



南開大學
Nankai University

OPERATING SYSTEM EXPERIMENTAL REPORT I

BOOTING A PC AND MEMORY MANAGEMENT

Yinhao Li 1611303

Jiajun Li 1613368

Qiaoqiao Pan 1613415

Xinghe Li 1611071

Hongbo Xu 1611336

Abstract

Through the study of Experiment 1, we learned how the operating system works when the computer starts and how the computer manages the memory.

In Chapter 1, we will discuss Computer Boot, Boot Loader and The Kernel in detail. We will start research on the X86 assembly instructions, BIOS execution, PC physical address space, how the Boot Loader loads the kernel, the connection address and load address, the use of virtual memory, the formatted output of the console, and the stack.

In Chapter 2, we will discuss physical page management, virtual memory, and kernel address space. Describe in detail the relationship between virtual addresses and linear and physical addresses, reference counts, page table management, permissions and fault isolation, initialization of linear portions of the kernel, and so on.

Key Words: Computer Boot Boot Loader The Kernel physical page management virtual memory kernel address space

CONTENTS

Abstract	I
Chapter1 Booting a PC	1
1.1 Computer Boot	1
1.1.1 Simulating The X86	1
1.1.2 Exercise I	1
1.1.3 The PC's Physical Address Space	4
1.1.4 BIOS	6
1.2 Boot Loader	7
1.2.1 Exercise II	7
1.2.2 Question I	8
1.2.3 Loading kernel	13
1.2.4 Connection address and load address	14
1.3 The Kernel	16
1.3.1 Question	16
1.3.2 Use virtual memory	17
1.3.3 Homework I	18
1.3.4 Exercise III	21
1.3.5 Homework II	23
1.3.6 Challenge homework	25
Chapter2 Memory Management	29
2.1 Physical Page Management	29
2.1.1 Homework III	29
2.2 Virtual Memory	36
2.2.1 Exercise IV	36
2.2.2 Virtual address, linear address and physical address	36
2.2.3 Exercise V	37
2.2.4 Question III	38
2.2.5 Reference counting	39
2.2.6 Page Table Management	39
2.2.7 Homework IV	39
2.3 Kernel address space	47
2.3.1 Homework V	48
2.3.2 Question IV	52

Chapter1 Booting a PC

1.1 Computer Boot

1.1.1 Simulating The X86

In this experiment, we chose QEMU to simulate a real computer. The use of QEMU is as follows.

```

——内存空间——
物理内存: 66556K
basemem = 640K,
extended = 65532K,
IO = 384K

——各类内存页数——
npages: 16639
npages_extmem: 16383
npages_basemem: 160
npages_IO: 96

OS页表起始时的虚拟地址: f0118000
OS页表结束、pages开始的虚拟地址: f0119000
OS的页表的第0项(指针kern_pgdir)所在的物理内存地址: 118005
pages页表结束的虚拟地址: f013a000

——数组pages初始化——
【BASEMEM】 0-159
0被占用, 1-159未被占用
【IO hole】 160-255
160-255全部被占用
【EXTMEM】 256-16638
256-313被占用, 314-16638未被占用

check_page_alloc() succeeded!
check_page() succeeded!
UPAGES
map_region_size =133112,(32 pages)
memory stack
map_region_size =32768,(8 pages)
kernel
map_region_size =268435456,(65536 pages)
check_kern_pgdir() succeeded!
check_page_installed_pgdir() succeeded!
Welcome to the JOS kernel monitor!

```

图 1.1: make qemu

```

K> kerninfo
Special kernel symbols:
_start      0010000c (phys)
entry       f010000c (virt) 0010000c (phys)
etext       f0101861 (virt) 00101861 (phys)
edata       f0112300 (virt) 00112300 (phys)
end         f0112944 (virt) 00112944 (phys)
Kernel executable memory footprint: 75KB
K>

```

图 1.2: kerninfo

Through the above two instructions, we can see that the kernel information has been output in the console. It can be seen that QEMU can function as a simulation computer and the system can run normally.

1.1.2 Exercise I

Question

Refer to the NASM assembler used in the book, and we are using the GNU assembler. The syntax is slightly different. Also, refer to Brennan's Guide to Inline Assembly for the differences.

Answer

Basic structure

AT&T and Intel syntax use the opposite order for source and destination operands. For example:

Intel: `mov eax, 4`

AT&T: `movl $4, %eax`

In AT&T syntax, immediate operands are preceded by `$`; in Intel syntax, immediate operands are not. For example:

Intel: `push 4`

AT&T: `pushl $4`

In AT&T syntax, register operands are preceded by `%`; in Intel syntax, they are not.

In AT&T syntax, the size of memory operands is determined by the last character of the opcode name. Opcode suffixes of `b`, `w`, and `l` specify byte (8-bit), word (16-bit), and long (32-bit) memory references. Intel syntax accomplishes this by prefixing memory operands (not the opcodes themselves) with `byte ptr`, `word ptr`, and `dword ptr`. Thus:

Intel: `mov al, byte ptr foo`

AT&T: `movb foo, %al`

Immediate form long jumps and calls are `lcall/ljmp $section, $offset` in AT&T syntax; the Intel syntax is `call/jmp far section:offset`. The far return instruction is `lret $stack-adjust` in AT&T syntax, whereas Intel uses `ret far stack-adjust`.

Variables and accessing memory

We can find the difference between the two grammars by the following code.

Line	NASM	GAS
1	001	1
2	002	2
3	003	3
4	004	4
5	005	5
6	006	6
7	007	7
8	008	8
9	009	9
10	010	10
11	011	11
12	012	12
13	013	13
14	014	14
15	015	15
16	016	16
17	017	17
18	018	18
19	019	19
20	020	20
21	021	21
22	022	22
23	023	23
24	024	24
25	025	25
26	026	26
27	027	27
28	028	28
29	029	29
30	030	30
31	031	31
	<pre> ; Data section begins section .data var1 dd 40 var2 dd 20 var3 dd 30 section .text global _start _start: ; Move the contents of variables mov ecx, [var1] cmp ecx, [var2] jg check_third_var mov ecx, [var2] check_third_var: cmp ecx, [var3] jg _exit mov ecx, [var3] _exit: mov eax, 1 mov ebx, ecx int 80h </pre>	<pre> // Data section begins .section .data var1: .int 40 var2: .int 20 var3: .int 30 .section .text .globl _start _start: # move the contents of variables movl (var1), %ecx cmpl (var2), %ecx jg check_third_var movl (var2), %ecx check_third_var: cmpl (var3), %ecx jg _exit movl (var3), %ecx _exit: movl \$1, %eax movl %ecx, %ebx int \$0x80 </pre>

We can see several differences above in the declaration of memory variables. NASM uses the `dd`, `dw`, and `db` directives to declare 32-, 16-, and 8-bit numbers, respectively,

whereas GAS uses the `.long`, `.int`, and `.byte` for the same purpose. GAS has other directives too, such as `.ascii`, `.asciz`, and `.string`. In GAS, you declare variables just like other labels (using a colon), but in NASM you simply type a variable name (without the colon) before the memory allocation directive (`dd`, `dw`, etc.), followed by the value of the variable.

Line 18 illustrates the memory indirect addressing mode. NASM uses square brackets to dereference the value at the address pointed to by a memory location: `[var1]`. GAS uses a circular brace to dereference the same value: `(var1)`. The use of other addressing modes is covered later in this article.

Using macros

NASM users might be familiar with declaring macros using the `%beginmacro` directive and ending them with an `%endmacro` directive. A `%beginmacro` directive is followed by the macro name. After the macro name comes a count, the number of macro arguments the macro is supposed to have. In NASM, macro arguments are numbered sequentially starting with 1. That is, the first argument to a macro is `%1`, the second is `%2`, the third is `%3`, and so on. For example:

```
%beginmacro macroname 2    mov eax, %1    mov ebx, %2    %endmacro
```

This creates a macro with two arguments, the first being `%1` and the second being `%2`. Thus, a call to the above macro would look something like this:

```
macroname 5, 6
```

Macros can also be created without arguments, in which case they don't specify any number.

Now let's take a look at how GAS uses macros. GAS provides the `.macro` and `.endm` directives to create macros. A `.macro` directive is followed by a macro name, which may or may not have arguments. In GAS, macro arguments are given by name. For example:

```
.macro macroname arg1, arg2
    movl arg1, %eax
    movl arg2, %ebx
.endm
```

A backslash precedes the name of each argument of the macro when the name is actually used inside a macro. If this is not done, the linker would treat the names as labels rather than as arguments and will report an error.

Functions, external routines, and the stack

In NASM, any label that is not part of our program that needs to be resolved during link time must be predefined, which is the function of the `extern` keyword (line 24).

GAS doesn't have such requirements. After this, the address of the string `usort_str` is

pushed onto the stack (line 30). In NASM, a memory variable such as `usort_str` represents the address of the memory location itself, and thus a call such as `push usort_str` actually pushes the address on top of the stack. In GAS, on the other hand, the variable `usort_str` must be prefixed with `$`, so that it is treated as an immediate address. If it's not prefixed with `$`, the actual bytes represented by the memory variable are pushed onto the stack instead of the address.

The addressing modes are represented in NASM and GAS as follows:

NASM: `mov al, byte [ebx + esi]`

GAS: `movb (%ebx, %esi, 1), %al`

1.1.3 The PC's Physical Address Space

The approximate physical distribution of the physical memory of the pc is as follows

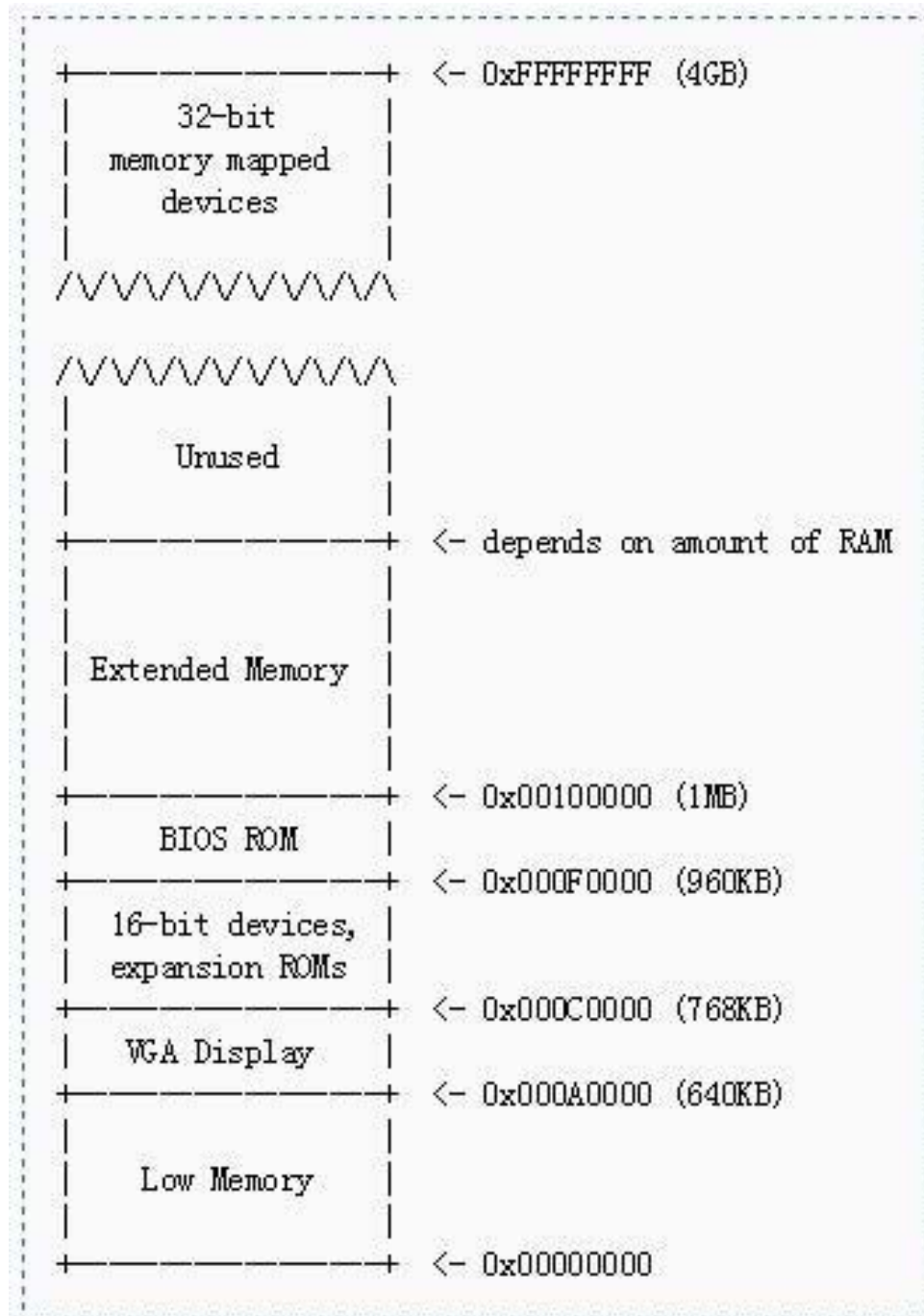


图 1.3: i 文件部分代码

The 384KB of physical address range 0x000A0000 0x000FFFFFFF is reserved for hardware devices such as VGA display buffers. The most important part of this part is as the basic input and output system (BIOS, Basic)Input/Output System) The 64KB area of 0x000F0000 0x000FFFFFFF. The BIOS handles the basic initialization tasks of the system, such as activating the display, checking memory, and so on. After completing these initialization tasks, the BIOS loads the operating system from a floppy disk, hard drive,

optical drive, or network and transfers control to the operating system.

On the modern PC, the 0x000A0000 0x00100000 segment is left as a "hole", which divides the memory into two segments. The lower 640K is called "low-segment memory", and the remaining high-address portion is expanded memory. At the same time, in the highest part of the physical address space, above all physical RAM, some are also reserved by the BIOS for 32-bit PCI devices.

