

# Zeos

Yolanda Becerra, Juan José Costa and Alex Pajuelo

Course 2016–2017



## PREFACE

The aim of this document is to serve as a basic documentation of an operating system called ZeOS, based on a Linux 2.4 kernel, and developed for the intel 80386 architecture. ZeOS was first developed by a group of students from the Barcelona School of Informatics (FIB), with the support of a number of professors from the Department of Computer Architecture (AC). Anyone is free to add more functionalities to this OS and to make further contributions.

The document will describe an initial basic implementation, to which a number of functionalities will be added. This basic implementation is responsible for booting up the OS and initialize all the necessary data structures it requires. Both high-level (C) and low-level (assembler) programming languages will be used in order to add more functionality to the system.

The first piece of work to be completed will be on an essential part of any OS: the boot process. This will be followed by a section on the management of some basic interrupt and exception handling. Work will then be undertaken on process management in the OS.

For this document to be understood, the following concepts must be borne in mind:

- Input/output mechanisms (exceptions/interrupts/system calls).
- Process management (data structures/algorithms/scheduling policies/context switch/related system calls).
- Input/output management (devices/file descriptors).
- Subroutines and exceptions.
- Memory management.

### What you should do with this document

You must first read all the documentation so that you have an overall vision of ZeOS. After that, it is advisable that you follow all the steps described in this document to design and implement some of the components of the OS such as new data structures, functions, algorithms, etc.

# CONTENTS

<b>1</b>	<b>Acknowledgements</b>	<b>6</b>
<b>2</b>	<b>INTRODUCTION TO ZEOS</b>	<b>6</b>
2.1	Contents . . . . .	6
2.2	Bootng . . . . .	8
2.3	Image construction . . . . .	8
2.4	Introductory session . . . . .	9
2.4.1	Getting started . . . . .	9
2.4.2	Working environment . . . . .	10
2.4.3	Prior knowledge . . . . .	11
2.4.4	Bochs . . . . .	12
2.4.5	Frequently used commands . . . . .	14
2.5	Work to do . . . . .	15
2.5.1	Where is ZeOS hanged? . . . . .	16
2.5.2	User code modification . . . . .	16
2.5.3	Use of inline assembler . . . . .	18
2.5.4	Additional questions . . . . .	19
<b>3</b>	<b>MECHANISMS TO ENTER THE SYSTEM</b>	<b>20</b>
3.1	Preliminary concepts . . . . .	20
3.2	Hardware management of an interrupt . . . . .	21
3.2.1	Task State Segment (TSS) . . . . .	21
3.3	Function name conventions . . . . .	22
3.4	Files . . . . .	23
3.5	Programming exceptions . . . . .	23
3.5.1	Exception parameters . . . . .	23
3.5.2	Initializing the IDT . . . . .	23
3.5.3	Writing the handler . . . . .	24
3.5.4	Writing the service routines . . . . .	25
3.6	Programming interrupts . . . . .	25
3.6.1	The keyboard interrupt . . . . .	25
3.6.2	Writing the handler . . . . .	26
3.6.3	Writing the service routine . . . . .	27
3.7	Programming system calls . . . . .	28

3.7.1	Independence from devices . . . . .	28
3.7.2	Returning results . . . . .	29
3.7.3	Writing the write wrapper . . . . .	29
3.7.4	IDT initialization . . . . .	29
3.7.5	Writing the handler . . . . .	30
3.7.6	Service Routine to the write system call . . . . .	31
3.7.7	Copying data from/to the user address space . . . . .	32
3.8	Work to do . . . . .	32
3.8.1	Complete Zeos Code . . . . .	32
3.8.2	Clock management . . . . .	33
3.8.3	Gettime system call . . . . .	34
<b>4</b>	<b>BASIC PROCESS MANAGEMENT</b>	<b>35</b>
4.1	Prior concepts . . . . .	35
4.2	Definition of data structures . . . . .	35
4.2.1	Process data structures . . . . .	35
4.2.2	Process identification . . . . .	38
4.3	Memory management . . . . .	38
4.3.1	MMU: Paging mechanism . . . . .	39
4.3.2	The directory and page table . . . . .	39
4.3.3	Translation Lookahead Buffer (TLB) . . . . .	41
4.3.4	Specific organization for ZeOS . . . . .	42
4.4	Initial processes initialization . . . . .	45
4.4.1	Idle process . . . . .	45
4.4.2	Init process . . . . .	46
4.4.3	Zeos implementation details . . . . .	47
4.5	Process switch . . . . .	47
4.5.1	Intial testing . . . . .	49
4.6	Process creation and identification . . . . .	49
4.6.1	<i>getpid</i> system call implementation . . . . .	49
4.6.2	<i>fork</i> system call implementation . . . . .	49
4.7	Process scheduling . . . . .	51
4.7.1	Main scheduling functions . . . . .	51
4.7.2	Round robin policy implementation . . . . .	52
4.8	Process destruction . . . . .	52
4.9	Statistical information of processes . . . . .	53
4.10	Work to do . . . . .	54

<b>5</b>	<b>LIGHTWEIGHT THREADS AND PROCESS SYNCHRONIZATION</b>	<b>54</b>
5.1	Work to do . . . . .	54
5.1.1	Lightweight processes: threads . . . . .	55
5.1.2	Synchronization between processes: semaphores . . . . .	55
<b>6</b>	<b>I/O MANAGEMENT AND DYNAMIC MEMORY</b>	<b>57</b>
6.1	MANAGEMENT OPERATIONS OF THE KEYBOARD DEVICE . . . . .	57
6.1.1	New data structures . . . . .	57
6.1.2	Keyboard read implementation . . . . .	58
6.1.3	Keyboard service routine . . . . .	58
6.2	DYNAMIC MEMORY . . . . .	59

## 1 ACKNOWLEDGEMENTS

This document was drawn up with the support of professors on previous courses: Julita Corbalán, Marisa Gil, Jordi Guitart, Gemma Reig, Amador Millán, Jordi García, Silvia LLorente, Pablo Chacín and Rubén González. The authors wish to thank A. Bartra, M. Muntanyá and O. Nieto for their contributions.

This document has been improved by the following people:

- Albert Batalle Garcia (2011)

## 2 INTRODUCTION TO ZEOS

The construction of an OS is similar to the construction of an ordinary executable. This document will show you how an OS is built from the source code. The construction is very similar to the Linux OS building process. With minor changes, this documentation may in fact be useful in explaining how to build a Linux OS.

After downloading and uncompressing the ZeOS source code, the fastest way to build your ZeOS is to type *make* in the directory with the files. The *Makefile* will guide you through the process, compiling all the source files and linking them together to build a final bootable image (a file called *zeos.bin*). This process is explained in detail in the following subsections.

### 2.1 Contents

The source code of your ZeOS contains files and directories. The content of the files can be divided into the following groups, depending on their extension:

- 1) Source files written in C language, whose extension is *.c*
- 2) Source files written in Intel 80386 assembler language with preprocessor sentences, with an *.S* (capital S) extension. The *.c* and *.S* files are the only ones that add code to the OS.
- 3) Scripts used by the *ld* linker to combine the various files in a single binary file, with an *.lds* extension.
- 4) Header files (extension *.h*) located in a dedicated directory, as occurs in Linux.
- 5) The *Makefile* that will guide you through the steps to build the OS. It is the only file without an extension.

Usually, there is a header file for each source file in C language, which includes all variables, functions, macros and type declarations used in this C file in case you wish to use them from another C file.

After the *make* process you should obtain a single file called *zeos.bin* with the bootable image of your ZeOS (which could be copied to a floppy disk to test that it effectively boots up in your computer).

All files that add code to ZeOS (*.c* and *.S* files) can be divided into several groups, depending on the final position of each file in the image of the OS. Figure 1 shows the final snapshot that the image of the OS should have once it is built<sup>1</sup>. It shows three main blocks:

- The **boot sector block**. This block only comprises one file: *bootsect.S*. It must weigh exactly 512 bytes, since it has to be fitted in sector 0 of a conventional floppy disk.

