) is a standard for passing boot information from the bootloader to a kernel.
2. The boot loader then switches the processor from real mode to *32-bit protected mode*, because it is only in this mode that software can access all the memory above 1MB in the processor's physical address space. Protected mode is described briefly in sections 1.2.7 and 1.2.8 of [PC Assembly Language](#), and in great detail in the Intel architecture manuals. At this point you only have to understand that translation of segmented addresses (segment:offset pairs) into physical addresses happens differently in protected mode, and that after the transition offsets are 32 bits instead of 16. One of the arcane x86 features the bootloader handles is properly configuring address bit 20. To better understand this, skim [this article](#).
3. The bootloader sets up the stack and starts executing the kernel's C code in boot/main.c.
4. Finally, the boot loader reads the kernel from the hard disk by directly accessing the IDE disk device registers via the x86's special I/O instructions. If you would like to understand better what the particular I/O instructions here mean, check out the "IDE hard drive controller" section on [the CSE 506 reference page](#). You will not need to learn much about programming specific devices in this class: writing device drivers is in practice a very important part of OS development, but from a conceptual or architectural viewpoint it is also one of the least interesting.

After you understand the boot loader source code, look at the files obj/boot/boot.asm. This file is a disassembly of the boot loader that our GNUmakefile creates *after* compiling the boot loader. This disassembly file makes it easy to see exactly where in physical memory all of the boot loader's code resides, and makes it easier to track what's happening while stepping through the boot loader in GDB.

You can set address breakpoints in GDB with the **b** command. You have to start hex numbers with 0x, so say something like **b \*0x7c00** sets a breakpoint at address 0x7C00. Once at a breakpoint, you can continue execution using the **c** and **si** commands: **c** causes QEMU to continue execution until the next breakpoint (or until you press **Ctrl-C** in GDB), and **si N** steps through the instructions *N* at a time.

To examine instructions in memory (besides the immediate next one to be executed, which GDB prints automatically), you use the **x/i** command. This command has the syntax **x/Ni ADDR**, where *N* is the number of consecutive instructions to disassemble and *ADDR* is the memory address at which to start disassembling.

**Exercise 3.** Take a look at the [lab tools guide](#), especially the section on GDB commands. Even if you're familiar with GDB, this includes some esoteric GDB commands that are useful for OS work.

Set a breakpoint at address 0x7c00, which is where the boot sector will be loaded. Continue execution until that breakpoint. Trace through the code in boot/boot.S, using the source code and the disassembly file obj/boot/boot.asm to keep track of where you are. Also use the x/i command in GDB to disassemble sequences of instructions in the boot loader, and compare the original boot loader source code with both the disassembly in obj/boot/boot.asm and GDB.

Trace into bootmain() in boot/main.c, and then into readsect(). Identify the exact assembly instructions that correspond to each of the statements in readsect(). Trace through the rest of readsect() and back out into bootmain(), and identify the begin and end of the for loop that reads the remaining sectors of the kernel from the disk. Find out what code will run when the loop is finished, set a breakpoint there, and continue to that breakpoint. Then step through the remainder of the boot loader.

Be able to answer the following questions:

- At what point does the processor start executing 32-bit code? What exactly causes the switch from 16- to 32-bit mode?
- What is the *last* instruction of the boot loader executed, and what is the *first* instruction of the kernel it just loaded?
- How does the boot loader decide how many sectors it must read in order to fetch the entire kernel from disk? Where does it find this information?

## Loading the Kernel

We will now look in further detail at the C language portion of the boot loader, in boot/main.c. But before doing so, this is a good time to stop and review some of the basics of C programming.

**Exercise 4.** Read about programming with pointers in C. The best reference for the C language is *The C Programming Language* by Brian Kernighan and Dennis Ritchie (known as 'K&R'). We recommend that students purchase this book (here is an [Amazon Link](#)). There are several copies on reserve in the Science and Engineering library as well.

Read 5.1 (Pointers and Addresses) through 5.5 (Character Pointers and Functions) in K&R. Then download the code for [pointers.c](#), run it, and make sure you understand where all of the printed values come from. In particular, make sure you understand where the pointer addresses in lines 1 and 6 come from, how all the values in lines 2 through 4 get there, and why the values printed in line 5 are seemingly corrupted.

There are other references on pointers in C, though not as strongly recommended. [A tutorial by Ted Jensen](#) that cites

K&R heavily is available in the course readings.