1.1.4 BIOS

```

lyh@ubuntu: ~/Downloads/lab1/src/lab1_1
GNU gdb (Ubuntu 7.11.1-0ubuntu1~16.5) 7.11.1
Copyright (C) 2016 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 "x86_64-linux-gnu".
Type "show configuration" for configuration details.
For bug reporting instructions, please see:
<http://www.gnu.org/software/gdb/bugs/>.
Find the GDB manual and other documentation resources online at:
<http://www.gnu.org/software/gdb/documentation/>.
For help, type "help".
Type "apropos word" to search for commands related to "word".
+ target remote localhost:26000
warning: A handler for the OS ABI "GNU/Linux" is not built into this configurati
on
of GDB. Attempting to continue with the default i8086 settings.

The target architecture is assumed to be i8086
[f000:ffff] 0xfffff0: jmp $0xf000,$0xe05b
0x0000ffff in ?? ()
+ symbol-file obj/kern/kernel
(gdb)
    
```

图 1.4: gdb

In the running content of gdb, we can see the first instruction that gdb runs (as shown below). In this instruction, we can see that pc starts from CS = 0xf000 and IP = 0xffff0. The execution of the first sentence is a JMP operation that jumps to CS = 0xf000 and IP = 0xe05b. Since the modern CPU is divided into real mode and protection mode, it runs in real mode at startup and runs in protected mode after startup. The BIOS is the software that runs when the PC first starts up, so it must work in real mode. So the real address of the jump is $0xf000 \ll 4 + 0xe05b = 0xfe05b$.

```

The target architecture is assumed to be i8086
[f000:ffff] 0xfffff0: jmp $0xf000,$0xe05b
0x0000ffff in ?? ()
+ symbol-file obj/kern/kernel
(gdb)
    
```

图 1.5: The first instruction of the BIOS

1.2 Boot Loader

For PCs, floppy disks and hard disks can be divided into 512-byte areas called sectors. A sector is the smallest granularity of a disk operation. Each read or write operation must be one or more sectors. If a disk can be used to boot the operating system, the first sector of the disk is called the boot sector. The boot loader program described in this section is located in this boot sector. When the BIOS finds a floppy disk or hard disk that can be booted, it loads the 512-byte boot sector into the memory address 0x7c00 0x7dff. The boot loader must perform two main functions.

1. First, the boot loader will convert the processor from real mode to 32-bit protected mode, because only in this mode the software can access more than 1MB of content.
2. The boot loader can then access the IDE disk device registers directly from the disk by using x86 specific IO instructions.

For the boot loader, there is a file that is important, `obj/boot/boot.asm`. This file is a disassembled version of our real-run boot loader program. So we can compare it with its source code, `boot.S` and `main.c`.

1.2.1 Exercise II

Question

Use the debug function of GDB to complete the boot loader tracking process. The `b` command sets a breakpoint, such as `b *0x7c00` to set a breakpoint at 0x7c00; the `c` command to execute to the next breakpoint; the `si` command is single-step.

Answer

```
0x0000fff0 in ?? ()
+ symbol-file obj/kern/kernel
(gdb) b *0x7c00
Breakpoint 1 at 0x7c00
(gdb) c
Continuing.
[ 0:7c00] => 0x7c00: cli
Breakpoint 1, 0x00007c00 in ?? ()
(gdb) si
[ 0:7c01] => 0x7c01: cld
0x00007c01 in ?? ()
(gdb) si
[ 0:7c02] => 0x7c02: xor    %ax,%ax
0x00007c02 in ?? ()
(gdb) █
```

According to the prompt, we use the `b` and `c` instructions to set the breakpoint and execute to the next breakpoint, and use `si` to step through the debug.

1.2.2 Question I

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.

- 1) At what point does the processor start executing 32-bit code? What exactly causes the switch from 16- to 32-bit mode?
- 2) What is the last instruction of the boot loader executed, and what is the first instruction of the kernel it just loaded?
- 3) Where is the first instruction of the kernel?
- 4) 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?

The BIOS will copy the boot sector to the address `0x7c00`, so the starting address of `boot.S` is `0x7c00`. So we enter `b *0x7c00` in the gdb window, then enter `c`, which means continue to run to the breakpoint, where we enter `x/30i 0x7c00`. This gdb instruction disassembles the instructions stored in `0x7c00` and the next 30 bytes of memory.

```

The target architecture is assumed to be i8086
[f000:fff0] 0xffff0: jmp $0xf000,$0xe05b
0x0000fff0 in ?? ()
+ symbol-file obj/kern/kernel
(gdb) b *0x7c00
Breakpoint 1 at 0x7c00
(gdb) c
Continuing.
[ 0:7c00] => 0x7c00: cli

Breakpoint 1, 0x00007c00 in ?? ()
(gdb) x/30i 0x7c00
=> 0x7c00: cli
0x7c01: cld
0x7c02: xor %ax,%ax
0x7c04: mov %ax,%ds
0x7c06: mov %ax,%es
0x7c08: mov %ax,%ss
0x7c0a: in $0x64,%al
0x7c0c: test $0x2,%al
0x7c0e: jne 0x7c0a
0x7c10: mov $0xd1,%al
0x7c12: out %al,$0x64
0x7c14: in $0x64,%al
0x7c16: test $0x2,%al
0x7c18: jne 0x7c14
0x7c1a: mov $0xdf,%al
0x7c1c: out %al,$0x60
0x7c1e: lgdtw 0x7c64
0x7c23: mov %cr0,%eax
0x7c26: or $0x1,%eax
0x7c2a: mov %eax,%cr0
0x7c2d: ljmp $0x8,$0x7c32
0x7c32: mov $0xd88e0010,%eax
0x7c38: mov %ax,%es
---Type <return> to continue, or q <return> to quit---

```

图 1.6: Disassembly

We take it directly with boot.S and at obj/boot/boot.asm

```

obj/boot/boot.out:      file format elf32-i386

Disassembly of section .text:

00007c00 <start>:
.set CR0_PE_ON,      0x1          # protected mode enable flag

.globl start
start:
.code16                        # Assemble for 16-bit mode
cli                          # Disable interrupts
7c00:      fa          cli
7c01:      fc          # String operations increment
                        cld

# Set up the important data segment registers (DS, ES, SS).
xorw  %ax,%ax                # Segment number zero
7c02:      31 c0          xor  %eax,%eax
movw  %ax,%ds                # -> Data Segment
7c04:      8e d8          mov  %eax,%ds
movw  %ax,%es                # -> Extra Segment
7c06:      8e c0          mov  %eax,%es
movw  %ax,%ss                # -> Stack Segment
7c08:      8e d0          mov  %eax,%ss

00007c0a <seta20.1>:
# Enable A20:
# For backwards compatibility with the earliest PCs, physical
# address line 20 is tied low, so that addresses higher than
# 1MB wrap around to zero by default. This code undoes this.
seta20.1:
inb  $0x64,%al              # Wait for not busy
7c0a:      e4 64          in   $0x64,%al
testb $0x2,%al
7c0c:      a8 02          test  $0x2,%al
jnz  seta20.1
7c0e:      75 fa          jne  7c0a <seta20.1>

movb  $0xd1,%al              # 0xd1 -> port 0x64
7c10:      b0 d1          mov  $0xd1,%al
outb  %al,$0x64
7c12:      e6 64          out  %al,$0x64

00007c14 <seta20.2>:

```

图 1.7: boot.asm

```

boot.asm

#include <inc/mmu.h>

# Start the CPU: switch to 32-bit protected mode, jump into C.
# The BIOS loads this code from the first sector of the hard disk into
# memory at physical address 0x7c00 and starts executing in real mode
# with %cs=0 %ip=7c00.

.set PROT_MODE_CSEG, 0x8      # kernel code segment selector
.set PROT_MODE_DSEG, 0x10     # kernel data segment selector
.set CR0_PE_ON,      0x1      # protected mode enable flag

.globl start
start:
.code16                        # Assemble for 16-bit mode
cli                          # Disable interrupts
cld                          # String operations increment

# Set up the important data segment registers (DS, ES, SS).
xorw  %ax,%ax                # Segment number zero
movw  %ax,%ds                # -> Data Segment
movw  %ax,%es                # -> Extra Segment
movw  %ax,%ss                # -> Stack Segment

# Enable A20:
# For backwards compatibility with the earliest PCs, physical
# address line 20 is tied low, so that addresses higher than
# 1MB wrap around to zero by default. This code undoes this.
seta20.1:
inb  $0x64,%al              # Wait for not busy
testb $0x2,%al
jnz  seta20.1

movb  $0xd1,%al              # 0xd1 -> port 0x64
outb  %al,$0x64

seta20.2:
inb  $0x64,%al              # Wait for not busy
testb $0x2,%al
jnz  seta20.2

movb  $0xdf,%al              # 0xdf -> port 0x60
outb  %al,$0x60

```

图 1.8: boot.s

By comparing we can find that the three are not different in the instruction, but in the source code, we specify a lot of identifiers such as set20.1. Initially, these identifiers are converted to real physical addresses after being assembled into machine code. For example, set20.1 is converted to 0x7c0a, then this correspondence is listed in OBJ's /boot/boot.asm, but in the real case, in the first case, you can't see set20. 1 identifier, completely real physical address.

Then according to the title indication, we first traced to the bootmain function, we found that bootmain is at address 0x7c45, so we set the breakpoint there, and run here, then set the breakpoint at the readsect (0x7c7c) and jump here After using the si command to step through, found that into the waitdisk function (0x7c6a)

```
(gdb) b* 0x7c45
Breakpoint 2 at 0x7c45
(gdb) c
Continuing.
The target architecture is assumed to be i386
=> 0x7c45:      call    0x7d15

Breakpoint 2, 0x00007c45 in ?? ()
(gdb) b *0x7c7c
Breakpoint 3 at 0x7c7c
(gdb) c
Continuing.
=> 0x7c7c:      push    %ebp

Breakpoint 3, 0x00007c7c in ?? ()
(gdb) si
=> 0x7c7d:      mov     %esp,%ebp
0x00007c7d in ?? ()
(gdb) si
=> 0x7c7f:      push    %edi
0x00007c7f in ?? ()
(gdb) si
=> 0x7c80:      mov     0xc(%ebp),%ecx
0x00007c80 in ?? ()
(gdb) si
=> 0x7c83:      call    0x7c6a
0x00007c83 in ?? ()
(gdb) █
```

图 1.9: call bootmain

In boot.asm we can find the corresponding content of 0x7d6b. Jumping to here means

`((void (*)(void)) (ELFHDR->e_entry))()` the function is executed, where the meaning of the `e_entry` field of the ELF header is the first instruction of the executable. Virtual address. So the meaning of this sentence is to transfer control to the operating system kernel. Then we find the end of the for statement in the file and jump to this location (0x7d6b).

```

        // note: does not return!
        ((void (*)(void)) (ELFHDR->e_entry))();
7d6b:      ff 15 18 00 01 00      call    *0x10018
    }

```

图 1.10: end loop

```

        // call the entry point from the ELF header
        // note: does not return!
        ((void (*)(void)) (ELFHDR->e_entry))();|
7d6b:      ff 15 18 00 01 00      call    *0x10018

```

图 1.11: the instruction after loop

Answer:

1)

As we discussed earlier, When the PC starts up, the CPU runs in real mode (real mode), and when it enters the operating system kernel, it will run in protected mode (protected mode). Then it will set PE. When entering protection mode, the CPU will switch from 16 bits to 32 bits.

```

# Switches processor into 32-bit mode.
ljmp    $PROT_MODE_CSEG, $protcseg
7c2d:      ea                      .byte 0xea
7c2e:      32 7c 08 00            xor    0x0(%eax,%ecx,1),%bh

```

图 1.12: mode change

In the `boot.S` file, the computer first works in real mode, this time is the 16bit working mode. As can be seen from the above figure, when the `"ljmp $PROT_MODE_CSEG, $protcseg"` statement is run, the 32-bit working mode is officially entered. The root cause is that the CPU is working in protected mode at this time.

2)


```

// call the entry point from the ELF header
// note: does not return!
((void (*)(void)) (ELFHDR->e_entry))();|
7d6b:      ff 15 18 00 01 00      call    *0x10018

```

图 1.13: the instruction after loop

As shown in the figure above, the last statement executed by the boot loader is the last statement in the bootmain subroutine “((void (*)(void)) (ELFHDR->e_entry))();”, that is, jump to the operating system. The starting instruction of the kernel program. This first instruction is located in the /kern/entry.S file, the first sentence movw \$0x1234, 0x472

3)

The first instruction is in the /kern/entry.S file.

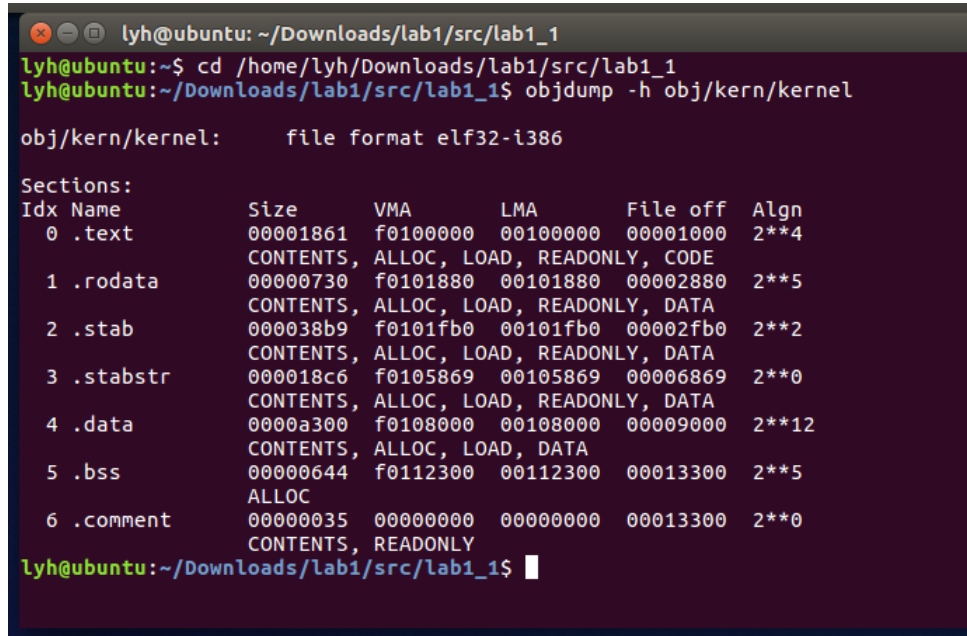
4)

First, how many segments are shared by the operating system, and how many sectors are in each segment are located in the Program Header Table in the operating system file. Each entry in this table corresponds to a segment of the operating system. And the content of each entry includes the size of the segment, the segment start address offset, and the like. So if we can find this table, we can determine how many sectors the kernel occupies by the information provided by the table entry. Then the information about where this table is stored is stored in the ELF header information of the operating system kernel image file.