1. Some of the file names that appear in Figure 1 may not match the ones from your currently downloaded source code

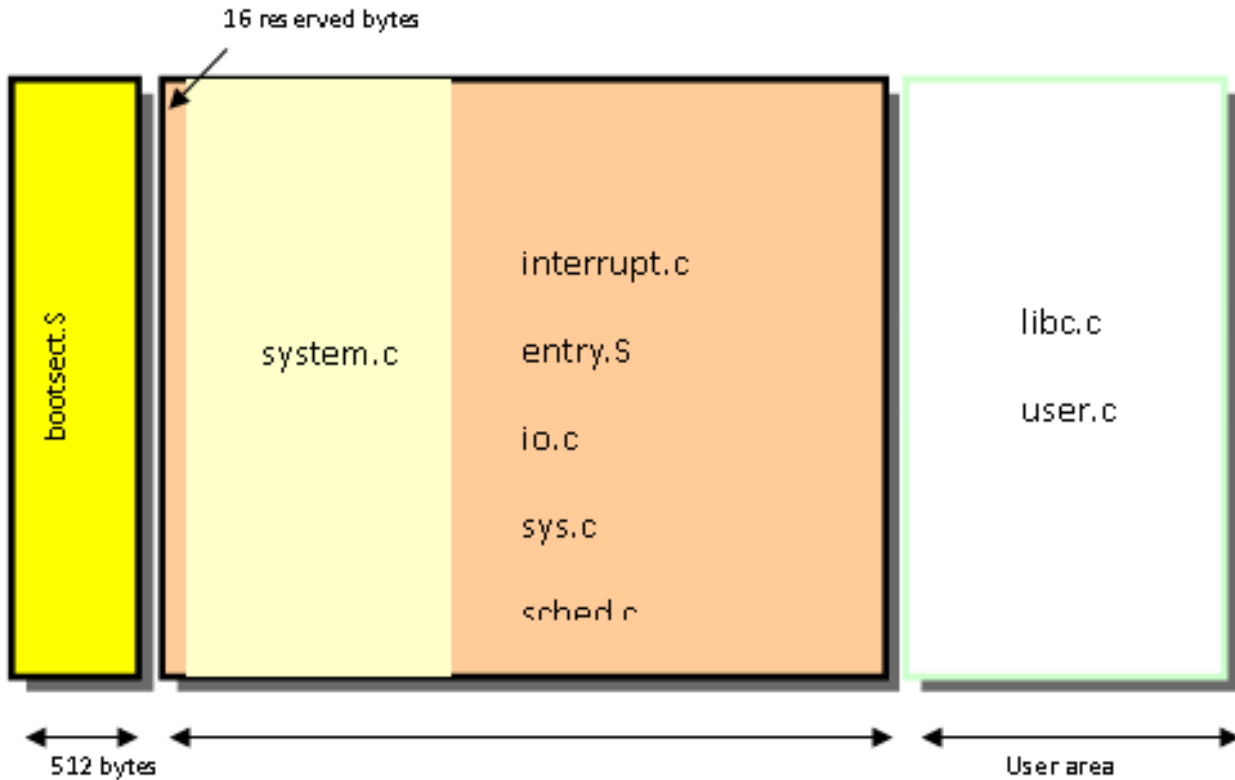


Fig. 1. Snapshot of Zeos

- The **system area block**. This block is the largest and contains the main files of the OS. This block is placed in a part of the memory that guarantees it will be executed with the processor's maximum privilege level.
- The **user area block**. This block has a user program that can make use of the operating system services (use system calls). The code inside is executed by the processor with a minimum privilege level.

One may ask why the user program is attached to the OS when this actually never happens. User programs are usually found on a hard disk inside a file system in a directory of this file system. The simple answer is that currently there is not either file system or an executable loader. Therefore, the user program is attached to the OS (in a position that is easy to calculate), and it is copied from its location in the image to a memory area with user privileges to which execution is transferred.

Figure 2 shows the contents of the system.lds file that enables the linker to locate on the memory the various sections of an executable file. In this case, it sets the starting address (0x10000) of the system code; reserves 24 bytes, using the BYTE directive, to store the user code size and other necessary data for managing the user code; then the system code (.text); the read-only data section (.rodata); the initialized data (.data); the not initialized data (.bss); and, finally, it leaves a gap to align the current address to a page-size (4096 bytes) address and adds a special section (.data.task) which must be highlighted because it will contain the task\_struct array.

```
ENTRY(main)
SECTIONS
{
    . = 0x10000; /* system code begins here */
    .text.main :
    {
        BYTE(24); /* reserved to store user code management data */
        *(.text.main) }
    .text : { *(.text) }
    .rodata : { *(.rodata) }
    .data : { *(.data) } /* Initialized data */
    .bss : { *(.bss) } /* Not initialized data */
    . = ALIGN(4096); /* task_structs array aligned to page */
    .data.task : { *(.data.task) }
```

Fig. 2. system.lids file. Linker script file for the system image.

```
union task_union task[NR_TASKS] __attribute__((__section__(".data.task")));
```

Fig. 3. Task\_struct array annotated with the section .data.task (sched.c)

## 2.2 Booting

The three blocks described in the previous section are appended to each other to create a single binary file with the content of all ZeOS operating system. If this file is copied in binary format to an unformatted floppy disk, beginning with sector 0, an image will be obtained of the OS that will boot automatically. It will be possible to run it on any computer with a floppy unit that can boot from a floppy disk.

The booting operation is simple. First turn on your computer. Before the OS is executed, a program that is placed in a read only memory (ROM) area will be executed, called the BIOS. This program checks that the PC is working properly. It is possible to see how the memory, the disk and other devices are checked during this process. It then looks for the preconfigured booting device and tries to access it. It does so by loading sector 0 from this device into the memory. It does not matter whether the device is a floppy disk, a hard disk, or a CD-ROM.

As the size of the first sector that the BIOS loads is quite small (512 bytes), there is not enough space to store the entire ZeOS. It is only possible to store a loader. Since ZeOS is very simple, the loader is designed to fit into 512 bytes and the BIOS will fully load it into the memory.

Once this 512 bytes of the boot sector are copied to the memory by the BIOS, it transfers the execution to the first byte of this small block. The boot loader is executed at this point. The main task of this boot sector code is to finish loading what is left in the image of the OS that is still on the floppy disk and load it into the memory. Basically the system and user areas. Finally, it transfers the execution to the first byte of the system code and starts the execution of your ZeOS.

## 2.3 Image construction

The way the OS image is built will be explained backwards. It will be assumed that the three blocks (boot sector, system area and user area) have been obtained and it will be seen how to put



them together into a single file like the snapshot in Figure 1. This process is done automatically using the Makefile that was provided with ZeOS.

The program that attaches the files is called *build* (*build.c*). This program is generated (by the *make* command) as follows:

```
$ gcc -Wall -Wstrict-prototypes -o build build.c
```

To attach the three blocks one after the other, it is only necessary to execute the following command (make does this for you):

```
$/build bootsect system.out user.out > zeos.bin
Boot sector 512 bytes.
System is 24 kB
User is 1 kB
Image is 25 kB
```

Where *bootsect* is the binary content of the boot sector; *system.out* is the binary content of the system area; *user.out* is the binary content of the user area; and *zeos.bin* is the resulting binary of the OS. *build* checks block sizes, adds the user and system sizes and writes this value in a specific boot sector position, specifically in bytes 500 and 501, which are labeled in the *bootsect.S* as *syssize*. This will be used by the bootloader to finish the operating system load in memory. Once in memory, the OS code must move the user code to its final position, the user code start, so it needs the total user size. That is why 16 bytes are reserved at the beginning of the system block, as can be seen in Figure 1. They are initially empty but the *build* program writes them with the sizes of the system and user images.

## 2.4 Introductory session

The main objectives of this section are:

- Become familiar with the working environment.
- Learn about the tools that must be used.
- Start analyzing and modifying the ZeOS code.
- Learn Bochs basic commands.
- Remember some of the concepts needed.

### 2.4.1 Getting started

This section describes the step to set up ZeOS. The following steps are performed supposing that the student works with an standard installation of Ubuntu 16.04 <sup>2</sup>:

- 1) Prepare the development environment:
  - a) Install basic development packages (*gcc*, *ld*, ...)

```
sudo apt-get install build-essential
```

- b) Install package for x86 tools (*as86* and *ld86*)

2. The development environment is already installed in the laboratory but it lacks the x86 tools, and you need to install them manually

```
sudo apt-get install bin86
```

c) Download and install Bochs 2.6.7

2) Download and install ZeOS:

- a) ZeOS is available as a compressed file (.tar.gz) at the web page.
- b) To install it, execute:

```
> tar xzfv zeos.tar.gz
```

3) Test the environment:

a) Generate your ZeOS (zeos.bin):

```
> make
```

b) Execute it using bochs with internal debugger:

```
> make emuldbg
```

c) If everything worked well, a new window similar to the one in Figure 4 will appear (it will show the output of your ZeOS) and the prompt for the internal debugger will show the current address to be executed and wait for new debug commands. Your emulated computer is ready to boot your ZeOS image. Press c and INTRO to check that it works.

### 2.4.2 Working environment

When an OS is being developed, a working environment is needed and it is usually offered by another operating system. In this case, to develop ZeOS we will use an Ubuntu system with the following configuration:

- OS. Ubuntu 16.04<sup>3</sup>
- Compiler. GCC 4.x
- Emulator. Bochs version 2.6.7 (bochs.sourceforge.net)

This environment will enable you to generate the ZeOS operating system. Once ZeOS has been generated, it can be used to boot your computer (after copying it to a floppy disk). However, you should bear in mind that the system will be recompiled and loaded many times, so it is advisable to run your system on an architecture emulator (like Bochs) in order to save time. This way, when modifications are necessary, only the emulator has to be rebooted rather than the computer.

Whenever you need to compress and extract your files use the following commands:

- Compress: `tar czvf file_name.tar.gz file_list_to_compress`
- Extract: `tar xzvf file_name.tar.gz`

These commands use *tar* and *gzip* at the same time. Not all shells support this, so, if it is not available, the commands would be:

- Compress: `tar cvf file.tar files_to_unify` and then `gzip file.tar`
- Extract: `gunzip file.tar.gz` and then `tar xvf file.tar`

In Ubuntu, you will be able to work with several editors (GVim, emacs, nedit, etc.).