We also recommend reading the [Ksplice pointer challenge](#) as a way to test that you understand how pointer arithmetic and arrays work in C.

*Warning:* Unless you are already thoroughly versed in C, do not skip or even skim this reading exercise. If you do not really understand pointers in C, you will suffer untold pain and misery in subsequent labs, and then eventually come to understand them the hard way. Trust us; you don't want to find out what "the hard way" is.

To make sense out of `boot/main.c` you'll need to know what an ELF binary is. When you compile and link a C program such as the JOS kernel, the compiler transforms each C source (`.c`) file into an *object* (`.o`) file containing assembly language instructions encoded in the binary format expected by the hardware. The linker then combines all of the compiled object files into a single *binary image* such as `obj/kern/kernel`, which in this case is a binary in the ELF format, which stands for "Executable and Linkable Format".

Full information about this format is available in [the ELF specification](#) on [our reference page](#), but you will not need to delve very deeply into the details of this format in this class. Although as a whole the format is quite powerful and complex, most of the complex parts are for supporting dynamic loading of shared libraries, which we will not do in this class.

For purposes of CSE 506, you can consider an ELF executable to be a header with loading information, followed by several *program sections*, each of which is a contiguous chunk of code or data intended to be loaded into memory at a specified address. The boot loader does not modify the code or data; it loads it into memory and starts executing it.

An ELF binary starts with a fixed-length *ELF header*, followed by a variable-length *program header* listing each of the program sections to be loaded. The C definitions for these ELF headers are in `inc/elf.h`. The program sections we're interested in are:

- `.text`: The program's executable instructions.
- `.rodata`: Read-only data, such as ASCII string constants produced by the C compiler. (We will not bother setting up the hardware to prohibit writing, however.)
- `.data`: The data section holds the program's initialized data, such as global variables declared with initializers like `int x = 5;`.

When the linker computes the memory layout of a program, it reserves space for *uninitialized* global variables, such as `int x;`, in a section called `.bss` that immediately follows `.data` in memory. C requires that "uninitialized" global variables start with a value of zero. Thus there is no need to store contents for `.bss` in the ELF binary; instead, the linker records just the address and size of the `.bss` section. The loader or the program itself must arrange to zero the `.bss` section.

You can display a full list of the names, sizes, and link addresses of all the sections in the kernel executable by typing:

```
kermit% objdump -h obj/kern/kernel
```

You will see many more sections than the ones we listed above, but the others are not important for our purposes. Most of the others are to hold debugging information, which is typically included in the program's executable file but not loaded into memory by the program loader.

Take particular note of the "VMA" (or *link address*) and the "LMA" (or *load address*) of the `.text` section. The load address of a section is the memory address at which that section should be loaded into memory. In the ELF object, this is stored in the `ph->p_pa` field (in this case, it really is a physical address, though the ELF specification is vague on the actual meaning of this field).

The link address of a section is the memory address from which the section expects to execute. The linker encodes the link address in the binary in various ways, such as when the code needs the address of a global variable, with the result that a binary usually won't work if it is executing from an address that it is not linked for. (It is possible to generate *position-independent* code that does not contain any such absolute addresses. This is used extensively by modern shared libraries, but it has performance and complexity costs, so we won't be using it in CSE 506.)

Typically, the link and load addresses are the same. For example, look at the `.text` section of the boot loader:

```
kermit objdump -h obj/boot/boot.out
```

The BIOS loads the boot sector into memory starting at address `0x7c00`, so this is the boot sector's load address. This is also where the boot sector executes from, so this is also its link address. We set the link address by passing `-Ttext 0x7C00` to the linker in `boot/Makefrag`, so the linker will produce the correct memory addresses in the generated code.

**Exercise 5.** Trace through the first few instructions of the boot loader again and identify the first instruction that would "break" or otherwise do the wrong thing if you were to get the boot loader's link address wrong. Then change the link address in `boot/Makefrag` to something wrong, run `make clean`, recompile the lab with `make`, and trace into the boot loader again to see what happens. Don't forget to change the link address back and `make clean` again afterward!

Look back at the load and link addresses for the kernel. Unlike the boot loader, these two addresses aren't the same: the kernel is telling the boot loader to load it into memory at a low address (1 megabyte), but it expects to execute from a high address. We'll dig in to how we make this work in the next section.