1.2.3 Loading kernel

An executable ELF file consists of three main parts: a file header with loading information, followed by a program segment table, followed by several program segments. Each of these segments is a piece of continuous code or data. They are first loaded into memory when they are run. The job of the boot loader is to load them into memory. We can use the following instructions to examine the names, sizes, and addresses of all segments in the JOS kernel.

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



```

lyh@ubuntu: ~/Downloads/lab1/src/lab1_1
lyh@ubuntu:~$ cd /home/lyh/Downloads/lab1/src/lab1_1
lyh@ubuntu:~/Downloads/lab1/src/lab1_1$ objdump -h obj/kern/kernel

obj/kern/kernel:      file format elf32-i386

Sections:
Idx Name              Size      VMA       LMA       File off  Algn
 0  .text              00001861  f0100000  00100000  00001000  2**4
    CONTENTS, ALLOC, LOAD, READONLY, CODE
 1  .rodata            00000730  f0101880  00101880  00002880  2**5
    CONTENTS, ALLOC, LOAD, READONLY, DATA
 2  .stab              000038b9  f0101fb0  00101fb0  00002fb0  2**2
    CONTENTS, ALLOC, LOAD, READONLY, DATA
 3  .stabstr           000018c6  f0105869  00105869  00006869  2**0
    CONTENTS, ALLOC, LOAD, READONLY, DATA
 4  .data              0000a300  f0108000  00108000  00009000  2**12
    CONTENTS, ALLOC, LOAD, DATA
 5  .bss               00000644  f0112300  00112300  00013300  2**5
    ALLOC
 6  .comment           00000035  00000000  00000000  00013300  2**0
    CONTENTS, READONLY
lyh@ubuntu:~/Downloads/lab1/src/lab1_1$

```

图 1.14: All segment information in the kernel

1.2.4 Connection address and load address

There are two more important fields in each segment, VMA (link address), LMA (load address). The load address represents the physical address of the segment after it is loaded into memory. The link address refers to the logical address to which this segment is expected to be stored.

Each ELF file has a Program Headers Table that indicates which parts of the ELF file are loaded into memory and the addresses that are loaded into memory. We get the information of the Program Headers Table of the kernel by entering the following instructions:

```
objdump -x obj/kern/kernel
```

```

                                CONTENTS, READONLY
lyh@ubuntu:~/Downloads/lab1/src/lab1_1$ objdump -x obj/kern/kernel

obj/kern/kernel:      file format elf32-i386
obj/kern/kernel
architecture: i386, flags 0x00000112:
EXEC_P, HAS_SYMS, D_PAGED
start address 0x0010000c

Program Header:
  LOAD off 0x00001000 vaddr 0xf0100000 paddr 0x00100000 align 2**12
        filesz 0x0000712f memsz 0x0000712f flags r-x
  LOAD off 0x00009000 vaddr 0xf0108000 paddr 0x00108000 align 2**12
        filesz 0x0000a300 memsz 0x0000a944 flags rw-
  STACK off 0x00000000 vaddr 0x00000000 paddr 0x00000000 align 2**4
        filesz 0x00000000 memsz 0x00000000 flags rwx

Sections:
Idx Name          Size      VMA       LMA       File off  Algn
  0 .text          00001861 f0100000 00100000 00001000  2**4
                CONTENTS, ALLOC, LOAD, READONLY, CODE
  1 .rodata        00000730 f0101880 00101880 00002880  2**5
                CONTENTS, ALLOC, LOAD, READONLY, DATA
  2 .stab          000038b9 f0101fb0 00101fb0 00002fb0  2**2
                CONTENTS, ALLOC, LOAD, READONLY, DATA
  3 .stabstr        000018c6 f0105869 00105869 00006869  2**0
                CONTENTS, ALLOC, LOAD, READONLY, DATA
  4 .data           0000a300 f0108000 00108000 00009000  2**12
                CONTENTS, ALLOC, LOAD, DATA
  5 .bss            00000644 f0112300 00112300 00013300  2**5
                ALLOC
  6 .comment        00000035 00000000 00000000 00013300  2**0
                CONTENTS, READONLY

SYMBOL TABLE:
f0100000 l d .text 00000000 .text
f0101880 l d .rodata 00000000 .rodata
f0101fb0 l d .stab 00000000 .stab
f0105869 l d .stabstr 00000000 .stabstr
f0108000 l d .data 00000000 .data
f0112300 l d .bss 00000000 .bss
00000000 l d .comment 00000000 .comment
00000000 l df *ABS* 00000000 obj/kern/entry.o
f010002f l .text 00000000 relocated
f010003e l .text 00000000 spin
00000000 l df *ABS* 00000000 entrypgdir.c
00000000 l df *ABS* 00000000 init.c
00000000 l df *ABS* 00000000 console.c
f0100177 l F .text 0000001f serial_proc_data
f0100196 l F .text 00000043 cons_intr
f0112320 l 0 .bss 00000208 cons

```

图 1.15: Kernel's Program Headers Table information

1.3 The Kernel

1.3.1 Question

问题:

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          memcpy 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      }
```

Answer:

1)

`Console.c` defines how to display a character to the console, which is on top of our display, which includes many operations on the IO port. What is defined in `printf.c` is the top-level formatted output subroutine we will use in programming, such as `printf`, `sprintf`, and so on. The function exported by `console.c` which is used by `printf.c` is `cputchar()`. That function prints a character in the parallel port and in the display.

Specific call relationship: `cprintf` -> `vcprintf` -> `putch` -> `cputchar`. Kernel's `cprintf()` function calls `vcprintfmt()` (from `lib/printfmt.c`) to actually print in the console, `vcprintfmt()` does the needed formatting and then call a function passed to it to actually print in the display.

```
#include <inc/types.h>
#include <inc/stdio.h>
#include <inc/stdarg.h>

static void
putch(int ch, int *cnt)
{
    cputchar(ch);
    *cnt++;
}
```

图 1.16: Function call

2)

Crt_buf: This is a character array buffer that holds the characters to be displayed on the screen.

Crt_pos: This indicates the position of the current last character displayed on the screen.

When `crt_pos >= CRT_SIZE`, where `CRT_SIZE = 80*25`, since we know that the value range of `crt_pos` is 0 ($80*25-1$), if this condition is true, it means that the content output on the screen has exceeded one page. . So at this point you have to scroll the page up one line, that is, put the original line 1 79 on the current 0 78 line, and then replace the line 79 with a line of spaces (of course not all spaces, 0 characters) To display the character `int c`) you entered. So the `memcpy` operation is to copy the contents of lines 1 79 in the `crt_buf` character array to the position of lines 0 78. The next for loop is to turn the last line, line 79 into a space. Finally, you need to modify the value of `crt_pos`.

1.3.2 Use virtual memory

As we discussed earlier, the computer is divided into real mode and protected mode when it starts up. When running the boot loader, the link address (virtual address) and the load address (physical address) in the boot loader are the same. But when you enter the kernel, the two addresses are no longer the same. In the virtual address space, we put the operating system at the high address `0xf0100000`, but in the actual memory we store the operating system in a low physical address space, such as `0x00100000`. Then when

the user program wants to access an instruction of an operating system kernel, the first is to give a high virtual address, and then the virtual address is mapped to a real physical address by a certain mechanism in the computer, thus solving the above problem. . Then such an organization is usually implemented through segment management and paging management.

1.3.3 Homework I

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.

Answer

To answer this question, we should first understand the three files \ kern \ printf.c, \ kern \ console.c, \ lib \ printfmt.c to understand the relationship between them. First of all, we should pay attention to the console.c file, the most important of which is the cputchar subroutine. We can find that this program is the IO control program of the high-level console. In addition, the implementation of cputchar is actually done by calling cons_putc. And the function of the cons_putc program is to output a character to the console.

```
// 'High'-level console I/O.  Used by readline and cprintf.
void
cputchar(int c)
{
    cons_putc(c);
}
```

图 1.17: cputchar

```
// output a character to the console
static void
cons_putc(int c)
{
    serial_putc(c);
    lpt_putc(c);
    cga_putc(c);
}
```

图 1.18: cons_putc

Then, let's focus on the `printfmt.c` file. By commenting, we can see that the subroutine defined in this file is the key to the information we can use to directly output information to the screen during programming using the `printf` function.

```
console.c
// Stripped-down primitive printf-style formatting routines,
// used in common by printf, sprintf, fprintf, etc.
// This code is also used by both the kernel and user programs.

#include <inc/types.h>
#include <inc/stdio.h>
#include <inc/string.h>
#include <inc/stdarg.h>
#include <inc/error.h>
```

图 1.19: printfmt

Finally, let's take a look at the `printf.c` file. By commenting, we can see that the function of this file is to implement the simple implementation of `cprintf` console output for the kernel.

```
// Simple implementation of cprintf console output for the kernel,
// based on printfmt() and the kernel console's cputchar().

#include <inc/types.h>
#include <inc/stdio.h>
#include <inc/stdarg.h>

static void
putch(int ch, int *cnt)
{
    cputchar(ch);
    *cnt++;
}

int
vcprintf(const char *fmt, va_list ap)
{
    int cnt = 0;

    vprintfmt((void*)putch, &cnt, fmt, ap);
    return cnt;
}

int
cprintf(const char *fmt, ...)
{
    va_list ap;
    int cnt;

    va_start(ap, fmt);
    cnt = vcprintf(fmt, ap);
    va_end(ap);

    return cnt;
}
```

图 1.20: printf

In this question, we follow the case 'u' to write the code.

```
// (unsigned) octal
case 'o':
    // Replace this with your code.
    /*putch('X', putdat);
    putch('X', putdat);
    putch('X', putdat);
    break;*/

    num = getuint(&ap, lflag);
    base = 8;
    goto number;
```

图 1.21: answer

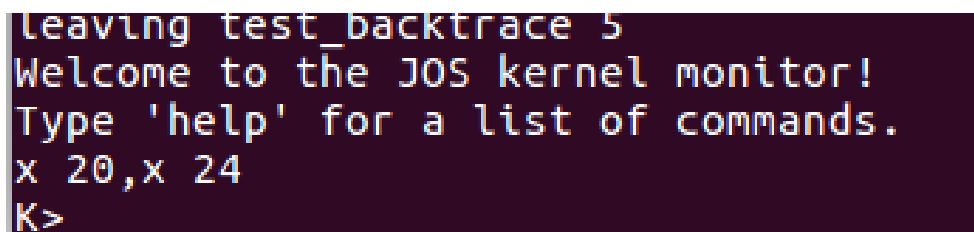
After that we modify the monitor.c file and then run qemu on the terminal to test our results.

```
void
monitor(struct Trapframe *tf)
{
    char *buf;

    cprintf("Welcome to the JOS kernel monitor!\n");
    cprintf("Type 'help' for a list of commands.\n");
    int x = 20;
    cprintf("x %d,x %o\n",x,x);

    while (1) {
        buf = readline("K> ");
        if (buf != NULL)
            if (runcmd(buf, tf) < 0)
                break;
    }
}
```

图 1.22: Change monitor



```
Leaving test_backtrace 5
Welcome to the JOS kernel monitor!
Type 'help' for a list of commands.
x 20,x 24
K>
```

图 1.23: Result test

As can be seen from the results, the 20 we defined is under %d, and the %o is correctly changed to octal 24

1.3.4 Exercise III

Question

To become familiar with the C calling conventions on the x86, 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 32-bits words does each recursive nesting level of `test_backtrace` push on the stack, and what are those words?

Answer

First, let's take a look at the source code of the `test_backtrace` function in `kernel.asm`

```
// test the stack backtrace function (tab 1 only)
void
test_backtrace(int x)
{
f0100040:      55                push    %ebp
f0100041:      89 e5             mov     %esp,%ebp
f0100043:      53                push    %ebx
f0100044:      83 ec 0c          sub     $0xc,%esp
f0100047:      8b 5d 08          mov     0x8(%ebp),%ebx
      cprintf("entering test_backtrace %d\n", x);
f010004a:      53                push    %ebx
f010004b:      68 00 18 10 f0    push    $0xf0101800
f0100050:      e8 a6 08 00 00    call    f01008fb <cstdio>
      if (x > 0)
f0100055:      83 c4 10          add     $0x10,%esp
f0100058:      85 db             test    %ebx,%ebx
f010005a:      7e 11             jle     f010006d <test_backtrace+0x2d>
      test_backtrace(x-1);
f010005c:      83 ec 0c          sub     $0xc,%esp
f010005f:      8d 43 ff          lea     -0x1(%ebx),%eax
f0100062:      50                push    %eax
f0100063:      e8 d8 ff ff ff    call    f0100040 <test_backtrace>
f0100068:      83 c4 10          add     $0x10,%esp
f010006b:      eb 11             jmp     f010007e <test_backtrace+0x3e>
      else
      mon_backtrace(0, 0, 0);
f010006d:      83 ec 04          sub     $0x4,%esp
f0100070:      6a 00             push    $0x0
f0100072:      6a 00             push    $0x0
f0100074:      6a 00             push    $0x0
f0100076:      e8 e5 06 00 00    call    f0100760 <mon_backtrace>
f010007b:      83 c4 10          add     $0x10,%esp
      cprintf("leaving test_backtrace %d\n", x);
f010007e:      83 ec 08          sub     $0x8,%esp
f0100081:      53                push    %ebx
f0100082:      68 1c 18 10 f0    push    $0xf010181c
f0100087:      e8 6f 08 00 00    call    f01008fb <cstdio>
}
f010008c:      83 c4 10          add     $0x10,%esp
f010008f:      8b 5d fc          mov     -0x4(%ebp),%ebx
f0100092:      c9                leave
f0100093:      c3                ret
f0100094 <i386_init>:
```

From the code we can see the address of `test_backtrace`, so we set a breakpoint here, then jump to here, as shown below:

```

lyh@ubuntu: ~/Downloads/lab1/src/lab1_1
(gdb) b *0xf0100040
Breakpoint 1 at 0xf0100040: file kern/init.c, line 13.
(gdb) c
Continuing.
The target architecture is assumed to be i386
=> 0xf0100040 <test_backtrace>: push  %ebp

Breakpoint 1, test_backtrace (x=5) at kern/init.c:13
13      {

```

Then we use the `i r` instruction to view the register contents and repeat the above operation

```

(gdb) i r
eax                0x0          0
ecx                0x3d4        980
edx                0x3d5        981
ebx                0x10094      65684
esp                0xf010ffdc    0xf010ffdc
ebp                0xf010fff8    0xf010fff8
esi                0x10094      65684
edi                0x0          0
eip                0xf0100040    0xf0100040 <test_backtrace>
eflags             0x46         [ PF ZF ]
cs                 0x8          8
ss                 0x10         16
ds                 0x10         16
es                 0x10         16
fs                 0x10         16
gs                 0x10         16
(gdb) c
Continuing.
=> 0xf0100040 <test_backtrace>: push  %ebp

Breakpoint 1, test_backtrace (x=4) at kern/init.c:13
13      {
(gdb) ir
Undefined command: "ir".  Try "help".
(gdb) i r
eax                0x4          4
ecx                0x3d4        980
edx                0x3d5        981
ebx                0x5          5
esp                0xf010ffbc    0xf010ffbc
ebp                0xf010ffd8    0xf010ffd8
esi                0x10094      65684
edi                0x0          0
eip                0xf0100040    0xf0100040 <test_backtrace>
eflags             0x92         [ AF SF ]
cs                 0x8          8
ss                 0x10         16
ds                 0x10         16
es                 0x10         16
fs                 0x10         16
gs                 0x10         16
(gdb) █

```

By comparison, we can find that the value of `ebp` is reduced by 20 (hexadecimal), so we can judge every time it pushes 8 4-byte words.

According to the structure of the `test_backtrace` function and the change of the stack when the function is called, each time the function is called, the `ebp`, the return address, the saved `ebx` value, and the parameters of the next call function must be pushed onto the stack. There are also 4 reserved words.

1.3.5 Homework II

Question

You can do `mon_backtrace()` entirely in C. 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:

```
Waring:
read_ebp
display format
```

Answer

First we call the `read_ebp` function to get the value of the current `EBP` register. We treat the entire call stack as an array, `EBP[0]` represents the `ebp` value of the previous function, and `EBP [1]` stores the function return address, `EBP [2]` What is stored later is the value of the input parameter.

Ebp of the previous function
ebx
Parameter 1
Parameter 2
Parameter 3
Parameter 4
Parameter 5

表 1.1: Stack structure diagram

The modified `mon_backtrace` function code is shown in the figure

```

int
mon_backtrace(int argc, char **argv, struct Trapframe *tf)
{

    uint32_t ebp=read_ebp();
    uint32_t *tmp_ebp=(uint32_t *)ebp;
    int i,j;

    cprintf("Stack backtrace:\n");
    while(tmp_ebp)
    {
        cprintf("ebp %x ",tmp_ebp);
        cprintf("eip %x ",tmp_ebp[1]);
        cprintf("args ");
        for(i=2;i<7;i++)
        {
            cprintf("%08x ",tmp_ebp[i]);
        }

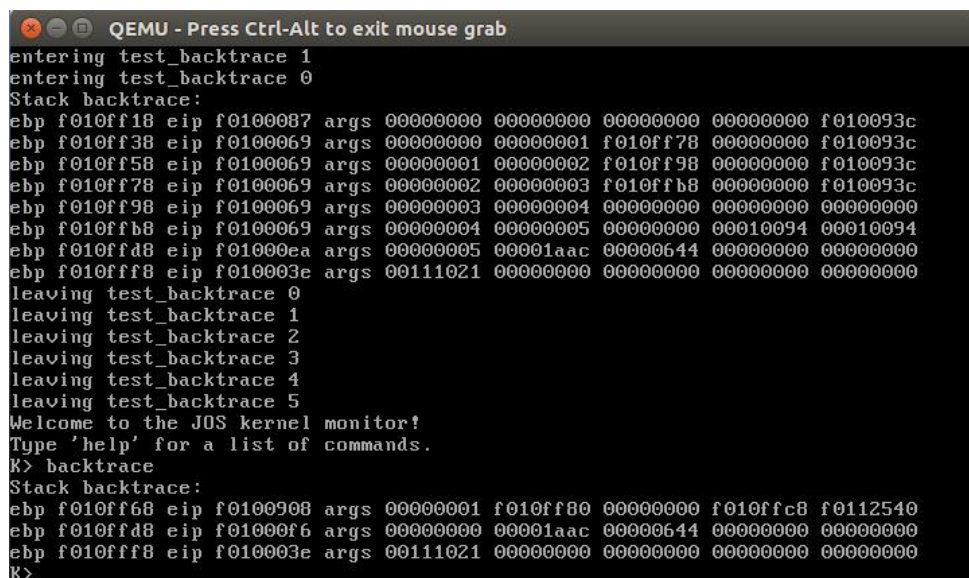
        cprintf("\n");
        tmp_ebp=(uint32_t *)(*tmp_ebp);
    }

    return 0;
}

```

After analysis, we can know that when the function is called, the parameters are first pushed onto the stack, then the eip is pushed onto the stack, and finally the ebp is pushed onto the stack, and the new ebp points to the position of the old ebp. Therefore, the location pointed to by ebp stores the value of the previous function ebp, and the location of ebp+1 stores the value of eip. So we first print out tmp_ebp as ebp, tmp_ebp[1] as eip. Then we print out the values of 5 parameters according to the requirements of the topic. Finally, we take the value in the memory space pointed to by ebp and repeat the above process until the first function is called.

The result of the operation is as follows



```

QEMU - Press Ctrl-Alt to exit mouse grab
entering test_backtrace 1
entering test_backtrace 0
Stack backtrace:
ebp f010ff18 eip f0100087 args 00000000 00000000 00000000 00000000 f010093c
ebp f010ff38 eip f0100069 args 00000000 00000001 f010ff78 00000000 f010093c
ebp f010ff58 eip f0100069 args 00000001 00000002 f010ff98 00000000 f010093c
ebp f010ff78 eip f0100069 args 00000002 00000003 f010ffb8 00000000 f010093c
ebp f010ff98 eip f0100069 args 00000003 00000004 00000000 00000000 00000000
ebp f010ffb8 eip f0100069 args 00000004 00000005 00000000 00010094 00010094
ebp f010ffd8 eip f01000ea args 00000005 00001aac 00000644 00000000 00000000
ebp f010fff8 eip f010003e args 00111021 00000000 00000000 00000000 00000000
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> backtrace
Stack backtrace:
ebp f010ff68 eip f0100908 args 00000001 f010ff80 00000000 f010ffc8 f0112540
ebp f010ffd8 eip f01000f6 args 00000000 00001aac 00000644 00000000 00000000
ebp f010fff8 eip f010003e args 00111021 00000000 00000000 00000000 00000000
K>

```

Add backtrace to the command list

```
static struct Command commands[] = {
    { "help", "Display this list of commands", mon_help },
    { "kerninfo", "Display information about the kernel", mon_kerninfo },
    { "backtrace", "Display stack backtrace", mon_backtrace },
};
```

1.3.6 Challenge homework

Question

Modify your stack backtrace function to display, for each eip, the function name, source file name, and line number corresponding to that eip. In `debuginfo_eip`, where do `__STAB__` come from?

Answer

After reading the contents of `stab.h` and using the `objdump -G kernel` command, you can roughly know the fields of the stab table that store debugging information and their meanings, as follows

`UInt32_t n_strx`; index, can be added to `stabstr` to get the corresponding address of the string information.

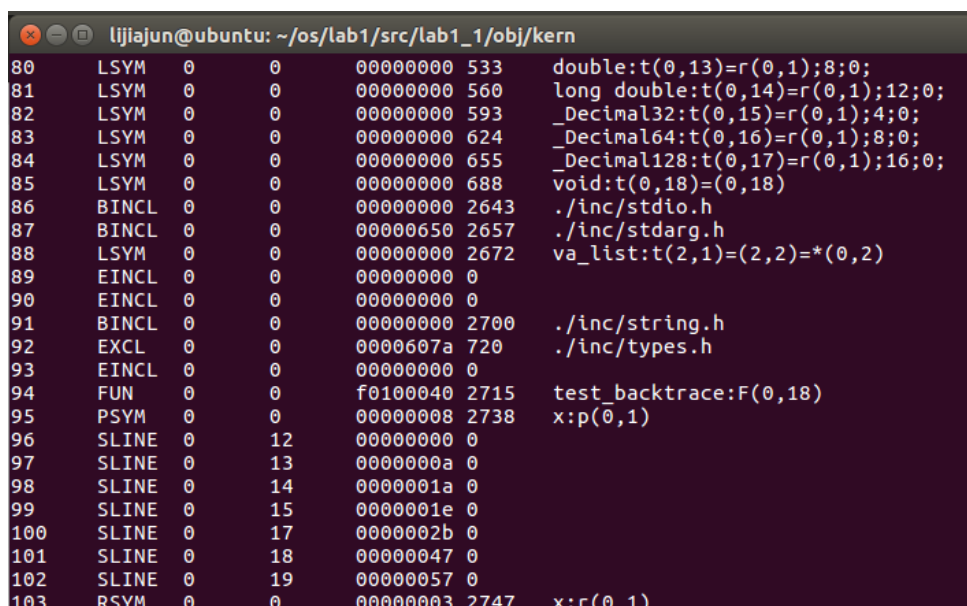
`UInt8_t n_type`; The type of the line content, such as functions, function parameters, a line in the source file, etc.

`UInt8_t n_other`; usually empty.

`UInt16_t n_desc`; Stores the specific number of rows when type is `n_sline`.

`UIntptr_t n_value`; The position of the corresponding content runtime in memory (or relative to the location where the function starts the instruction).

Collect the contents of the header file and the stab table, and analyze the meaning of each field as follows:



```
lijiajun@ubuntu: ~/os/lab1/src/lab1_1/obj/kern
80  LSYM  0  0  00000000 533  double:t(0,13)=r(0,1);8;0;
81  LSYM  0  0  00000000 560  long double:t(0,14)=r(0,1);12;0;
82  LSYM  0  0  00000000 593  _Decimal32:t(0,15)=r(0,1);4;0;
83  LSYM  0  0  00000000 624  _Decimal64:t(0,16)=r(0,1);8;0;
84  LSYM  0  0  00000000 655  _Decimal128:t(0,17)=r(0,1);16;0;
85  LSYM  0  0  00000000 688  void:t(0,18)=(0,18)
86  BINCL 0  0  00000000 2643  ./inc/stdio.h
87  BINCL 0  0  00000650 2657  ./inc/stdarg.h
88  LSYM  0  0  00000000 2672  va_list:t(2,1)=(2,2)=*(0,2)
89  EINCL 0  0  00000000 0
90  EINCL 0  0  00000000 0
91  BINCL 0  0  00000000 2700  ./inc/string.h
92  EXCL  0  0  0000607a 720  ./inc/types.h
93  EINCL 0  0  00000000 0
94  FUN  0  0  f0100040 2715  test_backtrace:F(0,18)
95  PSYM  0  0  00000008 2738  x:p(0,1)
96  SLINE 0  12  00000000 0
97  SLINE 0  13  0000000a 0
98  SLINE 0  14  0000001a 0
99  SLINE 0  15  0000001e 0
100 SLINE 0  17  0000002b 0
101 SLINE 0  18  00000047 0
102 SLINE 0  19  00000057 0
103 RSYM 0  0  00000003 2747  x:r(0,1)
```

The function for finding symbol information is `debuginfo_eip(uintptr_t addr, struct Eipdebuginfo *info)`. Analyze `kdebug.h`, we found that the function of this function is to fill the corresponding fields of the `Eipdebuginfo` structure. Analyze the `debuginfo_eip` function, we find that its execution logic is roughly: the structure of the `eip` value to be searched and the information to save the result. In the `stab` table, a binary search is performed on the item whose `n_type` is `file`, and it is found that the value of the `eip` is in the middle of the label of the two source files (the first item in the figure). Then, the labels of the two source files are used as the left and right intervals, and the items in which the type is `fun` are binary searched, and the value is found between the two function labels. Finally, the `eip` value is subtracted from the left interval function value, and a binary search is performed to find the corresponding `eip` line number. The added code is as follows:

```
// Your code here.
stab_binsearch(stabs,&lline,&rline, N_SLINE,addr);
if(lline==rline)
{
    info->eip_line=stabs[lline].n_desc;
}
else{
    return -1;
}
```

Modify the `mon_backtrace` code to output the search result

```
int
mon_backtrace(int argc, char **argv, struct Trapframe *tf)
{
    uint32_t ebp=read_ebp();
    uint32_t *tmp_ebp=(uint32_t *)ebp;
    int i,j;
    struct Eipdebuginfo *info;
    struct Eipdebuginfo infos;
    info=&infos;
    cprintf("Stack backtrace:\n");
    while(tmp_ebp)
    {
        cprintf("ebp %x ",tmp_ebp);
        cprintf("eip %x ",tmp_ebp[1]);
        cprintf("args ");
        for(i=2;i<7;i++)
        {
            cprintf("%08x ",tmp_ebp[i]);
        }

        debuginfo_eip(tmp_ebp[1],info);
        cprintf("\n");
        cprintf("\t%s:%d: ",info->eip_file,info->eip_line);
        for(j=0;j<info->eip_fn_namelen;j++)
            cprintf("%c",info->eip_fn_name[j]);
        cprintf(" +%d",tmp_ebp[1]-info->eip_fn_addr);
        cprintf("\n");
        tmp_ebp=(uint32_t *)(*tmp_ebp);
    }

    return 0;
}
```

Eip_fn_namelen is used to exclude colon information after the function name

Tmp_ebp[1](eip)-info->eip_fn_addr outputs the byte difference between the current instruction and the first instruction address of the function.

Source of `__STAB__` * in `debuginfo_eip`:

```
extern const struct Stab __STAB_BEGIN__[];    // Beginning of stabs table
extern const struct Stab __STAB_END__[];      // End of stabs table
extern const char __STABSTR_BEGIN__[];        // Beginning of string table
extern const char __STABSTR_END__[];          // End of string table
```

In `kdebug.c`, these variables are declared as `extern`, and their definitions are in `kernel.ld` as follows:

```
.stab : {
    PROVIDE(__STAB_BEGIN__ = .);
    *(.stab);
    PROVIDE(__STAB_END__ = .);
    BYTE(0) /* Force the linker to allocate space
             for this section */
}

.stabstr : {
    PROVIDE(__STABSTR_BEGIN__ = .);
    *(.stabstr);
    PROVIDE(__STABSTR_END__ = .);
    BYTE(0) /* Force the linker to allocate space
             for this section */
}
```

Taking the `.stab` section as an example, first define the `__STAB_BEGIN__` variable through `PROVIDE`, its address points to the beginning of the `stab` section, and then merge the `stab` sections of each file by the `*(.stab)` command (the section is compiled by the compiler when compiling each file). Generate, there are debugging information related to each file), and finally define the `__STAB_END__` variable, its position points to the end. The `.stabstr` section has the same definition as `.stab`.

The result is as follows:

```

QEMU
ebp f010ff98 eip f0100069 args 00000003 00000004 00000000 00000000 00000000
kern/init.c:15: test_backtrace+41
ebp f010ffb8 eip f0100069 args 00000004 00000005 00000000 00010094 00010094
kern/init.c:15: test_backtrace+41
ebp f010ffd8 eip f01000ea args 00000005 00001aac 00000644 00000000 00000000
kern/init.c:42: i386_init+77
ebp f010fff8 eip f010003e args 00111021 00000000 00000000 00000000 00000000
kern/entry.S:83: <unknown>+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> backtrace
Stack backtrace:
ebp f010ff68 eip f0100976 args 00000001 f010ff80 00000000 f010ffc8 f0112540
kern/monitor.c:156: monitor+251
ebp f010ffd8 eip f01000f6 args 00000000 00001aac 00000644 00000000 00000000
kern/init.c:42: i386_init+89
ebp f010fff8 eip f010003e args 00111021 00000000 00000000 00000000 00000000
kern/entry.S:83: <unknown>+0

```


Chapter2 Memory Management

2.1 Physical Page Management

2.1.1 Homework III

Question

In the file kern/pmap.c, you need to implement the code for the following function (see below, given in order):

Boot_alloc()

Mem_init() (before calling check_page_free_list(1))

Page_init()

Page_alloc()

Page_free()

Check_page_free_list() and check_page_alloc() will test your physical page allocator. You need to guide JOS

Then check the success report for check_page_alloc(). It would be helpful to add your own assert() to verify that your assumptions are correct.

Theoretical preparation

The operating system must keep track of which memory regions are free and which are occupied. The JOS kernel manages memory with the minimum granularity of pages (pages), which uses the MMU to map and protect the memory allocated for each block.

Here we have to write the allocation subfunction of the physical memory page. It uses a linked list of structure PageInfo to record which pages are free, and each node in the linked list corresponds to a physical page.

Analysis & Answer

The operating system must keep track of which physical RAM is free and which is in use. This exercise mainly writes a physical page allocator. It uses a linked list of PageInfo structures to record which pages are free, and each structure corresponds to a physical page.

Because the implementation of the page table requires the allocation of physical memory to store the page table, we need to write the physical page allocator before the implementation of the virtual memory.

```
struct PageInfo {  
    // pp_link 表示空闲内存列表中的下一页。注意，非空闲页的 pp_link 总是为 NULL。  
    struct PageInfo *pp_link;  
  
    //有多少个指针指向该页  
    uint16_t pp_ref;  
};
```

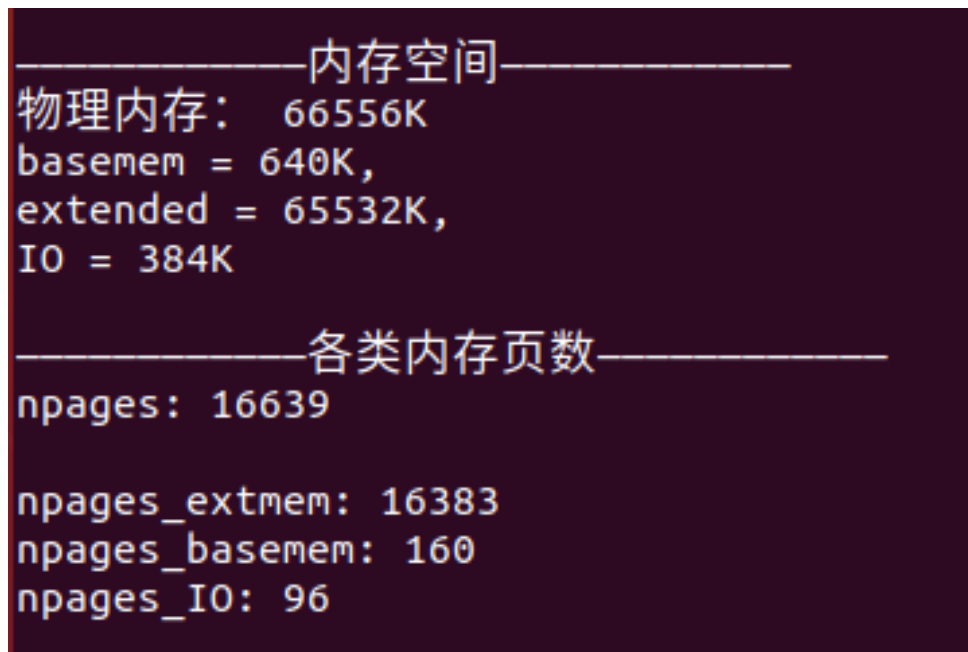
图 2.1: Definition of PageInfo in memlayout.h

In the file kern/pmap.c, the function `mem_init()` is called when the kernel first starts running, and some initialization settings are made for the memory management system of the entire operating system.

[Step 1] Enter `mem_init()`. The first step is to call the sub-function `i386_detect_memory()` to check how much total memory space is available in the system and how much the three parts of the memory space are, and divide the page according to the set `PGSIZE`. The total number of pages and the number of pages in each of the three parts are calculated in the function.

JOS divides the entire physical memory space into three parts:

- 1 Basemen: Available from 0x00000 to 0xA0000
- 2 IO hole: From 0xA0000 to 0x100000, not available, mainly used to allocate to external devices.
- 3 Extmen: From 0x100000 to 0x, available, the most important memory area.



```
———内存空间———  
物理内存: 66556K  
basemem = 640K,  
extended = 65532K,  
IO = 384K  
  
———各类内存页数———  
npages: 16639  
  
npages_extmem: 16383  
npages_basemem: 160  
npages_IO: 96
```

图 2.2: `i386_detect_memory()` run result

[Step 2] (1) Call the function `boot_alloc` function to allocate a virtual memory of size `PGSIZE` followed by the operating system kernel bss for storing the operating system page table. When the operating system is working in virtual memory mode, this page directory table is required for address translation (virtual address→physical address). The `boot_alloc` function can be used to allocate virtual memory, and its parameter is the virtual memory byte size to be allocated. When 0 is input, the currently used virtual memory tail can be queried. The core idea of this function is to maintain a static variable `nextfree`, which stores the virtual address corresponding to the next free memory space that can be used. The first time you enter the function you need to initialize `nextfree`. Note: The allocated memory size needs to be aligned with `PGSIZE` by the `ROUNDUP` function.

```
// 分配一个足够大的可以保存n字节的区域
// 更新nextfree保存下一个可用地址，现在存着，下次调用该函数时使用
// 确保nextfree保持与PGSIZE(4096)倍数对齐
result = nextfree;
nextfree = ROUNDUP(nextfree+n, PGSIZE);

//ROUNDUP(number,num_digits)|
//number 为需要向上舍入的任意实数。
//num_digits 四舍五入后的数字的位数。

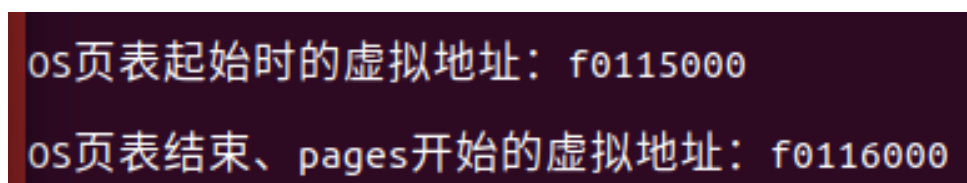
//超出内存，报错
if((uint32_t)nextfree - KERNBASE > (npages*PGSIZE))
    panic("Out of memory!\n");

return result;
```

图 2.3: The code supplemented by the `boot_alloc()` function

(2) Initialize a pointer in physical memory `kern_pgdir` points to the above operating system page table.

(3) Call the `memset` function to clear the virtual memory of the operating system page table area.



```
os页表起始时的虚拟地址: f0115000
os页表结束、pages开始的虚拟地址: f0116000
```

图 2.4: By information printing, the OS page table end address - OS page table start address = `PGSIZE`

[Step 3] Add the first page directory entry for the page directory table just created, occupying the 0th position of the page table. The content of the entry is the real address of the pointer `kern_pgdir`. UVPT is defined as the starting address of a virtual address, 0xef400000. Starting from this virtual address, the page table `kern_pgdir` of this operating system is stored, so we must map it to the physical address of the page table `kern_pgdir`, `PADDR(kern_pgdir)`. It is calculating the real physical address corresponding to `kern_pgdir`.

```
os的页表的第0项(指针kern_pgdir)所在的物理内存地址: 115005
```

图 2.5: The actual physical address of the pointer is known by information printing.

[Step 4] Use the `boot_alloc` function to allocate a virtual memory of size `npages * sizeof(struct PageInfo)`, which is used to store an array of `struct PageInfo`. The pointers in memory point to the virtual memory. The OS uses this array to track the usage of all memory pages. Each `PageInfo` represents a page in memory. `PageInfo` has two attributes: 1. Whether the current page is occupied. 2. Pointer to the next page. Then call the `memset` function to clear the corresponding area memory.

```
os页表结束、pages开始的虚拟地址: f0116000
os的页表的第0项(指针kern_pgdir)所在的物理内存地址: 115005
pages页表结束的虚拟地址: f0137000
```

图 2.6: Allocate virtual memory for pages

Note: It is found through information printing that the memory address at the end of the pages page table minus the memory address at the beginning of the pages is not exactly equal to `npages * sizeof(struct PageInfo)`, but slightly larger than when the virtual memory is allocated in the `boot_alloc` function. It needs to be noted that it is aligned with the `PGSIZE` multiple, so the allocated virtual memory size is appropriately allocated upwards.

[Step 5] Call the function `page_init()` to initialize the linked list of two `struct pageinfo` pages and `page_free_list`. The `pages_free_list` linked list stores all the free page information, and the pages store the information of the free non-free pages.

Jos divides the entire physical memory space into three parts, namely `BASEMEM`, `IO HOLE`, and `EXTMEM`.

According to the first item of the known basemem is already occupied, and other pages of basemem are not occupied. Since IO HOLE is mainly used to be allocated to external devices and is not available, it is equivalent to occupying all pages. EXTMEM is occupied in the first part of the page in step four, and part of it is idle.

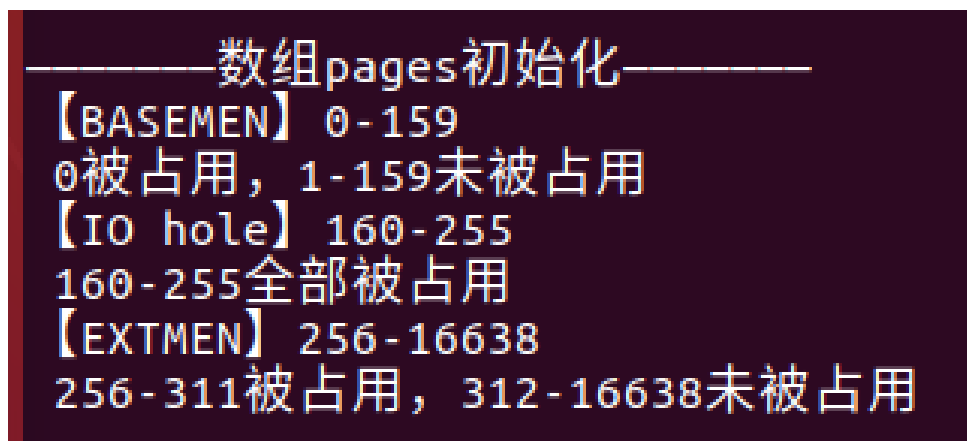


图 2.7: Pages initialization result

```
size_t i;
printf("\n\n\n——数组pages初始化——\n");
//已知pages[0]已被占用
pages[0].pp_ref = 1;
printf("【BASEMEM】 0-159\n 0被占用, ");
//basemem全都可用且未被占用
for (i = 1; i < npages_basemem; i++) {
    pages[i].pp_ref = 0;
    pages[i].pp_link = page_free_list;
    page_free_list = &pages[i];
}
printf("1-%d未被占用\n", npages_basemem-1);
//IO全都不可用
for (i = IOPHYSMEM/PGSIZE; i < EXTPHYSMEM/PGSIZE; i++) {
    pages[i].pp_ref = 1;
}
printf("【IO hole】 160-255\n %d-%d全部被占用\n", IOPHYSMEM/PGSIZE, EXTPHYSMEM/PGSIZE-1);
//extmen可用, 且部分已被占用
// (1) 已被占有部分
size_t first_free_address = PADDR(boot_alloc(0));
for (i = EXTPHYSMEM/PGSIZE; i < first_free_address/PGSIZE; i++) {
    pages[i].pp_ref = 1;
}
printf("【EXTMEM】 256-16638\n %d-%d被占用, %d-%d未被占用\n\n\n", EXTPHYSMEM/PGSIZE, first_free_address/PGSIZE-1, first_free_address/PGSIZE, npages-1);
// (2) 未被占用部分
for (i = first_free_address/PGSIZE; i < npages; i++) {
    pages[i].pp_ref = 0;
    pages[i].pp_link = page_free_list;
    page_free_list = &pages[i];
}
}
```

图 2.8: page_init() supplemental code screenshot

[Step 6] Call check_page_free_list() and check_page_alloc() to test the above physical page allocator.

Check_page_free_list() checks whether the so-called free pages of the page_free_list linked list are really legal and free. When the input parameter is 1, this function needs to perform one additional operation before checking, and modify the free page list free_page_list. After page_init, free_page_list has stored all the free page tables, but their order is according to the page table. The numbers are arranged from large to small. In the page directory table entry_pgdir (not kern_pgdir) used by the current operating system, the page table of the large number is not mapped, so we cannot operate this part of the page table. However, the small numbered page table, that is, from the page table 0 to the page table 1023, has been mapped, so this page table can be operated. Then check_page_free_list(1) is to complete the PageInfo structure corresponding to this part of the page table to the front end of the free_page_list for the operating system to use now. The rest of the operation is to check the free_page_list.

The function of the check_page_alloc() function is to check if page_alloc() and page_free() functions correctly.

One of the job requirements is to implement the page_alloc() function. By commenting we can know that the function of this function is to allocate a physical page. The return value of the function is the PageInfo structure corresponding to this physical page. So the general steps of this function should be:

1. Take a PageInfo structure of a free page from free_page_list
2. Modify the free_page_list related information, such as modifying the linked list header
3. Modify the PageInfo structure information of the taken free page to initialize the memory of the page.