3. The system has the user alumne (pass: sistemas).

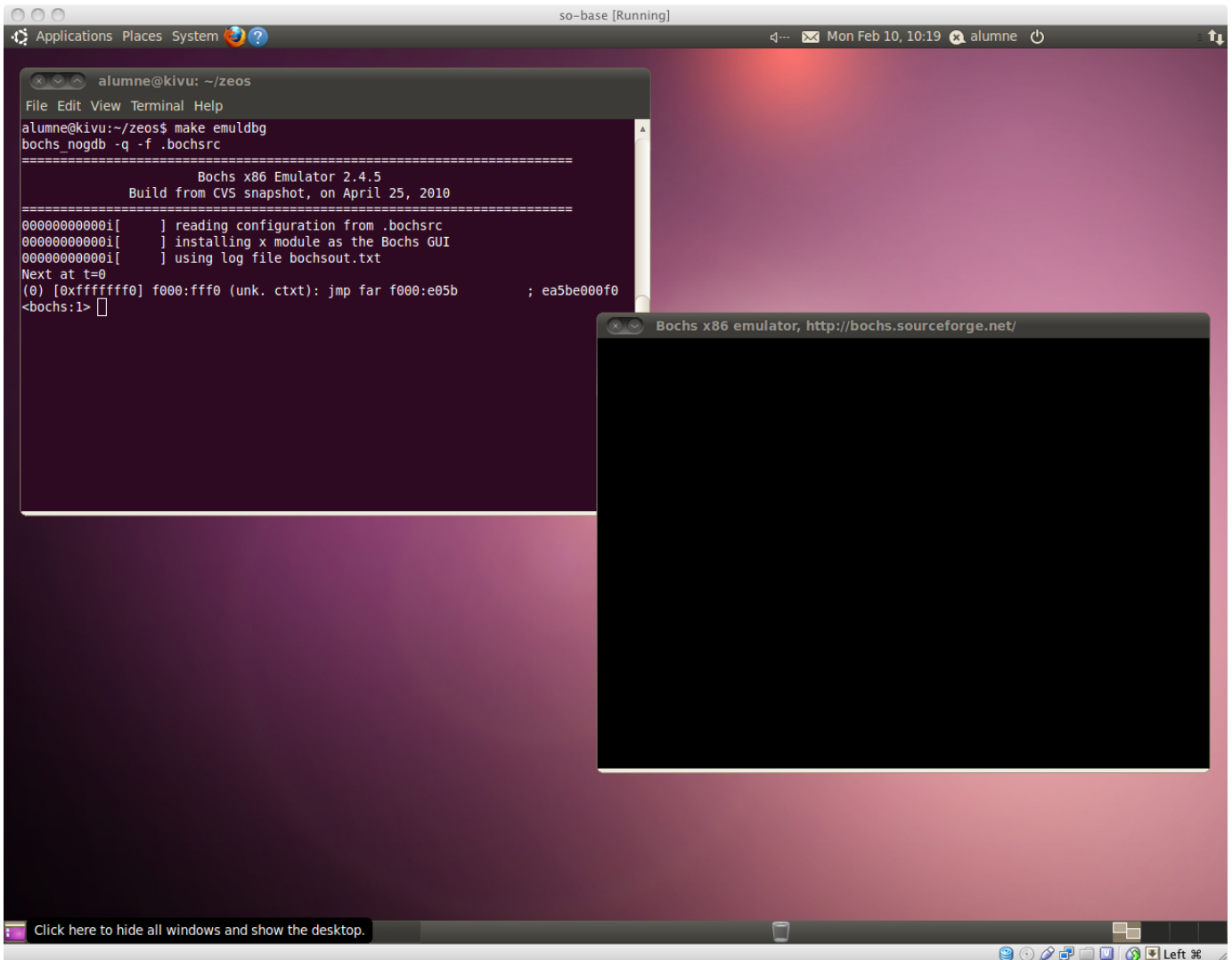


Fig. 4. Working environment with the Bochs commands window (left) and Bochs emulation window (right)

### 2.4.3 Prior knowledge

For this document to be understood, it is assumed that you have acquired knowledge from other courses and that you are able to work in specific environments. More specifically, it is assumed that you are able to:

- Write programs in C language.
- Write programs in Linux i386 assembler language.
- Add assembler code to a C file.
- Modify a Makefile to add new rules or modify existing ones.

If you lack any of these skills, you should find additional information to that given on this course. The information in the bibliography section on the course's web page may also be useful.

## 2.4.4 Bochs

Bochs is a PC Intel x86 emulator written in C++. It was created in around 1994. Originally it was not free, but when Mandrake bought it, it was granted a GNU LGPL license. It is a little slow, although not to any noticeable extent for the purposes of this project. However, it is very reliable. In this document we will use **version 2.6.7**.

### Execution

The Bochs executable is `/usr/local/bin/bochs`. This will read a configuration file to prepare the emulated computer and will start its execution.

### Debugging using Bochs

In order to debug your operating system, it is necessary to control execution using a debugger. Bochs offers two options for debugging code: using an external debugger (such as gdb) or using an internal debugger that is part of the emulator. The two options are *exclusive*, namely, the Bochs executable can support only one. Both versions are available in the laboratory (*bochs* and *bochs\_nogdb*), so you can choose whichever version you prefer. The internal debugger is recommended when you want to debug assembler code or access or consult specific information about the emulated hardware (namely physical addresses translation, special registers, ...). GDB is more suitable to debug programming problems not related with the underlying hardware. In summary, start using GDB and use only the bochs internal debugger when everything else fails.

**If you are working at the lab or you use the provided virtual image you can skip the rest of this section.** Otherwise, if you are working at home, you need to compile the Bochs source code in order to use the Bochs debugging facilities. It is necessary to activate some debugger options during the compilation because the option to use the debugger is disabled by default in the Bochs Binary Standard Package.

To enable the external GDB debugger you need to execute *configure* with parameter *-enable-gdb-stub*, and recompile Bochs:

```
$ ./configure --enable-gdb-stub
$ make all install
```

If you want to activate the internal debugger you must recompile Bochs using:

```
$ ./configure --enable-debug --enable-disasm
$ make all install
```

### Bochs configuration file

Bochs uses a file to configure the features of the emulated computer<sup>4</sup>. We will focus on the following three points:

- **Image location.** Line 5 defines the file *zeos.bin* as ZeOS image to boot. The file name must correspond to the path that points to the *zeos.bin*. Currently it is set to load the image *zeos.bin* from the current directory.

4. <http://bochs.sourceforge.net/doc/docbook/user/bochsrc.html>

```
floppya: 1_44= ./zeos.bin, status=inserted
```

- **Boot device.** Line 6 defines the device used to boot the machine. In our case, the floppy device.

```
boot: floppy
```

- **External debugger configuration.** In order to use the gdb debugger it is necessary to add the following line to configure it. If this line is not added, even Bochs has been compiled to use the external debugger, it will execute without any debugging support.

```
gdbstub: enabled=1, port=1234, text_base=0, data_base=0, bss_base=0
```

It should be highlighted that this line is only interpreted when the Bochs version is compiled to use the external debugger. If this line is present and Bochs does not have support for the external debugger it will fail with an error. Due to this, we provide you with two configuration files: 1) *.bochsrc* and 2) *.bochsrc\_gdb*.

In order to use the needed configuration file you can use the Bochs option `-f`<sup>5</sup>:

```
$ bochs -f .bochsrc_gdb -q
```

If you do not use this option, Bochs assumes a default configuration file named *.bochsrc*.

### Controlling execution through internal Bochs debugger

- Compile the ZeOS code with debug symbols, using `"-g"` flag in CC (it is enabled by default in the Makefile).
- As the provided configuration file for this version of bochs is called *.bochsrc*, it is not necessary to use the `-f` option.
- Execute the bochs version without GDB<sup>6</sup>.

```
$ bochs_nogdb -q
```

Once the version of bochs with the debugger option is executed, you should see two windows like the ones shown in previous Figure 4: the emulation window and the commands window (the same where bochs has been executed).

The commands window displays a prompt in which commands can be introduced. This prompt shows information about the memory address and its contents that the debugger is about to execute. A summary of the commands that can be executed may be found at <http://bochs.sourceforge.net/doc/docbook/user/internal-debugger.html>.

5. The target *gdb* from the Makefile has the same effect.

6. The target *emuldbg* from the Makefile has the same effect.

## Controlling execution through an external GDB debugger

To use GDB (GNU debugger) to debug your operating system, follow the steps below:

- Compile the ZeOS code with debug symbols, using "-g" flag in CC (it is enabled by default in the Makefile).
- In order to use the gdb version, the configuration file needs the line enabling the gdbstub. This means using the provided *.bochsrc\_gdb*.
- Execute the bochs version with the external debugger activated<sup>7</sup>.

```
$bochs -f .bochsrc_gdb -q
```

The Bochs virtual machine will start (opening the emulation window) and will wait for the connection request from GDB.

- After executing bochs, a new terminal is needed to execute GDB or a front-end to GDB (for example ddd or gdbtui<sup>8</sup>). GDB accepts an executable program as the argument to be debugged. The system binary file of our Zeos will be usually used.

```
$gdb system
```

- Connect GDB to Bochs by executing the command below in GDB (the port should be the same as the one you use in *.bochsrc*).

```
(gdb) target remote localhost:1234
```

- Add user symbols.

```
(gdb) add-symbol-file user
```

Figure 5 shows the appearing windows when bochs is started with gdb enabled. The Makefile's target *gdb* automatically does all this under the hood and therefore it appears directly at the gdb prompt already connected to the bochs instance.