Besides the section information, there is one more field in the ELF header that is important to us, named `e_entry`. This field holds the link address of the *entry point* in the program: the memory address in the program's text section at which the program should begin executing. You can see the entry



point:

```
kermit% objdump -f obj/kern/kernel
```

You should now be able to understand the minimal ELF loader in `boot/main.c`. It reads each section of the kernel from disk into memory at the section's load address and then jumps to the kernel's entry point.

**Exercise 6.** We can examine memory using GDB's `x` command. The [GDB manual](#) has full details, but for now, it is enough to know that the command `x/Nx ADDR` prints *N* words of memory at *ADDR*. (Note that both 'x's in the command are lowercase.) *Warning:* The size of a word is not a universal standard. In GNU assembly, a word is two bytes (the 'w' in `xorw`, which stands for word, means 2 bytes).

Reset the machine (exit QEMU/GDB and start them again). Examine the 8 words of memory at `0x00100000` at the point the BIOS enters the boot loader, and then again at the point the boot loader enters the kernel. Why are they different? What is there at the second breakpoint? (You do not really need to use QEMU to answer this question. Just think.)

## Link vs. Load Address

The *load address* of a binary is the memory address at which a binary is *actually* loaded. For example, the BIOS is loaded by the PC hardware at address `0xf000`. So this is the BIOS's load address. Similarly, the BIOS loads the boot sector at address `0x7c00`. So this is the boot sector's load address.

The *link address* of a binary is the memory address for which the binary is linked. Linking a binary for a given link address prepares it to be loaded at that address. The linker encodes the link address in the binary in various ways, for example when the code needs the address of a global variable, with the result that a binary usually won't work if it is not loaded at the address that it is linked for.

In one sentence: the link address is the location where a binary *assumes* it is going to be loaded, while the load address is the location where a binary is loaded. It's up to us to make sure that they turn out to be the same.

Look at the `-Ttext` linker command in `boot/Makefrag`, and at the address mentioned early in the linker script in `kern/kernel.ld`. These set the link address for the boot loader and kernel respectively.

When object code contains no absolute addresses that encode the link address in this fashion, we say that the code is *position-independent*: it will behave correctly no matter where it is loaded. GCC can generate position-independent code using the `-fpic` option, and this feature is used extensively in modern shared libraries that use the ELF executable format. Position independence typically has some performance cost, however, because it restricts the ways in which the compiler may choose instructions to access the program's data. We will not use `-fpic` in CSE 506.

## Part 3: The Kernel

We will now start to examine the minimal JOS kernel in a bit more detail. (And you will finally get to write some code!). Like the boot loader, the kernel begins with some assembly language code that sets things up so that C language code can execute properly.

The initial assembly code of the kernel does the following:

1. When the kernel starts executing, the processor is in the 32-bit protected mode. The first thing the kernel does is it tests whether the CPU supports long (64-bit) mode.
2. The kernel initializes a simple set of page tables for the first 4GB of memory. These pages map virtual addresses in the lowest 3GB to the same physical addresses, and then map the upper 256 MB back to the lowest 256 MB of memory. At this point, the kernel places the CPU in long mode. The kernel could determine dynamically whether to run in 64 or 32-bit mode based on whether the CPU supports long mode. Of course, this would substantially complicate the kernel.
3. Finally, the kernel sets up the stack and a few other things to start executing C code.

Look through the kernel source files `kern/entry.S` and `kern/bootstrp.S` and be able to answer the following question

- At what point does the processor start executing 64-bit code? What exactly causes the switch from 32- to 64-bit mode?

## Using segmentation to work around position dependence

Did you notice above that while the boot loader's link and load addresses match perfectly, there appears to be a (rather large) disparity between the *kernel's* link and load addresses? Go back and check both and make sure you can see what we're talking about.

Operating system kernels often like to be linked and run at very high *virtual address*, such as `0x800410000`, in order to leave the lower part of the processor's virtual address space for user programs to use. The reason for this arrangement will become clearer in the next lab.

Many machines don't have any physical memory at address `0x800410000` so we can't count on being able to store the kernel there. Instead, we will use the processor's memory management hardware to map virtual address `0x800410000` - the link address at which the kernel code *expects* to run - to physical address `0x100000` - where the boot loader loaded the kernel. This way, although the kernel's virtual address is high enough to leave plenty of address space for user processes, it will be loaded in physical memory at the 1MB point in the PC's RAM, just above the BIOS ROM. This approach requires that the PC have at least a few megabytes of physical memory (so that address `0x00100000` works), but this is likely to be true of any PC built after about 1990.