```

//功能: 分配一个物理页, 返回值是这个物理页所对应的PageInfo结构体
struct PageInfo *
page_alloc(int alloc_flags)
{

    //    1. 从free_page_list中取出一个空闲页的PageInfo结构体
    //    2. 修改free_page_list相关信息, 比如修改链表表头
    //    3. 修改取出的空闲页的PageInfo结构体信息, 初始化该页的内存

    struct PageInfo *result;

    if (page_free_list == NULL)
        return NULL;

    //取出free_page_list的表头, 存在result里, 之后free_page_list失去该表头
    result= page_free_list;
    page_free_list = result->pp_link;
    result->pp_link = NULL;

    if (alloc_flags & ALLOC_ZERO)
        memset(page2kva(result), 0, PGSIZE);
    //修改取出的空闲页的PageInfo结构体信息, 初始化该页的内存

    return result;
}
    
```

图 2.9: Page_alloc() function complement code

Similarly, the job requires page_free(). According to the comment, the function of this method is to return the PageInfo structure of a page to the page_free_list free page list, which means that the page is recycled. Mainly complete the following operations:

1. Modify the corresponding information of the PageInfo structure of the page being recycled.
2. Insert the structure back into the page_free_list free page list

```

// 把一个页的PageInfo结构体再返回给page_free_list空闲页链表, 代表回收了这个页
void
page_free(struct PageInfo *pp)
{
    assert(pp->pp_ref == 0);
    assert(pp->pp_link == NULL);

    pp->pp_link = page_free_list;
    page_free_list = pp;
}

//
// Decrement the reference count on a page,
// freeing it if there are no more refs.
//
void
page_decref(struct PageInfo* pp)
{
    if (--pp->pp_ref == 0)
        page_free(pp);
}
    
```

图 2.10: Page_free function supplement code

2.2 Virtual Memory

2.2.1 Exercise IV

Question

Read Chapters 5 and 6 of the Intel 80386 Reference Manual for reading about page conversion and page-based protection(5.2 and 6.4). We recommend that you briefly understand the chapter on segmentation mode, because although the use of JOS is based on Page memory for virtual memory and protection, but segment conversion and segment-based protection cannot be disabled on x86, so you need Have a basic understanding..

Answer

The logical address of the page consists of the page number and the in-page address

The physical address of the page is spliced by the block number and the address within the page.

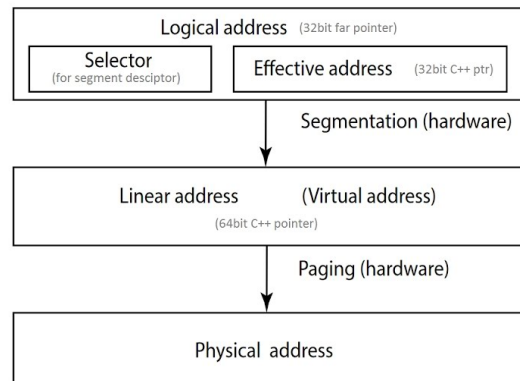
The purpose of the page table is to implement address mapping from page number to physical block number. The page table is retrieved by the page number of the logical address to obtain the physical block number of the page; and the in-page address d is directly sent into the block address field of the physical address register. In this way, the physical block number and the intra-block address are spliced into an address that actually accesses the memory, thereby completing the conversion from the logical address to the physical address.

So the calculation formula for the physical address is:

Physical address = block size (ie page size L)' block number f + page address d

2.2.2 Virtual address, linear address and physical address

In the x86 architecture, a virtual address is composed of two parts, one is a segment selector and the other is a segment offset. A Linear Address refers to an address obtained by converting a virtual address by a segment address translation mechanism. A physical address (Physical Addresses) is the real memory address obtained by the paging address translation mechanism after converting the linear address. This address will eventually be sent to the address bus of your memory chip. The specific relationship of the three addresses is as follows



2.2.3 Exercise V

Question

GDB can only access QEMU's memory through virtual addresses, but when learning to create virtual addresses, we also need to check the physical address at the same time. Learn QEMU's monitor command, especially the `xp` command, which allows you to check physical memory.

Answer

Open Terminal, enter `qemu-system-i386 -hda obj/kern/kernel.img -monitor stdio -gdb tcp::26000 -D qemu.log`, after the correct installation, we can use the monitor command.

```

QEMU
SeaBIOS (version Ubuntu-1.8.2-1ubuntu1)

iPXE (http://ipxe.org) 00:03.0 C980 PCI2.10 PnP PMM+07F92460+07ED2460 C980

Booting from Hard Disk...
6828 decimal is 15254 octal!
Physical memory: 66556K available, base = 640K, extended = 65532K
kernel panic at kern/pmap.c:124: mem_init: This function is not finished

Welcome to the JOS kernel monitor!
Type 'help' for a list of commands.
x 20,x 24
K> _

lyh@ubuntu: ~/Downloads/lab1/src/lab1_2
ES =0010 00000000 ffffffff 00cf9300 DPL=0 DS [-WA]
CS =0008 00000000 ffffffff 00cf9a00 DPL=0 CS32 [-R-]
SS =0010 00000000 ffffffff 00cf9300 DPL=0 DS [-WA]
DS =0010 00000000 ffffffff 00cf9300 DPL=0 DS [-WA]
FS =0010 00000000 ffffffff 00cf9300 DPL=0 DS [-WA]
GS =0010 00000000 ffffffff 00cf9300 DPL=0 DS [-WA]
LDT=0000 00000000 0000ffff 00008200 DPL=0 LDT
TR =0000 00000000 0000ffff 00008b00 DPL=0 TSS32-busy
GDT= 00007c4c 00000017
IDT= 00000000 000003ff
CR0=80010011 CR2=00000000 CR3=00110000 CR4=00000000
DR0=00000000 DR1=00000000 DR2=00000000 DR3=00000000
DR6=ffff0ff0 DR7=00000400
EFER=0000000000000000
FCW=037f FSW=0000 [ST=0] FTW=00 MXCSR=00001f80
FPR0=0000000000000000 0000 FPR1=0000000000000000 0000
FPR2=0000000000000000 0000 FPR3=0000000000000000 0000
FPR4=0000000000000000 0000 FPR5=0000000000000000 0000
FPR6=0000000000000000 0000 FPR7=0000000000000000 0000
XMM00=00000000000000000000000000000000 XMM01=00000000000000000000000000000000
XMM02=00000000000000000000000000000000 XMM03=00000000000000000000000000000000
XMM04=00000000000000000000000000000000 XMM05=00000000000000000000000000000000
XMM06=00000000000000000000000000000000 XMM07=00000000000000000000000000000000
    
```

图 2.11: qemu command

2.2.4 Question III

Question

Assuming the following kernel code is correct, then what type of variable x will be, `uintptr_t` or `Physaddr_t`?