- Now you can now add breakpoints, continue the execution, etc. You can find a reference guide with the main GDB commands on the following link: <http://refcards.com/docs/peschr/gdb/gdb-refcard-a4.pdf>

### 2.4.5 Frequently used commands

Table 1 presents some of the most frequently used commands for bochs and gdb debuggers.

7. The target *gdb* from the Makefile has the same effect.

8. <https://sourceware.org/gdb/onlinedocs/gdb/TUI.html>

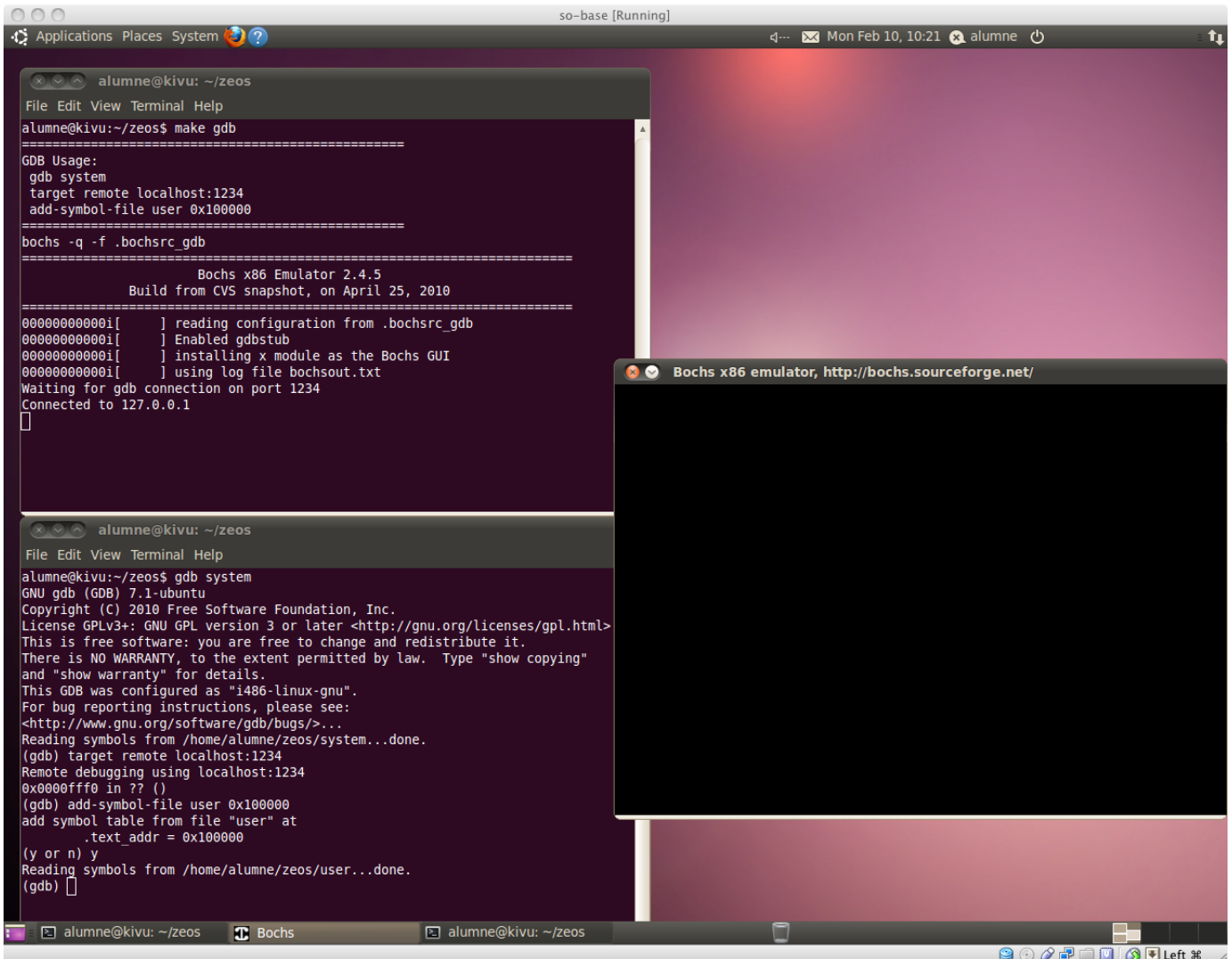


Fig. 5. Bochs commands windows (up), Bochs emulation window (right) and gdb window (down) after establishing a connection between gdb and bochs.

## 2.5 Work to do

The `user.c` file is the place in which the user code must be written. **We strongly recommend that you only modify the `user.c` file.**

In order to familiarize with the environment is advisable to carefully follow and try the steps described in this section. You will know how to:

- Use the debugger to control and view information about your code.
- Visualize the generated object code.
- Locate variables and functions in object code.
- Modify the user code.
- Mix assembler and C in your code.

Bochs	GDB	Description
<b>help</b>		Show the commands that can be executed.
<b>continue</b>		Execute the image normally until the next breakpoint (if you have one). While executing the image, you can not issue new debugger commands, so in order to <b>stop the execution you should press the Ctrl-C</b> key combination, and the debugger prompt will appear again.
<b>step</b> [ <i>num</i> ]	<b>stepi</b>	Execute <i>num</i> assembler instructions. Execute a single one if no <i>num</i> is provided.
<b>next</b>	<b>nexti</b>	Execute an assembler instruction. If the instruction is a <i>call</i> , the whole function is executed.
<b>break</b> <i>address</i>	<b>b</b> <i>*address</i> <b>b</b> <i>function</i>	Insert a breakpoint in the instruction indicated in the address. The address can be written in decimal (123459), octal (0123456) or, as is usual, in hexadecimal (0x123abc). For example, if you write "b 0x100000" (or b *0x100000 in the GDB), it will insert a breakpoint in the first code line of the user process. GDB also accepts a function name or a specific point in a file ( <i>filename:linenumber</i> )
<b>r</b>	<b>info r</b>	Show the content of the registers in the mode currently being used (user or system)
<b>print-stack</b> [ <i>num_words</i> ]	<b>x/16 \$esp</b>	Print <i>num_words</i> from the bottom of the stack. If <i>num_words</i> is not specified, the default value is 16. It is only reliable if you are in system mode, when the base address of the stack segment is 0.
<b>info tab</b>		Show paging address translation.
<b>quit</b>		Exit debugger.

TABLE 1  
Frequently used commands for bochs and gdb

### 2.5.1 Where is ZeOS hanged?

In this section you must prove that after running the ZeOS image it is executing the infinite loop in the user code. To do that you must:

- 1) Generate a ZeOS image. Check whether there are any warnings when it is compiled. Correct any errors that may have occurred.
- 2) Run the bochs debugger and execute this image.
- 3) When the ZeOS seems hanged, which instruction is executing?
- 4) Where is that instruction in the code?
  - To interpret the code execution, it is useful to know the memory location of each section of the image. For that, the commands *objdump* and *nm* can be used.
  - The *objdump* command shows the compiler-generated code and the corresponding memory addresses, which can be useful to insert breakpoints (Figure 6 shows the result of using *objdump* with the system file). The *nm* command shows the location of variables and functions. For instance, Figure 7 shows that the task struct array is located at address 0x11000.

### 2.5.2 User code modification

- Modify the *main()* of the *user.c* file as Figure 8, adding some local variables and a call to an *outer* function.
- Add the functions *inner* and *outer* to the *user.c* file just before the *main* as in Figure 9.
- Generate a new ZeOS image. Check whether there are any warnings when it is compiled. Correct any errors that may have occurred.
- Execute the command "*objdump -d user.o*" to see the code in the assembler language of the file *user.c*. Note that the addresses of the symbols that appear are in PIC (Position Independent Code) and that they start from zero.



```

00010000 <main-0x4>:
    10000:      10 00                adc    %al, (%eax)
    ...
00010004 <main>:
    10004:      55                push   %ebp
    10005:      89 e5            mov    %esp, %ebp
    10007:      83 ec 08        sub    $0x8, %esp
    1000a:      83 e4 f0        and    $0xffffffff0, %esp
    1000d:      6a 00            push   $0x0
    1000f:      9d                popf
    10010:      ba 18 00 00 00    mov    $0x18, %edx
    10015:      b8 18 00 00 00    mov    $0x18, %eax
    1001a:      fc                cld
    1001b:      8e da            mov    %edx, %ds
    1001d:      8e c2            mov    %edx, %es
    1001f:      8e e2            mov    %edx, %fs
    ...

```

Fig. 6. "objdump -d system" output sample

```

0001042c T set_ss_pag
00010224 T set_task_reg
000102b8 T setTrapHandler
000100f8 T setTSS
00011000 D task
0001c000 D taula_pagusr
0001d020 B tss
00010538 D usr_main
0001d8a0 B x
00010544 D y

```

Fig. 7. "nm system" output sample

- Now execute "objdump -d user". *Do the addresses match? Why?*
- Run the whole program (*continue* command). To regain control in the debugger you have to type *Ctrl+C*. After doing so, *which line is going to be executed?*
- Re-execute the program step by step. To do this, add a BREAKPOINT to the first instruction of the inner routine and continue the execution using the *continue* command. Check that the execution stops at the breakpoint line. Then, continue the code execution line by line (*step* and/or *next*<sup>9</sup>).