We will eventually map physical address 0x0 to virtual address .

We will eventually map the *entire* bottom 256MB of the PC's physical address space, from 0x00000000 through 0x0ffffff, to virtual addresses 0x800400000 through 0x8013ffffff, respectively.

For now, the kernel will initially set up a hand-written, statically-initialized page tables in kern/bootstrap.S. For now, you don't have to understand the details of how this works, just the effect that it accomplishes. Up until kern/bootstrap.S sets the CR0\_PG flag, memory references are treated as physical addresses (strictly speaking, they're linear addresses, but boot/boot.S set up an identity mapping from linear addresses to physical addresses and we're never going to change that). Once CR0\_PG is set, memory references are virtual addresses that get translated by the virtual memory hardware to physical addresses. pm14 translates virtual addresses in the range 0x800400000 through 0x8013ffffff to physical addresses 0x00000000 through 0x0ffffff, as well as virtual addresses 0x00000000 through 0xffffffff to physical addresses 0x00000000 through 0xffffffff.

**Exercise 7.** Use QEMU and GDB to trace into the early JOS kernel boot code (in the kern/bootstrap.S directory) and find where the new virtual-to-physical mapping takes effect. Then examine the Global Descriptor Table (GDT) that the code uses to achieve this effect, and make sure you understand what's going on.

What is the first instruction *after* the new mapping is established that would fail to work properly if the old mapping were still in place? Comment out or otherwise intentionally break the segmentation setup code in kern/entry.S, trace into it, and see if you were right.

## Formatted Printing to the Console

Most people take functions like printf() for granted, sometimes even thinking of them as "primitives" of the C language. But in an OS kernel, we have to implement all I/O ourselves.

Read through kern/printf.c, lib/printfmt.c, and kern/console.c, and make sure you understand their relationship. It will become clear in later labs why printfmt.c is located in the separate lib directory.

**Exercise 8.** We have omitted a small fragment of code - the code necessary to print octal numbers using patterns of the form "%o". Find and fill in this code fragment.

Be able to answer the following questions:

1. Explain the interface between printf.c and console.c. Specifically, what function does console.c export? How is this function used by printf.c?
2. Explain the following from console.c:

```
1   if (crt_pos >= CRT_SIZE) {
2       int i;
3       memmove(crt_buf, crt_buf + CRT_COLS, (CRT_SIZE - CRT_COLS) * sizeof(uint16_t));
4       for (i = CRT_SIZE - CRT_COLS; i < CRT_SIZE; i++)
5           crt_buf[i] = 0x0700 | ' ';
6       crt_pos -= CRT_COLS;
7   }
```

3. Trace the execution of the following code step-by-step:

```
int x = 1, y = 3, z = 4;
cprintf("x %d, y %x, z %d\n", x, y, z);
```

- In the call to cprintf(), to what does fmt point? To what does ap point?
- List (in order of execution) each call to cons\_putc, va\_arg, and vprintf. For cons\_putc, list its argument as well. For va\_arg, list what ap points to before and after the call. For vprintf list the values of its two arguments.

4. Run the following code.

```
unsigned int i = 0x00646c72;
cprintf("H%x Wo%s", 57616, &i);
```

What is the output? Explain how this output is arrived at in the step-by-step manner of the previous exercise. [Here's an ASCII table](#) that maps bytes to characters.

The output depends on that fact that the x86 is little-endian. If the x86 were instead big-endian what would you set i to in order to yield the same output? Would you need to change 57616 to a different value?

[Here's a description of little- and big-endian](#) and [a more whimsical description](#).

5. In the following code, what is going to be printed after 'y='? (note: the answer is not a specific value.) Why does this happen?

```
cprintf("x=%d y=%d", 3);
```

6. Let's say that GCC changed its calling convention so that it pushed arguments on the stack in declaration order, so that the last argument is pushed last. How would you have to change `cprintf` or its interface so that it would still be possible to pass it a variable number of arguments?