```
Mystery_t x;
```

```
Char* value = return_a_pointer();
```

```
*value = 10;
```

```
x=(mystery_t) value;
```

Answer

Since the `*` operator is used here to resolve the address, the variable `x` should be of type `uintptr_t`.

2.2.5 Reference counting

In the previous experiment, we encountered `pp_ref`, which records how many different virtual addresses exist on each physical page to reference it. When this value becomes 0, the physical page can be released. In general, the `pp_ref` value of any physical page `p` is equal to the number of times it is mapped by the virtual page under the virtual address `UTOP` in all page table entries (the address range above `UTOP` is already mapped at startup) Finished, will not be changed afterwards).

2.2.6 Page Table Management

Now we can start writing programs that manage page tables: including inserting and deleting linear address-to-physical address mappings, and creating page tables.

2.2.7 Homework IV

Question

In the file `kern/pmap.c`, you must implement code for the following functions.

`pgdir_walk()`

`boot_map_region()`

`page_lookup()`

`page_remove()`

`page_insert()`

`check_page()`, called from `mem_init()`, tests your page table management routines. You should make sure it reports success before proceeding.

Theoretical preparation

In the previous description, we have a basic understanding of the concept of virtual address, linear address, and physical address.

Physical address

The unit addressing for the memory chip level corresponds to the address bus to which the processor and the CPU are connected. - This concept should be the best understood of these concepts. However, it is worth mentioning that although the physical address can be directly understood as the memory itself inserted in the machine, the memory is treated as a large array of numbers from 0 bytes up to the maximum amount

of bytes, and then this is An array is called a physical address. However, in fact, this is just an abstraction provided by the hardware to the software. The way memory is addressed is not the case. So, to say that it is "corresponding to the address bus" is more appropriate. Just aside from the consideration of the physical memory addressing mode, it is acceptable to directly compare the physical address with the physical memory. Perhaps the wrong understanding is more conducive to metaphysical abstraction.

Virtual memory

This is a description of the abstraction of the entire memory (not to the top of the machine).

It is relative to physical memory and can be directly understood as "not straight" and "fake" memory. For example, a 0x08000000 memory address. It is not correct for the address element of 0x08000000 -1 in the large array on the physical address.

The reason is this. Because modern operating systems provide a kind of memory management, that is, virtual memory. The process uses the address in virtual memory, which is assisted by the operating system to "convert" it into a real physical address.

This "conversion". It is the key to all issues discussed. With this kind of abstraction. A program can use a much larger address space than the real physical address.

Even multiple processes can use the same address. Not surprisingly. Since the converted physical address is not the same. —Can decompile the connected program and find that the connector has assigned an address to the program. For example, to call a function A, the code is not call A, but call 0x08111111, that is, function A The address has been fixed. Without such a "conversion", there is no concept of a virtual address, and doing so is simply not feasible.

Logical address

In order to be compatible, Intel has retained the segmental memory management method of ancient times.

A logical address is an address used in a machine language instruction to specify an operand or an instruction. In the above example, we say that the address assigned to A by the connector 0x08111111 is the logical address. . "A logical address is an offset of a segment identifier plus a relative address within a specified segment. It is expressed as [segment identifier: offset within segment], that is, the 0x08111111 in the above example should indicate For [A's code segment identifier: 0x08111111], this is complete."

Linear address or virtual address (virtual address)

Similar to the logical address, it is also an unreal address. If the logical address is the corresponding hardware platform segment management pre-conversion address, then the linear address corresponds to the pre-conversion address of the hardware page memory.

We also need to know the knowledge about the page table.

Our physical memory has a total of 4GB, and we assign it to the page, each page is 4KB in size.

The 4GB (2 to the 32th power) linear address space can be divided into 1048576 (2 to the 20th power, that is, 1M, can also be regarded as $1024 * 1024$) pages, so you can randomly extract these pages, every 1024 The pages are a group and can be divided into 1024 groups. For each group of 1024 pages of physical addresses, arranged in a certain order can constitute a table (each entry is the physical address of a page), this table is the page table. The size of the page table is $1024 * 4B = 4KB$, which is exactly the size of a physical page.

Because it has been divided into 1024 groups, each group has a page table (size 4KB), so these 1024 page tables can be pointed to by a table, this is the page directory. Similar to the page table, the page directory has a total of 1024 entries (called page directory entries), and the content of each page directory entry is the physical address of a page table. The size of the page table is $1024 * 4B = 4KB$, which is exactly the size of a physical page.

The conversion of three addresses and the page table mechanism are shown in the figure below.

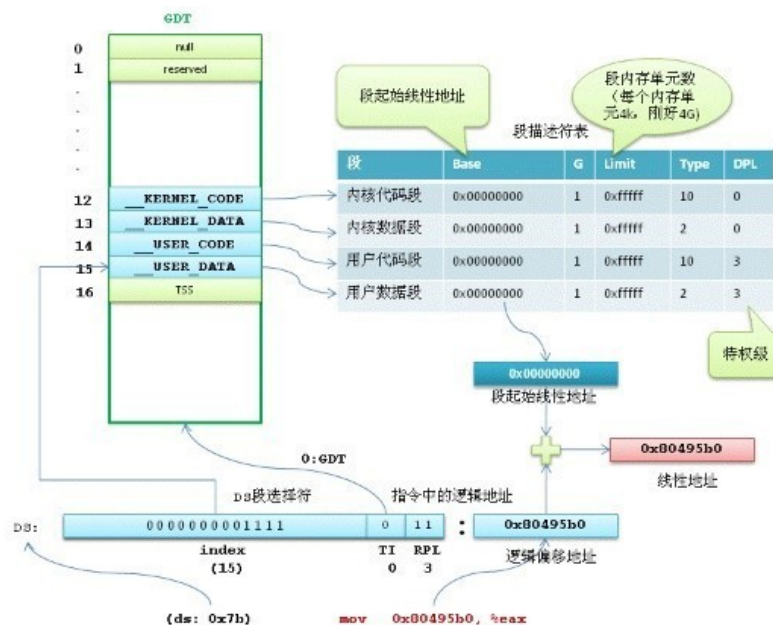


图1 逻辑地址转线性地址

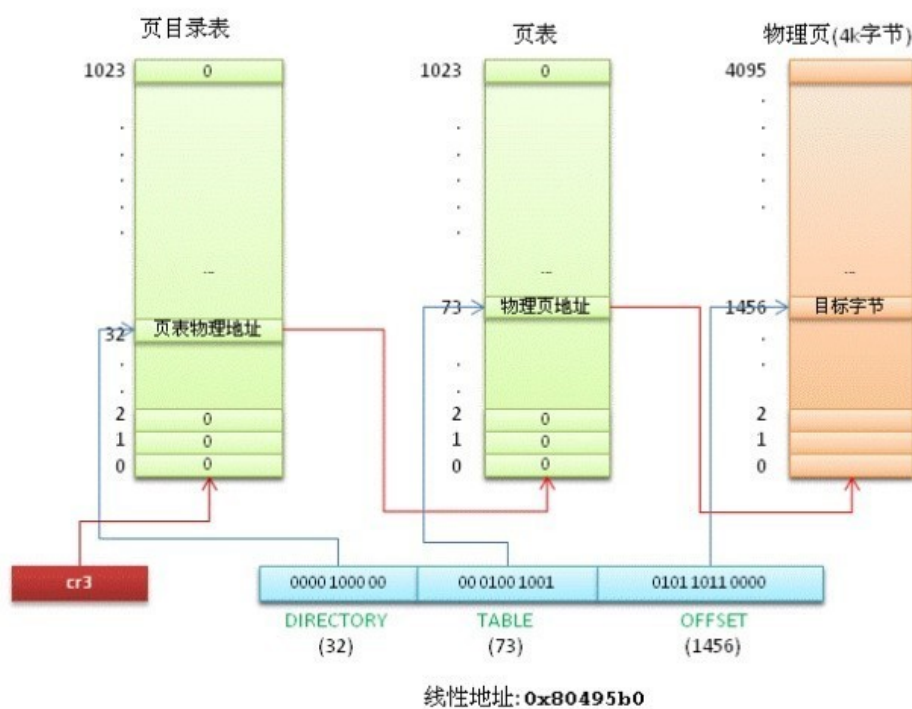
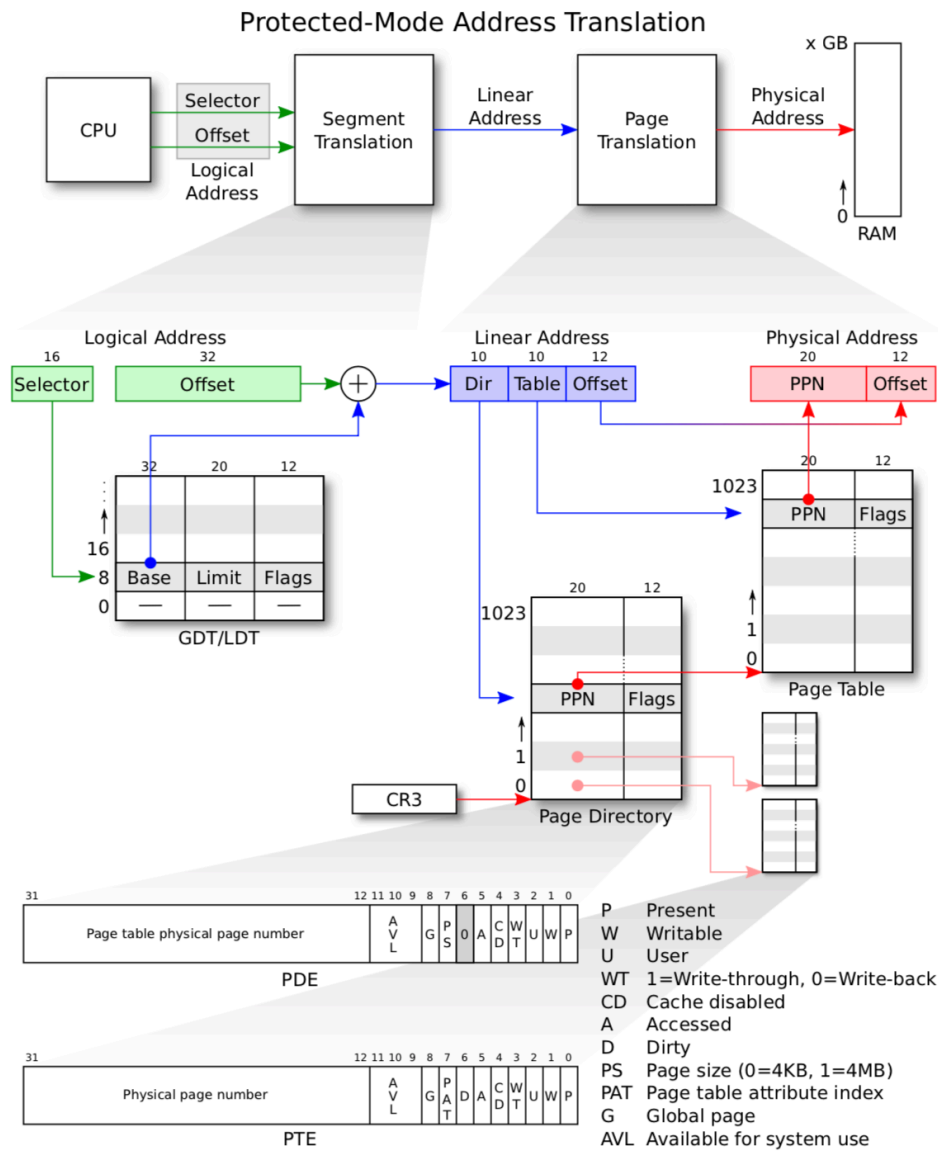
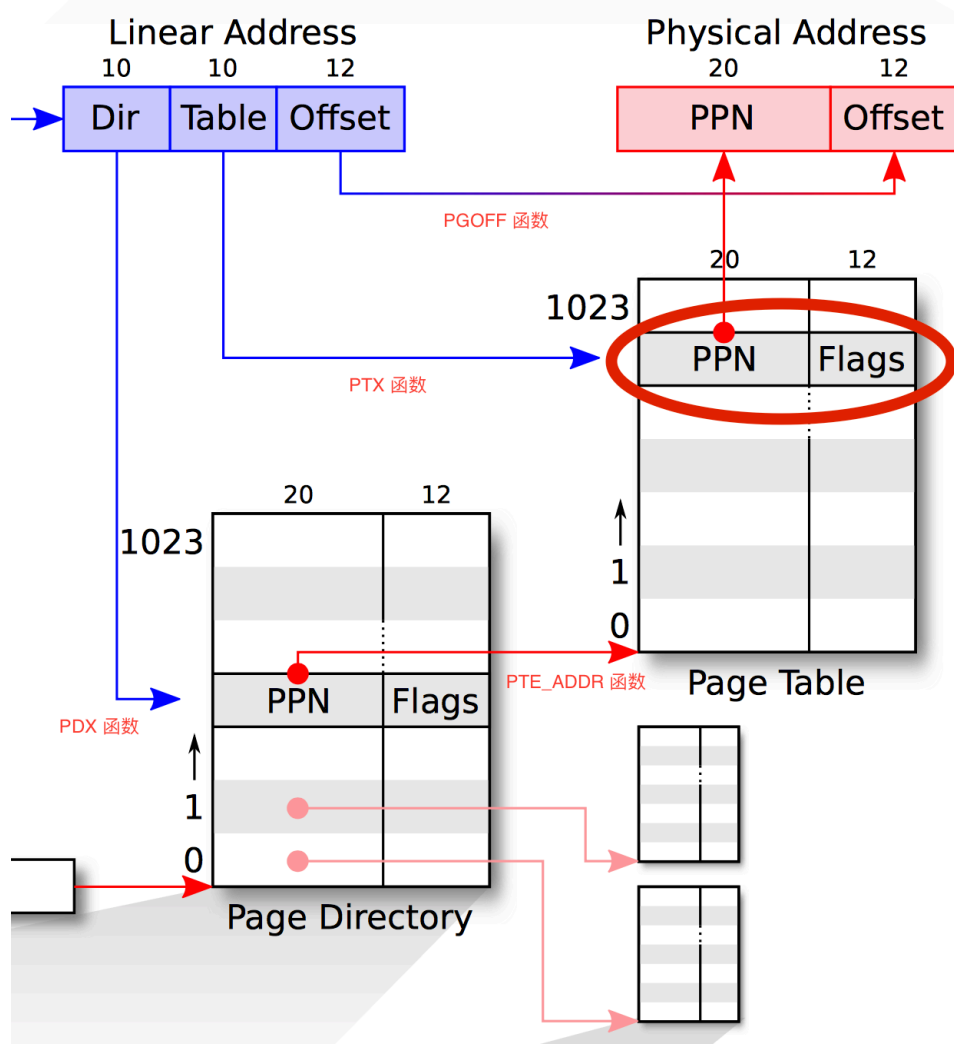


图2 线性地址转物理地址



Analysis & Answer

By commenting, we can know that the function of this function is given a page directory table pointer `pgdir`, which should return the page table entry pointer corresponding to the linear address `va`. We need to complete the transformation as shown in the figure, return the corresponding page table address, that is, the virtual address of the part circled by the red circle:



In addition, we need to understand the meaning of the three parameters. In addition, we need to understand the meaning of the three parameters, pgdir means the page directory item pointer, va means linear address, JOS equals virtual address, create means if page directory entry does not exist or not.

Here we should complete the following steps:

1. Find the page table page where the virtual address is located by the page directory table for the page directory entry address dic_entry_ptr in the page directory. (7-8)
2. Determine if the page table page corresponding to this page directory entry is already in memory. (10)
3. If yes, calculate the base address page_base of this page table page, and then return the address of the page entry corresponding to va &page_base[page_off] (23-25)
4. If not, and create is true, a new page is allocated, and the information of this page is added to the page directory entry dic_entry_ptr. (11-18)
5. If create is false, it returns NULL. (19-20)

The modified code is shown in the figure

```

uint32_t page_dir_idx=PDX(va);
uint32_t page_tab_idx=PTX(va);
pte_t *pgtab;
if(pgdir[page_dir_idx]==(pde_t)NULL)
{
    if(create==0)
        return NULL;
    else{
        struct PageInfo *new_pageinfo=page_alloc(ALLOC_ZERO, 1);
        if(new_pageinfo){
            new_pageinfo->pp_ref+=1;
            pgtab=(pte_t *)page2kva(new_pageinfo);
            pgdir[page_dir_idx]=PADDR(pgtab) | PTE_P | PTE_W;
            return &pgtab[page_tab_idx];
        }
        else
        {
            return NULL;
        }
    }
}
pgtab=KADDR(PTE_ADDR(pgdir[page_dir_idx]));
return &pgtab[page_tab_idx];
-}
    
```

图 2.12: pgdir_walk

Next we complete the `boot_map_region` function. By comment, we can see that the function of this function is to map $[va, va+size)$ of virtual address space to physical $[pa, pa+size)$ in the page table rooted at `pgdir`. Size is a multiple of `PGSIZE`, and `va` and `pa` are both page-aligned. Use permission bits `perm|PTE_P` for the entries.

The mapping of the virtual address space range $[va, va+size)$ to the physical space $[pa, pa+size)$ is added to the page table `pgdir`. The main purpose of this function is to set the address range above the virtual address `UTOP`. The address mapping of this part is static and will not change during the operation of the operating system, so the value of the `pp_ref` field in the `PageInfo` structure of this page. No change will happen.

The steps to be completed by this function are as follows:

Need to complete a loop, using the `pgdir_walk` function we just completed, we can use a loop to map all the memory in size bytes. The modified code is shown in the figure:

```

409 boot_map_region(pde_t *pgdir, uintptr_t va, size_t size, physaddr_t pa, int perm)
410 {
411     // Fill this function in
412     pte_t *pgtab;
413     size_t pg_num=PGNUM(size);
414     printf("map region size = %d, %d pages\n",size,pg_num);
415     size_t i;
416     for( i=0;i<pg_num;i++){
417         pgtab=pgdir_walk(pgdir,(void *)va,1);
418         if(!pgtab)return;
419         *pgtab=pa|perm|PTE_P;
420         va+=PGSIZE;
421         pa+=PGSIZE;
422     }
423 }

```

图 2.13: boot_map_region

Next, we continue to view page_insert(). The function prototype is as follows: page_insert(pde_t *pgdir, struct PageInfo *pp, void *va, int perm), which is functionally completed: mapping a page pp in a physical memory to a virtual address va .

The main steps of this function are as follows:

1. First, find the page table entry corresponding to the virtual address va by using the pgdir_walk function. (4)
2. Modify the value of pp_ref. (8)
3. Check the page table entry to determine if va has been mapped. If it is mapped, delete the mapping. (9-13)
4. Add the mapping between va and pp to the page table entry. (14-15)

```

450 int
451 page_insert(pde_t *pgdir, struct PageInfo *pp, void *va, int perm)
452 {
453     // Fill this function in
454     pte_t *pgtab=pgdir_walk(pgdir,va,1);
455     if(!pgtab)return -E_NO_MEM;
456     pp->pp_ref++;
457     if(*pgtab&PTE_P){
458         page_remove(pgdir,va);
459     }
460     *pgtab=page2pa(pp)|perm|PTE_P;
461     return 0;
462 }
463

```

图 2.14: page_insert

Next, we complete the page_lookup function, which functions as return the page mapped at virtual address 'va'. If the pte_store parameter is not 0, the page table entry address of this physical page is stored in pte_store.

We only need to call the pgdir_walk function to get the page table entry corresponding to this va, and then determine whether the page is already in memory, and if so, return the PageInfo structure pointer of this page. And store the contents of this page table entry in pte_store. The modified code is shown in the figure:

```

struct PageInfo *
page_lookup(pde_t *pgdir, void *va, pte_t **pte_store)
{
    // Fill this function in
    pte_t *pgtab=pgdir_walk(pgdir,va,0);
    if(!pgtab) return NULL;
    if(pte_store) *pte_store=pgtab;
    return pa2page(PTE_ADDR(*pgtab));
}

```

图 2.15: page_lookup

The last one is the page_remove function. Its prototype is: void page_remove(pde_t *pgdir, void *va). The function is to delete the mapping between the virtual address va and the physical page.

The note also hints at a few details to be aware of:

1. The pp_ref value should be reduced by one.
2. If pp_ref is reduced to 0, this page should be recycled
3. The page table entry corresponding to this page should be set to 0.