**Note:** Adding a BREAKPOINT is often the only useful way to reach to a subroutine address. In this case, the number of STEP commands that may be needed to reach the *inner()* routine is extremely high. In the ZeOS project, one BREAKPOINT is the only way to ensure that the execution jumps to an interrupt service routine (ISR).

- Examine the content of the registers.
- Which debugger command enables you to see the content of the memory? Does it work to see the content of the *acum* variable from the *user.c* file?

9. It has been seen that after typing *Ctrl+C*, Bochs is unable to continue the execution with the debugger instructions *step* or *next*. The commands *continue*, *stepi* or *nexti* must be used

```

int __attribute__ ((__section__(".text main")))

main()
{
    long count, acum;
    count = 75;
    acum = 0;
    acum = outer(count);
    while (1);
    return 0;
}

```

Fig. 8. Modification to the main() function in the user.c file

```

long inner(long n)
{
    int i;
    long suma;
    suma = 0;
    for (i=0; i<n; i++) suma = suma + i;
    return suma;
}
long outer(long n)
{
    int i;
    long acum;
    acum = 0;
    for (i=0; i<n; i++) acum = acum + inner(i);
    return acum;
}
int __attribute__ ((__section__(".text.main")))
main()
. . .

```

Fig. 9. New functions *inner* and *outer* to the added to the user.c file

- How can you see the value of acum in an easy way? (TIP: think about where C stores the returning value of a function)
- How can you see the parameters passed to the outer function?

### 2.5.3 Use of inline assembler

Add one function to the user.c file using the following header:

```

int add(int par1,int par2);

```

This function must perform (in assembler) the same as:

```

return par1+par2;

```

See the document in the bibliography about how to add assembler code to a C program. To do the sum, you must access the parameters (remember that they are on the stack), add them and return the result. Two different versions will be written:

- The entire code of the function will be written in assembler language. This means that you must review the way a value is returned in assembler; (in which register will the result have to be placed?). You cannot use the input/output parameters of inline assembler.
- Use assembler inline. A local variable must be defined in C and the result of the addition will be put in this variable. See the appendix to ascertain how to modify a defined variable from C in an assembler code. In this case, you must use the input/output parameters of inline assembler.

Add a call to the function `add()` in the `main()` function of the `user.c` file by putting the result of the addition in a local variable.

- Check with the debugger that both versions work correctly.
- Which address did you enter in the Bochs debugger to access the variable in which you put the result of the `add` function?

#### 2.5.4 Additional questions

Try to understand the relevant parts of the boot mechanism (`bootsect.S`) and system initialization (`system.c`). You can draw a structure diagram of the files and their contents. You should be ready to answer the following questions:

- 1) What are the main differences between real mode and protected mode?
- 2) In which mode is the system when it is booting?
- 3) What are the execution privilege levels?
- 4) What is the bootloader? In which part of the disk is it found?
- 5) What is the BIOS? What is it for?
- 6) What steps does ZeOS Makefile perform? What are the uses of the various compilation and linkage options?
- 7) Optionally, you can try to scroll the screen if the cursor surpasses the maximum coordinates.

### 3 MECHANISMS TO ENTER THE SYSTEM

To execute code from the OS you must use the interrupt mechanism. Interrupts are events that break the sequential execution of a program and they can be classified as synchronous or asynchronous depending on when they are generated. Synchronous interrupts are generated by the CPU at the end of the execution of an instruction. Asynchronous interrupts are generated by other hardware devices.

From now on, we will refer to synchronous interrupts as **exceptions** and to asynchronous interrupts as **interrupts**, although both follow the same mechanism.

Each interrupt is identified by a number between 0 and 255 and associated with a specific function using a system table known as an interrupt descriptor table (IDT). Each entry in the table has the information necessary to do the actions required when an interrupt is produced. Among other data, it contains the address of the routine to be executed (the **handler**) and the minimum code privilege level to execute it.

The interrupt id numbers are as follows:

0-31 are exceptions and unmasked interrupts. 32-47 are masked interrupts. 48-255 are software interrupts. For example: Linux uses id 0x80 (128) for the system calls.

These are the same numbers that appear in the Intel manuals. For more information, see chapter 4 of *Understanding the Linux Kernel*.

To write an exception or interrupt (except the system call), you must:

- 1) **Initialize** the IDT entry that matches the exception or interrupt id number.
- 2) **Write the handler** to be executed when the interrupt or the exception is generated. The handler is the assembler code implemented to manage the call to the service routine.
- 3) **Write the service routine** to the interrupt or the exception, which will be the code to perform the associated service to the interrupt or exception.

An **additional** step is added for the system calls to isolate user codes from low-level and non-portable code. A function responsible for passing the system call parameters, generating the interrupt and recovering its result should be written for each system call. These functions will be named **wrappers**. The user code invokes any of this wrappers (as any other user function) that then causes the actual enter to the system.

In this section you will find:

- Hardware management of an interrupt.
- The code for managing exceptions.
- The code for managing the keyboard interrupt.
- The generic code for making system calls.
- The code for implementing a write system call.

The following sections include a how-to guide to add an exception, an interrupt and a system call. **Remember that they are almost the same, since they are all managed by the IDT.** Although they are described separately in this guide, it is important to understand that they are quite similar.

#### 3.1 Preliminary concepts

- Interruption, exception and system call.

- Context of a process.
- Hardware management of an interrupt.
- Checking parameters and returning results in a function.

### 3.2 Hardware management of an interrupt

Once the system has been initialized, the hardware automatically performs the following steps when the CPU detects that an interrupt has been generated (for more information, see *Understanding the Linux Kernel*):

- The  $i$  index of the interrupt vector is determined and the corresponding IDT entry is accessed.
- The hardware verifies that the interrupt has enough privileges to execute the handler, by comparing the current privilege level with the one stored in the IDT entry. If access is unauthorized, a general protection exception is generated.
- The privilege level of the handler routine is checked to see if it is different from the current execution level (our case), in which case the stack will have to be changed.
- To make the change from user stack to system stack, the system takes the address of system stack from the Task State Segment (TSS). Once the `ss` and the `sp` of the system are known, the hardware saves: 1) the values of the `ss` and `esp` of the previous stack in the new stack (the user stack in this case); 2) the value of the `eflags` (state word); and 3) the values of `cs` and `eip` (address that generates the interrupt). The resulting stack is shown in Figure 10
- The address saved in the  $i^{th}$  entry of the IDT is executed.

**The hardware performs the aforementioned steps automatically.**

Once completed the interrupt management code, the control must be returned to the instruction that generated the interrupt by executing the `iret` instruction. The `iret` instruction takes the value of the `eip` and `cs` registers from the top of the stack, loads the `eflags` registers with the stored value in the stack and modifies the `esp` and `ss` registers to point to the stack that was in used before the interrupt happened (user stack in our case).

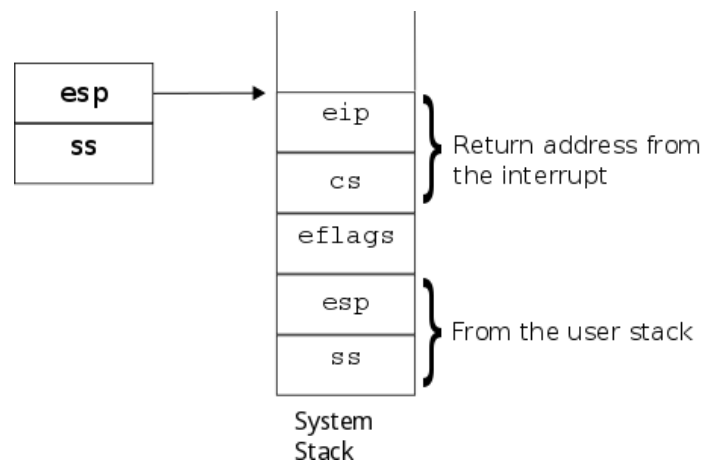


Fig. 10. State of the system stack when the execution of the interrupt handler begins

#### 3.2.1 Task State Segment (TSS)

The 80x86 architecture defines a specific segment called the task state segment (TSS) to store task related information in memory. It is mainly used to implement a context switch in hardware.

Neither ZeOS nor Linux want to use the TSS but the architecture forces them to define a TSS for each cpu (1 in our case). **We use the TSS to know the address of the system stack when making a user to system mode switch.** You can check the ZeOS files to see the fields that the TSS has and how it is initialized for the initial process.