**Challenge 1** (5 bonus points) Enhance the console to allow text to be printed in different colors. The traditional way to do this is to make it interpret [ANSI escape sequences](#) embedded in the text strings printed to the console, but you may use any mechanism you like. There is plenty of information on [the CSE 506 reference page](#) and elsewhere on the web on programming the VGA display hardware. If you're feeling really adventurous, you could try switching the VGA hardware into a graphics mode and making the console draw text onto the graphical frame buffer.

## The Stack

In the final exercise of this lab, we will explore in more detail the way the C language uses the stack on the x86, and in the process write a useful new kernel monitor function that prints a *backtrace* of the stack: a list of the saved Instruction Pointer (IP) values from the nested `call` instructions that led to the current point of execution.

**Exercise 9.** Determine where the kernel initializes its stack, and exactly where in memory its stack is located. How does the kernel reserve space for its stack? And at which "end" of this reserved area is the stack pointer initialized to point to?

The x86-64 stack pointer (`rsp` register) points to the lowest location on the stack that is currently in use. Everything *below* that location in the region reserved for the stack is free. Pushing a value onto the stack involves decreasing the stack pointer and then writing the value to the place the stack pointer points to. Popping a value from the stack involves reading the value the stack pointer points to and then increasing the stack pointer. In 64-bit mode, the stack can only hold 64-bit values, and `rsp` is always divisible by eight. Various x86-64 instructions, such as `call`, are "hard-wired" to use the stack pointer register.

The `rbp` (base pointer) register, in contrast, is associated with the stack primarily by software convention. On entry to a C function, the function's *prologue* code normally saves the previous function's base pointer by pushing it onto the stack, and then copies the current `rsp` value into `rbp` for the duration of the function. If all the functions in a program obey this convention, then at any given point during the program's execution, it is possible to trace back through the stack by following the chain of saved `rbp` pointers and determining exactly what nested sequence of function calls caused this particular point in the program to be reached. This capability can be particularly useful, for example, when a particular function causes an `assert` failure or `panic` because bad arguments were passed to it, but you aren't sure *who* passed the bad arguments. A stack backtrace lets you find the offending function.

**Exercise 10.** To become familiar with the C calling conventions on the x86-64, find the address of the `test_backtrace` function in `obj/kern/kernel.asm`, set a breakpoint there, and examine what happens each time it gets called after the kernel starts. How many 64-bit words does each recursive nesting level of `test_backtrace` push on the stack, and what are those words?

Note that, for this exercise to work properly, you should be using the patched version of QEMU available on the [tools](#) page or on your course virtual machine. Otherwise, you'll have to manually translate all breakpoint and memory addresses to linear addresses.

The above exercise should give you the information you need to implement a stack backtrace function, which you should call `mon_backtrace()`. A prototype for this function is already waiting for you in `kern/monitor.c`. You can do it entirely in C, but you may find the `read_rbp()` function in `inc/x86.h` useful. You'll also have to hook this new function into the kernel monitor's command list so that it can be invoked interactively by the user.

The backtrace function should display a listing of function call frames in the following format:

```
Stack backtrace:
  rbp 00000000f0111f20  rip 00000000f01000be
  rbp 00000000f0111f40  rip 00000000f01000a1
  ...
```

The first line printed reflects the *currently executing* function, namely `mon_backtrace` itself, the second line reflects the function that called `mon_backtrace`, the third line reflects the function that called that one, and so on. You should print *all* the outstanding stack frames. By studying `kern/entry.S` you'll find that there is an easy way to tell when to stop.

Within each line, the `rbp` value indicates the base pointer into the stack used by that function: i.e., the position of the stack pointer just after the function was entered and the function prologue code set up the base pointer. The listed `rip` value is the function's *return instruction pointer*: the instruction address to which control will return when the function returns. The return instruction pointer typically points to the instruction after the `call` instruction (why?).

Read [this article](#) on how arguments are mapped to registers on the x86-64 architecture. This is called a *calling convention*, as software developers agree on this mapping by convention. As the article explains, different compilers may have different standards. For instance, the Windows x86-64 calling convention differs from the Linux x86-64 calling convention.

Here are a few specific points you read about in K&R Chapter 5 that are worth remembering for the following exercise and for future labs.

- If `int *p = (int*)100`, then `(int)p + 1` and `(int)(p + 1)` are different numbers: the first is 101 but the second is 104. When adding an

integer to a pointer, as in the second case, the integer is implicitly multiplied by the size of the object the pointer points to.