```

500 void
501 page_remove(pde_t *pgdir, void *va)
502 {
503     // Fill this function in
504     pte_t *pte=pgdir_walk(pgdir,va,0);
505     pte_t **pte_store=&pte;
506     struct PageInfo *pp=page_lookup(pgdir,va,pte_store);
507     if(!pp) return;
508     page_decref(pp);
509     **pte_store=0;
510     tlb_invalidate(pgdir,va);
511 }
512

```

图 2.16: page_remove

2.3 Kernel address space

JOS divides the 32-bit linear address virtual space into two parts. The user environment (process running environment) usually occupies the part of the low address, called the user address space. The operating system kernel always occupies the part of the high address, called the kernel address space. The dividing line between these two parts is a macro ULIM defined in the memlayout.h file. JOS reserves nearly 256MB of virtual address space for the kernel. This can be understood, why in the experiment 1 to design a high address address space for the operating system. If you don't do this, the address space of the user environment is not enough.

2.3.1 Homework V

Question

Fill in the missing code in `mem_init()` after the call to `check_page()`.

Theoretical preparation

Since the kernel and user memory are present in the address space of each environment, we need to use the permission bits in the x86 page table to allow the user code to access only the user portion of the address space. Otherwise, bugs in the user code may overwrite the kernel data, causing a crash; or the user code can steal private data from other environments. The user environment will have no access to any memory above ULIM, and the kernel can read and write this portion of memory. For the address range [UTOP, ULIM), the kernel and the user environment have the same permissions: they can only be read and not written. Under UTOP The address space is used by the user environment, and the user environment will set permissions to access this part of the memory.

```

/*
** Virtual memory map: ..... Permissions
** ..... kernel/user
**
** 4 Gig -----> +-----+
** ..... | ..... | RW/--
** ..... ~~~~~
** ..... : ..... :
** ..... : ..... :
** ..... : ..... :
** ..... | ..... | RW/--
** ..... | ..... | RW/--
** ..... | Remapped Physical Memory | RW/--
** ..... | ..... | RW/--
** ..... KERNBASE, -----> +-----+ 0xf0000000 .....+
** ..... KSTACKTOP ..... CPU0's Kernel Stack ..... | RW/-- KSTKSIZE |
** ..... | ..... | ..... |
** ..... | Invalid Memory (*) ..... | --/-- KSTKGAP |
** ..... +-----+
** ..... | CPU1's Kernel Stack ..... | RW/-- KSTKSIZE |
** ..... | ..... | ..... | PTSIZE
** ..... | Invalid Memory (*) ..... | --/-- KSTKGAP |
** ..... +-----+
** ..... : ..... :
** ..... : ..... :
** ..... MMIO LIM -----> +-----+ 0xefc00000 .....+
** ..... | Memory-mapped I/O ..... | RW/-- PTSIZE
** ..... ULIM, MMIOBASE -----> +-----+ 0xef800000
** ..... | Cur. Page Table (User R-) ..... | R-/R- PTSIZE
** ..... UVPT -----> +-----+ 0xef400000
** ..... | RO PAGES ..... | R-/R- PTSIZE
** ..... UPAGES -----> +-----+ 0xef000000
** ..... | RO ENV$ ..... | R-/R- PTSIZE
** ..... UTOP, UENVS -----> +-----+ 0xsec00000
** ..... UXSTACKTOP -/ ..... | User Exception Stack ..... | RW/RW PGSIZE
** ..... +-----+ 0xeebfff00
** ..... | Empty Memory (*) ..... | --/-- PGSIZE
** ..... USTACKTOP -----> +-----+ 0xeebfe000
** ..... | Normal User Stack ..... | RW/RW PGSIZE
** ..... +-----+ 0xeebfd000
** ..... | ..... |
** ..... | ..... |
** ..... ~~~~~
** ..... : ..... :
** ..... : ..... :

```

Analysis & Answer

In this exercise, three virtual addresses are mapped to the physical page.

First, we will complete the UPAGES mapping UPAGES (0xef000000 0xef400000) up to 4MB, which is the data structure of JOS recording physical page usage.

Currently only one page directory is created, kernel_pgdir, so the first parameter is obviously kernel_pgdir. The second parameter is the virtual address, and UPAGES is originally given as a virtual address. The third parameter is the mapped memory block size. The fourth parameter is the physical address mapped to, and the physical address of pages can be taken directly. Permissions PTE_U indicates that the user has permission to read. Currently only one page directory is created, kernel_pgdir, so the first parameter is obviously kernel_pgdir. The second parameter is the virtual address, and UPAGES is originally given as a virtual address. The third parameter is the mapped memory block size. The fourth parameter is the physical address mapped to, and the physical address of

pages can be taken directly. Permissions PTE_U indicates that the user has permission to read. The modified code is shown in the figure

```

→ //////////////////////////////////////////////////
→ // Now we set up virtual memory
→ //////////////////////////////////////////////////
→ // Map 'pages' read-only by the user at linear address UPAGES
→ // Permissions:
→ //----- the new image at UPAGES -- kernel R, user R
→ //----- (ie. perm = PTE_U | PTE_P)
→ //----- pages itself -- kernel RW, user NONE
→ // Your code goes here:
→ cprintf("UPAGES\n");
→ boot_map_region(kern_pgdir, (uintptr_t) UPAGES, npages*sizeof(struct PageInfo), PADDR(pages), PTE_U | PTE_P);

```

Then there is the memory stack, the kernel stack (0xffff8000 0xf0000000) 32k-B Bootstrap represents the lowest address of the stack. Since the stack grows to the lower address, it is actually the top of the stack. We will map the address space in [KSTACKTOP-KSTKSIZE, KSTACKTOP) The modified code is shown in the figure

```

→ //////////////////////////////////////////////////
→ // Use the physical memory that 'bootstack' refers to as the kernel
→ // stack. The kernel stack grows down from virtual address KSTACKTOP.
→ // We consider the entire range from (KSTACKTOP-PTSIZE, KSTACKTOP)
→ // to be the kernel stack, but break this into two pieces:
→ //-----* (KSTACKTOP-KSTKSIZE, KSTACKTOP) -- backed by physical memory
→ //-----* (KSTACKTOP-PTSIZE, KSTACKTOP-KSTKSIZE) -- not backed; so if
→ //-----the kernel overflows its stack, it will fault rather than
→ //-----overwrite memory. Known as a "guard page".
→ //-----Permissions: kernel RW, user NONE
→ // Your code goes here:
→ cprintf("memory stack\n");
→ boot_map_region(kern_pgdir, (uintptr_t) (KSTACKTOP-KSTKSIZE), KSTKSIZE, PADDR(bootstack), PTE_W | PTE_P);

```

Finally, the kernel part, we will map the kernel (0xf0000000 0xffffffff) 256MB The modified code is shown in the figure

```

→ //////////////////////////////////////////////////
→ // Map all of physical memory at KERNBASE.
→ // Ie. the VA range [KERNBASE, 2^32) should map to
→ //----- the PA range [0, 2^32 - KERNBASE)
→ // We might not have 2^32 - KERNBASE bytes of physical memory, but
→ // we just set up the mapping anyway.
→ // Permissions: kernel RW, user NONE
→ // Your code goes here:
→ cprintf("kernel\n");
→ boot_map_region(kern_pgdir, (uintptr_t) KERNBASE, ROUNDUP(0xffffffff - KERNBASE, PGSIZE), 0, PTE_W | PTE_P);

```

As a result, the code works successfully.

```

———内存空间———
物理内存: 66556K
basemem = 640K,
extended = 65532K,
IO = 384K

———各类内存页数———
npages: 16639

npages_extmem: 16383
npages_basemem: 160
npages_IO: 96

os页表起始时的虚拟地址: f0118000
os页表结束、pages开始的虚拟地址: f0119000
os的页表的第0项(指针kern_pgdir)所在的物理内存地址: 118005
pages页表结束的虚拟地址: f013a000

———数组pages初始化———
【BASEMEM】 0-159
0被占用, 1-159未被占用
【IO hole】 160-255
160-255全部被占用
【EXTMEM】 256-16638
256-313被占用, 314-16638未被占用

check_page_alloc() succeeded!
check_page() succeeded!
UPAGES
map_region_size =133112,(32 pages)
memory stack
map_region_size =32768,(8 pages)
kernel
map_region_size =268435456,(65536 pages)
check_kern_pgdir() succeeded!
check_page_installed_pgdir() succeeded!
Welcome to the JOS kernel monitor!

```



```

nikoni@ubuntu:~/Documents/lab1/src/lab1_1$ make grade
make clean
make[1]: Entering directory '/home/nikoni/Documents/lab1/src/lab1_1'
rm -rf obj .gdbinit jos.in qemu.log
make[1]: Leaving directory '/home/nikoni/Documents/lab1/src/lab1_1'
./grade-lab2
make[1]: Entering directory '/home/nikoni/Documents/lab1/src/lab1_1'
+ as kern/entry.S
+ cc kern/entrypgdir.c
+ cc kern/init.c
+ cc kern/console.c
+ cc kern/monitor.c
+ cc kern/pmap.c
+ cc kern/kclock.c
+ cc kern/printf.c
+ cc kern/kdebug.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 390 bytes (max 510)
+ mk obj/kern/kernel.img
make[1]: Leaving directory '/home/nikoni/Documents/lab1/src/lab1_1'
running JOS: (1.1s)
  Physical page allocator: OK
  Page management: OK
  Kernel page directory: OK
  Page management 2: OK
Score: 70/70

```

2.3.2 Question IV

1)

What entries (rows) in the page directory have been filled in at this point? What addresses do they map and where do they point? In other words, fill out this table as much as possible:

Entry	Base Virtual Address	Points to (logically):
1023	?	Page table for top 4MB of phys memory
1022	?	?
.	?	?
.	?	?
.	?	?
2	0x00800000	?
1	0x00400000	?
0	0x00000000	[see next question]

Answer

Entry	Base Virtual Address	Points to(logically)
1023	0xffc00000 Points to(logically)	Page table for top 4MB of physical memory.This is the last address finding page table that the kernel can use.
1022	0xff800000	Page table for 248MB–(252MB-1)physical memory
...	...	Page table for physical memory
960	0xf0000000(KERNBASE)	static data 0–(4MB-1) physical memory
959	0xefc00000(VPT)	Page directory self (kernel RW).This is the first page table and it ' s in the bottom of the physical memory.
958	0xef800000(ULIM)	Page table for kernel stack.It is mapped into the physical memory which is the same as bootstack.We only map the memory that the same as KSTACKSIZE.The rest memory that wasn' t mapped is used to avoid the overflow of the kernel stack.
957	0xef400000(UVPT)	Same as 959(user kernel R)
956	0xef000000(UPAGES)	Page table for structure pages[]
...	...	NULL
2	0x00800000	NULL
1	0x00400000	NULL
0	0x00000000	The start of the virtual memory .The same as 960(then turn to NULL)

表 2.1: Answer

2) We have placed the kernel and user environment in the same address space. Why will user programs not be able to read or write the kernel's memory? What specific mechanisms protect the kernel memory?

Answer

User is not allowed to access kernel memory for safety reasons. If user have the permission, bugs in user code may led to crash. It is the paging mechanism that protects kernel address in JOS. If the flag bit PTE_U is 0 in a page, that means user have no permission to read or write the page.

3) What is the maximum amount of physical memory that this operating system can support? Why?

Answer

```
// Page directory and page table constants.
#define NPENTRIES 1024 // page directory entries per page directory
#define NPTENTRIES 1024 // page table entries per page table

#define PGSIZE 4096 // bytes mapped by a page
#define PGSHIFT 12 // log2(PGSIZE)

#define PTSIZE (PGSIZE*NPTENTRIES) // bytes mapped by a page directory entry
#define PTSHIFT 22 // log2(PTSIZE)

#define PTXSHIFT 12 // offset of PTX in a linear address
#define PDXSHIFT 22 // offset of PDX in a linear address
```

图 2.17: PTSIZE

```
// User read-only virtual page table (see 'uvpt' below)
#define UVPT (ULIM - PTSIZE)
// Read-only copies of the Page structures
#define UPAGES (UVPT - PTSIZE)
// Read-only copies of the global env structures
#define UENVS (UPAGES - PTSIZE)
..
```

图 2.18: UPAGES

```
struct PageInfo {
    // Next page on the free list.
    struct PageInfo *pp_link;

    // pp_ref is the count of pointers (usually in page table entries)
    // to this page, for pages allocated using page_alloc.
    // Pages allocated at boot time using pmap.c's
    // boot_alloc do not have valid reference count fields.

    uint16_t pp_ref;
};
```

图 2.19: Page

2GB.

Pages use up to 4MB space, and each PageInfo use 8Byte. $4M / 8 * 4kB = 2GB$ (4kB per page).

4)How much space overhead is there for managing memory, if we actually had the maximum amount of physical memory? How is this overhead broken down?

Answer

totally 6MB+4KB

The maximum physical memory is 2G, so there's 0.5M pages in total.

$8 * 0.5M = 4MB$ (struct PageInfo use 8bytes per page)

$4 * 0.5M = 2MB$ (each page table uses 4bytes)

4KB (Page director uses 4kB)

Chapter3 About Us

	Yinhao Li	Jiajun Li	Qiaoqiao Pan	Xinghe Li	Hongbo Xu
contribution	20%	20%	20%	20%	20%

表 3.1: Contribution

	BOOTING A PC	MEMORY MANAGE- MENT	Writing	Code sorting
Yinhao Li	✓	✓	✓	
Jiajun Li	✓	✓		✓
Qiaoqiao Pan	✓	✓		✓
Xinghe Li		✓		
Hongbo Xu		✓		

表 3.2: Group division