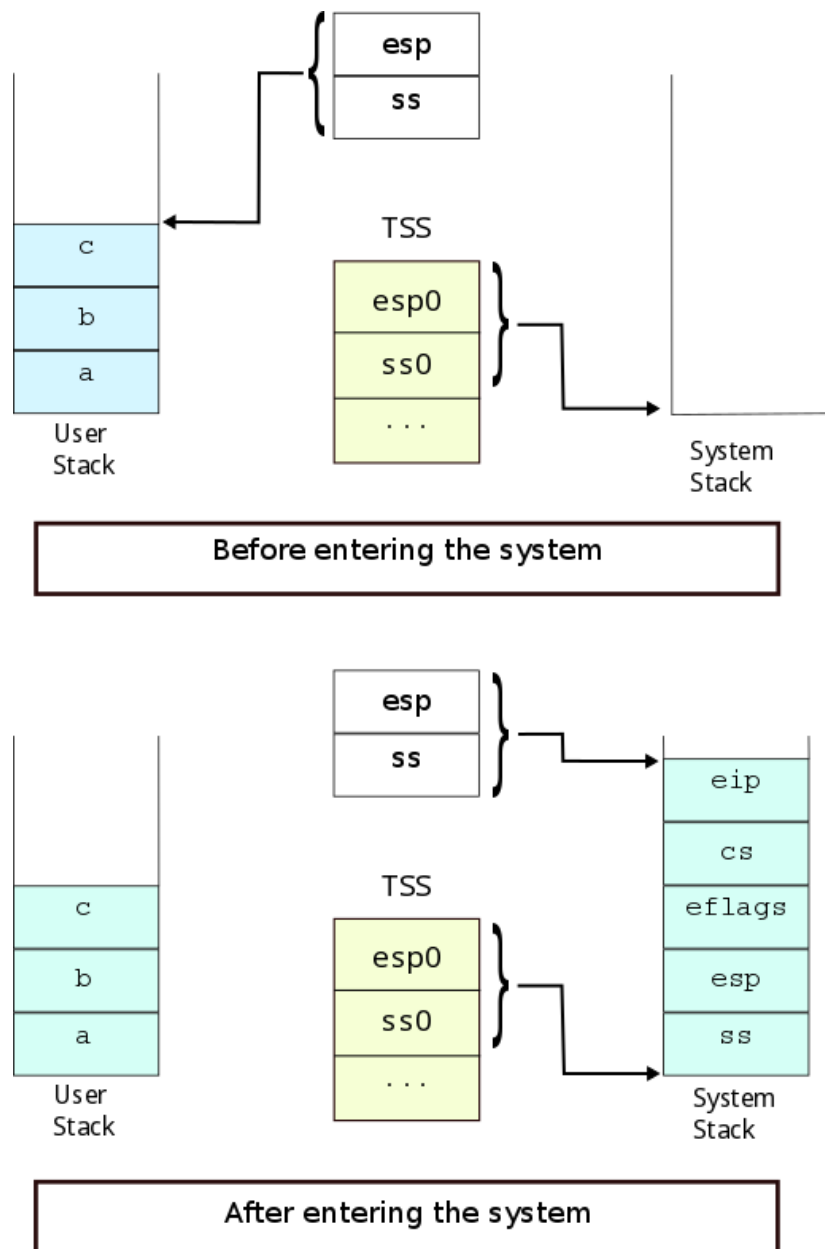


Fig. 11. The esp0 and ss0 fields of the TSS show the location of the system stack used in switch from user mode to system mode

Figure 11 shows how the ss0 and esp0 fields of the TSS are used when switching from user to system mode. These fields always point to the bottom of the system stack.

### 3.3 Function name conventions

The functions will be given similar names to make the code easier to follow. The handlers will be given the same name as the interrupt followed by "\_handler" and the routine services will

be given the same name as the interrupt or exception, followed by "\_routine". For example, the handler name of the exception "divide error" will be "divide\_error\_handler" and its routine service name will be "divide\_error\_routine".

As regards the system calls, the handler will be named "system\_call" and the service routines "sys\_namesyscall", for example "sys\_write".

### 3.4 Files

ZeOS files generally have functions that share a common objective or type. The main files (probably not the only ones) in this section are:

- system.c: System initialization
- entry.S: Entry system points (handlers)
- interrupt.c: Interruption service routines and exceptions
- sys.c: System calls
- device.c: Device dependent part of the system calls
- libc.c: System call wrappers
- io.c: Basic management of input-output

### 3.5 Programming exceptions

The system exceptions will be programmed in this section (except for the 15th exception, which is reserved for Intel). You will have to enter the system in order to catch an exception and execute the routine that you have created to handle it, as they are protected operations. The OS usually tries to recover from the exception and restore the system but in ZeOS the exception service routines are very simple: the screen will show that an exception has been generated and specify which one, and the system will remain in an infinite loop state (it will never return to the user mode).

#### 3.5.1 Exception parameters

Some exceptions receive additional information from the hardware in order to be solved (for instance, information on the type of access that generates the page fault exception). This information has a fixed size (4 bytes) and is pushed into the system stack automatically. Only the exceptions indicated in Table 2 receive these parameters. Figure 12 shows the system stack contents and the location pointed by the esp and ss registers when an exception with an error parameter is raised.

#### 3.5.2 Initializing the IDT

To initialize a position in the IDT, the following function is provided:

```
setInterruptHandler(int position, void (*handler)(), int privLevel)
```

The parameters required to write the exception are:

- int position. Entry in the IDT.
- void (\*handler)(). Address of the handler that will handle the exception<sup>10</sup>.

10. If you do not know how to declare and use the function pointers, consult the C Programming appendix

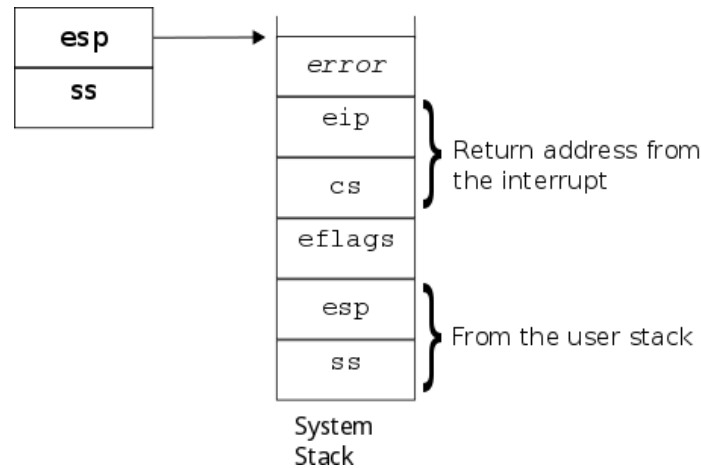


Fig. 12. System stack contents when an exception with an error parameter is raised

- `int privLevel`. Privilege level necessary to execute the handler. There are two possibilities: 0- Kernel, 3- User.

Table 2 shows the exception positions in the IDT, their names and the number of bytes of their parameters (or error code<sup>11</sup>) if they exist.

# IDT	Exception	Parameter
0	Divide error	No
1	Debug	No
2	NMI	No
3	Breakpoint	No
4	Overflow	No
5	Bounds check	No
6	Invalid opcode	No
7	Device not available	No
8	Double fault	4 bytes
9	Coprocessor segment overrun	No
10	Invalid TSS	4 bytes
11	Segment not present	4 bytes
12	Stack exception	4 bytes
13	General protection	4 bytes
14	Page fault	4 bytes
15	Intel reserved	No
16	Floating point error	No
17	Alignment check	4 bytes

TABLE 2  
List of system exceptions

So, to handle all those exceptions, one call to `setInterruptHandler` for each exception has been added in the file `interrupt.c` (implemented in function `set_handlers` in `libzeos.a`).

### 3.5.3 Writing the handler

**A handler must be written for each exception.** All exception handlers have a common scheme and must follow these steps (in assembler):

11. For some exceptions, the CPU generates a hardware error code and puts it in the system stack before it starts the execution of the exception handler



- 1) Define the function header in assembler. The **ENTRY** macro<sup>12</sup> can be used. The header is defined by calling it at the beginning of the code of the function and passing the name of the exception handler as a parameter (**ENTRY(handler\_name)**).
  - From this point on, the **ENTRY** macro will always be used to define a function header in the assembler.
- 2) Save the context in the stack (use the **SAVE\_ALL** macro).
- 3) Call the service routine.
- 4) Restore the context from the stack (use the **RESTORE\_ALL** macro). Notice that the order of restoring the values of the registers from the stack must match the order in the **SAVE\_ALL** macro.
- 5) Clear the error code (if the exception has one), removing it from the stack. Check Table 2.
- 6) Return from the exception. Since it is not a "normal" return because it has to change mode, an **iret** rather than a **ret**<sup>13</sup> will be used for the return.

All exception handlers are already implemented inside the *libzeos.a* library.

### 3.5.4 Writing the service routines

The exception management in this operating system will be very simple. Whenever an exception is raised, the OS will show a message with a short description and will wait in an infinite loop, stopping the system.

For example, the code of the general protection service routine is:

```
void general_protection_routine()
{
    printk("\nGeneral protection fault\n");
    while(1);
}
```

All exception service routines are already implemented inside the *libzeos.a* library.

## 3.6 Programming interrupts

Figure 13 shows the main steps for the clock interrupt. This interrupt can arrive at any time and deal with any point in the code. If the IDT has been correctly programmed, the `clock_handler` function programmed in the `entry.S` file will be the first part to be executed. Inside this function we find a call to the `clock_routine` function who will make the real interrupt management. The steps for the keyboard interrupt are similar and will be programmed in this section.

### 3.6.1 The keyboard interrupt

The keyboard interrupt will display the character that corresponds to the pressed key. Figure 14 shows the expected result with the character displayed in a fixed place on the screen.

To write an interrupt such as the keyboard interrupt, you must:

- Initialize the IDT as it is done for the exceptions. The keyboard interrupt is contained in entry 33.