- `p[i]` is defined to be the same as `*(p+i)`, referring to the *i*'th object in the memory pointed to by `p`. The above rule for addition helps this definition work when the objects are larger than one byte.
- `&p[i]` is the same as `(p+i)`, yielding the address of the *i*'th object in the memory pointed to by `p`.

Although most C programs never need to cast between pointers and integers, operating systems frequently do. Whenever you see an addition involving a memory address, ask yourself whether it is an integer addition or pointer addition and make sure the value being added is appropriately multiplied or not.

**Exercise 11.** Implement the backtrace function as specified above. Use the same format as in the example, since otherwise the grading script will be confused. When you think you have it working right, run `make grade` to see if its output conforms to what our grading script expects, and fix it if it doesn't. After you have handed in your Lab 1 code, you are welcome to change the output format of the backtrace function any way you like.

At this point, your backtrace function should give you the addresses of the function callers on the stack that lead to `mon_backtrace()` being executed. However, in practice you often want to know the function names corresponding to those addresses. For instance, you may want to know which functions could contain a bug that's causing your kernel to crash.

The final exercise will have you list the arguments to each function in the backtrace, if the function has any input values.

On x86-64, arguments are generally passed in registers. If function `a` calls function `b`, `a` may save some register values on the stack and reload them, depending on whether a value will be reused. These decisions are made by the compiler. Thus, in order to find arguments on the stack, we need some help from the compiler.

When JOS is compiled with debugging flags, the compiler outputs a variety of debugging information in the DWARF2 format. Read [this article](#) for a brief introduction to DWARF and how debugging symbols work. The key intuition is that DWARF gives the debugger (or monitor backtrace function) enough information to programmatically identify saved arguments on the stack.

To help you implement this functionality, we have provided the function `debuginfo_rip()`, which looks up `rip` in the symbol table and returns the debugging information for that address. This function is defined in `kern/kdebug.c`.

**Exercise 12.** Modify your stack backtrace function to display, for each `rip`, the function arguments, the function name, source file name, and line number corresponding to that `rip`.

Add a backtrace command to the kernel monitor, and extend your implementation of `mon_backtrace` to call `debuginfo_rip` and print a line for each stack frame of the form:

```
K> backtrace
Stack backtrace:
rbp 000000800421af00 rip 00000080042010ff
    kern/monitor.c:86: mon_backtrace+0000000000000035 args:3 0000000000000000 000000000421b909 0000000000000000
rbp 000000800421afb0 rip 000000800420144d
    kern/monitor.c:163: runcmd+00000000000001d3 args:2 0000000000000001 0000000000000002
rbp 000000800421afe0 rip 0000008004201508
    kern/monitor.c:185: monitor+000000000000007d args:1 0000000000000080
rbp 000000800421aff0 rip 0000008004200196
    kern/init.c:172: i386_init+00000000000000ba args:0
```

Each line gives the file name and line within that file of the stack frame's `rip`, followed by the name of the function and the offset of the `rip` from the first instruction of the function (e.g., `monitor+106` means the return `rip` is 106 bytes past the beginning of `monitor`), followed by the number of function arguments and then the actual arguments themselves.

Hint: for the function arguments, take a look at the struct `Ripdebuginfo` in `kern/kdebug.h`. This structure is filled by the call to `debuginfo_rip`. The x86\_64 calling convention states that the function arguments are pushed onto the stack. Refer to [this article](#) on the calling convention to figure out how to read the actual function arguments on the stack.

Be sure to print the file and function names on a separate line, to avoid confusing the grading script.

Tip: `printf` format strings provide an easy, albeit obscure, way to print non-null-terminated strings like those in the DWARF2 tables. `printf("%.5s", length, string)` prints at most `length` characters of `string`. Take a look at the `printf` man page to find out why this works.

You may find that some functions are missing from the backtrace. For example, you will probably see a call to `monitor()` but not to `runcmd()`. This is because the compiler in-lines some function calls. Other optimizations may cause you to see unexpected line numbers. If you get rid of the `-O2` from `GNUmakefile`, the backtraces may make more sense (but your kernel will run more slowly).

**This completes the lab.** Type `make handin` in the lab directory. This will only work after submitting your partner selection to the instructor. After successful submission, you should receive a confirmation email (although you may need to try twice to get the email). If submission fails, double check that you have committed all of your changes, and read any error messages carefully before emailing the TAs for help.