12. The corresponding annex shows how the macro mechanism of the assembler works

13. You need to know which data are saved in the stack when going into system mode and what the `iret` instruction does.

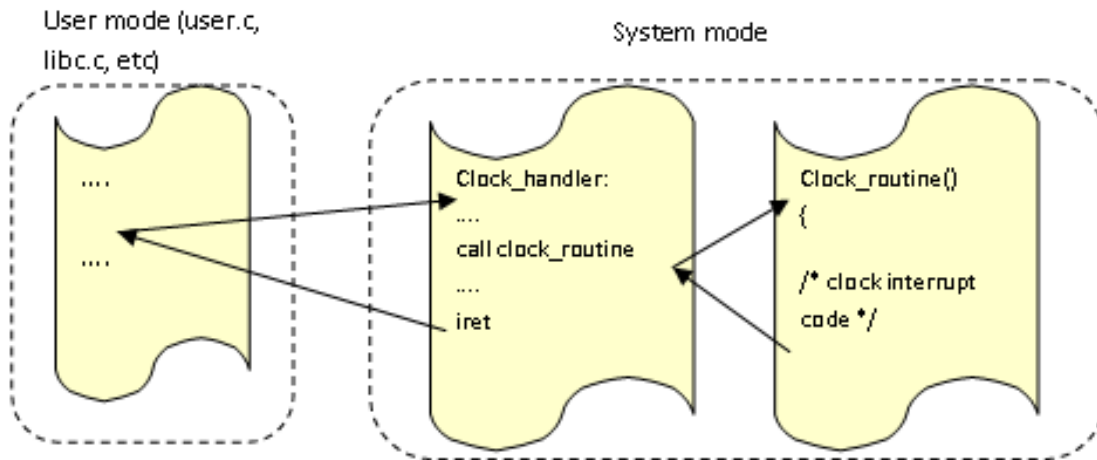


Fig. 13. Snapshot of the clock interrupt



Fig. 14. Keyboard interrupt

```
setInterruptHandler(33, keyboard_handler, 0);
```

- Write the handler.
- Write the service routine.
- Enable the interrupt. The key interrupt is a *masked* interrupt that is disabled in the code provided, which means that it is ignored by default. It is necessary to modify the mask that enables the interrupts, located in the **enable\_int()** routine (file **hardware.c**).

Notice that just after enabling the interrupts, one of them can be raised at any time. When this occurs, you must ensure that the system is fully initialized because the interrupt service routines may access any system structure. For this reason, the best place to enable the interruptions is just after all the OS services have been started and before the return to user mode is performed.

### 3.6.2 Writing the handler

An interrupt handler is written more or less the same way as an exception handler. Follow the steps below:

- 1) Define an assembler header for the handler.
- 2) Save the context.
- 3) Perform the **end of interruption (EOI)**. You must notify the system that you are treating the interrupt. For this, a new macro, named EOI, has been created:

```
#define EOI      \
movb $0x20, %al ; \
outb %al, $0x20 ;
```

- 4) Call the service routine.
- 5) Restore the context.
- 6) Return from the interrupt (be careful, since you are returning to user mode).

### 3.6.3 Writing the service routine

The service routine will display on the screen the character that corresponds to the pressed key.

This service routine performs the following steps:

- 1) Read the port of the data register (0x60) with the routine of the io.c file:

```
unsigned char inb(int port)
```

- 2) Once the value is read from this port, it must distinguish between a make (key pressed) or a break (key released). The contents of this register is shown in Figure 15. Bit number 7 specifies whether it is a make or a break. Bits 0..6 contain the code to be translated into a character.
- 3) If it is a make, the translation table `char_map` at `interrupt.c` must be used to obtain the character that matches the scan code.
- 4) Print the character on the upper left of the screen. For this the `putc_xy` function in `io.c` is used.
- 5) If the pressed key does not match any ASCII character (Control, Enter, Backspace, ...) a capital C will be displayed.

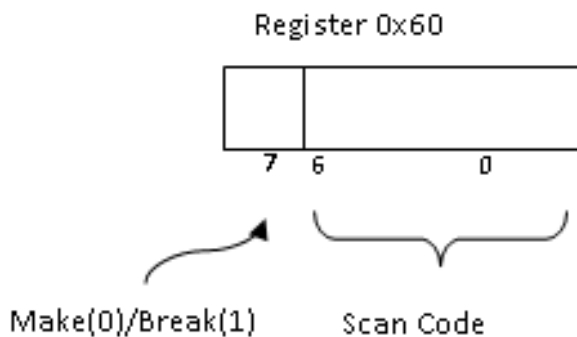


Fig. 15. Contents of the keyboard 0x60 register

### 3.7 Programming system calls

The first system call **write** will be written in this section. It will enable strings to be printed on the screen from user mode. For more information, see Chapter 8 of *Understanding the Linux Kernel*.

As you have seen in other courses, system calls have a common entry point: the 0x80 interrupt. This means that with one sole entry point (1 handler) there will be N system calls. In this section we will present all the common functions to all the system calls (and also the data structures) and the specific functions of the write call.

Figure 16 shows the steps followed by a system call. The user code calls what it believes to be the system call, but it is actually only an adapter. The library code implements the adapter, which you will call a **wrapper**, that is responsible for **passing parameters** between the user and the system, generating the **trap** (int \$0x80) and processing the result.

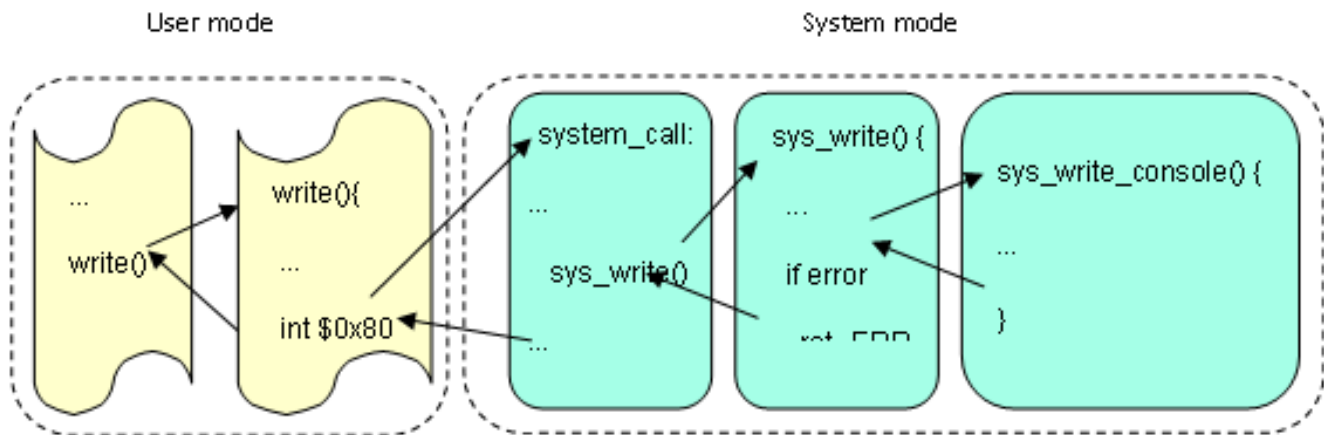


Fig. 16. Snapshot of a system call

Once the interrupt is generated, it is executed as another interrupt (the execution of the handler and the execution of the service routine). In this case, interrupt 0x80 does not need to be enabled since it is unmasked. However, the IDT does need to be initialized. As explained above, all system calls have a common handler, which will be named **system\_call**. In this handler, specify the system call you want to make and execute the corresponding service routine (**sys\_write** in the example).

#### 3.7.1 Independence from devices

As explained in other courses, one of the basic principles of an OS design is independence between devices; this is, the interface for input and output system calls is the same regardless of the device that is requesting its service.

Keeping this in mind, you will design each call considering a part that is dependent and a part that is independent from the device. Thus, the system will be prepared to grow in complexity, or to have more devices added to it. The design will be similar to a real OS.

For example, the **sys\_write** call (independent part) will call **sys\_write\_console** (dependent part), in order to print on the screen.

### 3.7.2 Returning results

The convention in Linux for returning the results of the system calls considers a positive or zero return when the execution of the call is correct and a negative value when an error is detected. In case of error, this negative value specifies the error.

Nevertheless, the returned value of the system call does not reach the user directly, but instead it is processed by the wrapper of the system call that is found in the library. Thus, if the returned value is positive or zero, it will be returned to the user as is. However, if the value is negative, the wrapper will save the absolute value of the return in an error variable defined in the library (called **errno**) and return -1 to the user to notify that the system call has an error. If the user requires further information about the error, he/she must consult the **errno** variable. If the system call does not return an error, the **errno** variable is not modified.

To use this convention in ZeOS, **all** system calls must return a negative number and specify the generated error in the event of a wrong execution. The constants used to identify these errors are contained in the Linux **errno.h** file. Moreover, your system call wrappers must process the negative values of the calls, update the **errno** variable with the absolute value of these errors and return -1 to the user.

A function should also be added for the users to obtain information about the error generated by a specific system call. Thus, the **void perror()** function will write an error message for the standard output according to the value of the **errno** variable.

**Note:** In the description of the system calls that appear in this document, the returned values are considered from the point of view of the user, which means that -1 will always be returned in the case of an error.

### 3.7.3 Writing the write wrapper

The write wrapper header is:

```
int write (int fd, char * buffer, int size);
```

Wrappers must carry out the following steps:

- 1) Parameter passing: ZeOS implements the parameter passing from user mode to system mode through the CPU registers, as occurs in Linux. The parameters of the stack must be copied to the registers **ebx**, **ecx**, **edx**, **esi**, **edi**. The correct order is the first parameter (on the left) in **ebx**, the second parameter in **ecx**, etc.
- 2) Put the identifier of the system call in the **eax** register (number 4 for write).
- 3) Generate the trap: **int \$0x80**.
- 4) Process the result (see section 3.7.2).
- 5) Return.

### 3.7.4 IDT initialization

Initialize the entry point of a system call using the call:

```
void setTrapHandler(int posició, void (*handler)(), int nivellPriv)
```

This function is similar to `setInterruptHandler`. In this case, since the system calls are invoked from the user privileged level, the value for the third parameter will be 3:

```
setTrapHandler(0x80, system_call_handler, 3);
```

### 3.7.5 Writing the handler

The steps are as follows:

- Save the context.
- Check that it is a correct system call number (that the system call identifier belongs to a defined range of valid system calls). Otherwise, the corresponding error must be returned.
- Execute the corresponding system call.
- Update the context so that, once restored, the result of the system call can be found in the `eax` register. Remember that C functions store their result in the `eax` register<sup>14</sup>.
- Restore the context.
- Return to the user mode.

A table is needed to relate the identifier of each call to its routine. The table will be filled in as new system calls are added.

This call table called `sys_call_table` will be defined in assembler. This table will be filled with entries like:

```
ENTRY(sys_call_table)
.long sys_ni_syscall // sys_ni_syscall address (not implemented)
.long sys_functionname // sys_functionname address
```

Bear in mind that non-implemented calls in a valid range must implement a call to `sys_ni_syscall`. This call will only return a negative identifier of the generated error. Each line is an entry of the table.

The instruction to execute the system call indexed by the register `eax` is:

```
call *sys_call_table(, %eax, 0x04);
```

So, the syscall table must be extended to enable the write system call. In particular, entries from 0 to 3 must be initialized to `sys_ni_syscall` to return the corresponding error (the syscall does not exist). Entry 4 points to the service routine (that will be written in section 3.7.6):

```
ENTRY (sys_call_table)
    .long sys_ni_syscall    // 0
    .long sys_ni_syscall    // 1
    .long sys_ni_syscall    // 2
    .long sys_ni_syscall    // 3
    .long sys_write         // 4
.globl MAX_SYSCALL
MAX_SYSCALL = (. - sys_call_table)/4
```

14. If the context is not modified, the return value will be the system call identifier

- Visit <http://lxr.linux.no> to look at the Linux declaration of the `sys_call_table` table.

Finally, if the number of the syscall is outside the correct range (from 0 to 4 in this case), the `sys_ni_syscall` will be called. This function must always return the corresponding error:

```
int sys_ni_syscall()
{
    return -38; /*ENOSYS*/
}
```

The code for the OS entry point (position 0x80 of the IDT) looks like:

```
ENTRY(system_call_handler)
    SAVE_ALL                // Save the current context
    cmpl $0, %eax           // Is syscall number negative?
    jl err                 // If it is, jump to return an error
    cmpl $MAX_SYSCALL, %eax // Is syscall greater than MAX_SYSCALL (4)?
    jg err                 // If it is, jump to return an error
    call *sys_call_table(, %eax, 0x04) // Call the corresponding service routine
    jmp fin                // Finish
err:
    movl $-ENOSYS, %eax     // Move to %eax the ENOSYS error
fin:
    movl %eax, 0x18(%esp)   // Change the %eax value in the stack
    RESTORE_ALL            // Restore the context
    iret
```

### 3.7.6 Service Routine to the write system call

In this section we will define the `sys_write` service routine, which checks that everything works correctly and shows the characters on the screen.

```
int sys_write(int fd, char * buffer, int size);
    fd: file descriptor. In this delivery it must always be 1.
    buffer: pointer to the bytes.
    size: number of bytes.
    return ' Negative number in case of error (specifying the kind of error) and
           the number of bytes written if OK.
```

System calls (in general) follow these steps:

- 1) **Check the user parameters:** fd, buffer and size. Bear in mind that the system has to be robust and assume that the parameters from the user space are unsafe by default (**lib.c is user code**).
  - a) Check the fd, we will use a new **int check\_fd (int fd, int operation)** function that checks whether the specified file descriptor and requested operation are correct. The operations can be `LECTURA` or `ESCRITURA`. If the operation indicated for this file descriptor is right, the function returns 0. Otherwise, it returns a negative identifier of the generated error.
  - b) Check buffer. In this case, it will be enough to verify that the pointer parameter is not `NULL`.

- c) Check size. Check positive values.
- 2) **Copy the data from/to the user address space if needed.** See the functions `copy_to_user` and `copy_from_user` (section 3.7.7).
- 3) **Implement the requested service. For the I/O system calls, this requires calling the device dependent routine.** In this particular case, the device dependent routine is `sys_write_console`:

```
int sys_write_console (char *buffer, int size);
```

This function displays *size* bytes contained in the *buffer* and returns the number of bytes written.

- 4) **Return the result.**

### 3.7.7 Copying data from/to the user address space

Copying data from/to the user is a critical operation because it could be a cause of kernel vulnerability. During the project, even it is possible to access the different address spaces because they are disjoint, you are asked to use a couple of operations (`copy_from_user` and `copy_to_user`) to explicitly mark the data transfers between both memory address spaces.

The Linux Kernel Module Programming Guide argues the need for these functions as follows:

*"The reason for `copy_from_user` or `get_user` is that Linux memory (on Intel architecture, it may be different under some other processors) is segmented. This means that a pointer, by itself, does not reference a unique location in memory, only a location in a memory segment, and you need to know which memory segment it is to be able to use it. There is one memory segment for the kernel, and one for each of the processes. The only memory segment accessible to a process is its own, so when writing regular programs to run as processes, there's no need to worry about segments. When you write a kernel module, normally you want to access the kernel memory segment, which is handled automatically by the system. However, when the content of a memory buffer needs to be passed between the currently running process and the kernel, the kernel function receives a pointer to the memory buffer which is in the process segment. The `copy_from_user` and `copy_to_user` functions allow you to access that memory".*

This means that `copy_to_user` and `copy_from_user` encapsulate complexity due to processor architectural differences. In our scenario, these functions will only be useful for copying data between the user and the OS address space.

## 3.8 Work to do

### 3.8.1 Complete Zeos Code

The released code lacks some features explained in the ZEOS document. You have to **complete** it:

- Implement the macro **RESTORE\_ALL**.
- Implement the macro **EOI**.
- Implement the keyboard management.
  - Implement the keyboard service routine.
  - Implement the keyboard handler.
  - Initialize the IDT with the keyboard handler
  - Enable the interrupt.



- Implement the *system\_call\_handler* routine.
- Initialize the IDT with the handler
- Implement the *write* system call.
  - Implement the *sys\_write* routine.
  - Modify the *sys\_call\_table* with the new routine.
  - Create a wrapper for the system call.
- Implement the *errno* and *perror* function.

### 3.8.2 Clock management

The clock management will display the seconds that have elapsed since the boot process. Figure 17 shows the expected result with the time displayed in a fixed place on the screen.

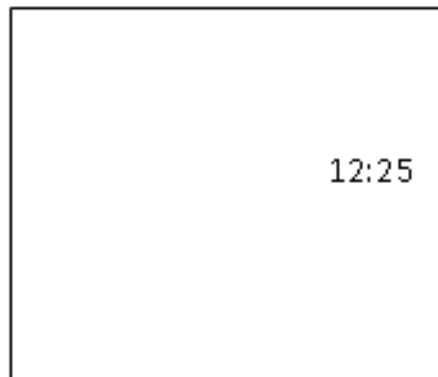


Fig. 17. Clock interrupt

To write the management of an interrupt such as the clock interrupt (that is, a masked type), you must:

- Initialize the entry of the clock interrupt in the IDT. The clock interrupt is contained in entry 32.
- Write the handler.
- Write the service routine.
- Enable the interrupt. The clock interrupt is a masked interrupt that is disabled in the code provided, which means that it is not dealt with. Modify the mask that enables the interrupts, located in the **enable\_int()** routine (file **hardware.c**).

Notice that just after enabling the interrupts, one of them can be raised at any time. When this occurs, you must ensure that the system is fully initialized because the interrupt service routines may access any system structure.

#### Writing the handler

An interrupt handler is written more or less the same way as an exception handler. Follow the steps below:

- 1) Define an assembler header for the handler.
- 2) Save the context.

- 3) Perform the EOI (remember that there is a macro for this available). You must notify the system that you are treating the interrupt.
- 4) Call the service routine.
- 5) Restore the context.
- 6) Return from the interrupt to user mode.

### *Writing the service routine*

Write the corresponding service routine of the clock interrupt. Inside this routine, you must call: **zeos\_show\_clock** which will display at the top-right part of the screen a clock with the elapsed time since the OS booted. Notice that in bochs the clock interrupt is also emulated. So, this showed time goes faster than the real one.

The header of this function is located at `zeos_interrupt.h`:

```
void zeos_show_clock();
```

The seconds that have elapsed since the OS boot process must appear on the screen. If it does not work, determine the problem with the debugger. Check whether the interrupt is executed, whether you return to the user mode, etc.

### *3.8.3 Gettime system call*

Extend ZeOS to incorporate this new system call. This syscall returns the number of clock ticks elapsed since the OS has booted. Its header is:

```
int gettime();
```

To implement this system call you will:

- 1) Create a global variable called `zeos_ticks`.
- 2) Initialize `zeos_ticks` to 0 at the beginning of the operating system (main).
- 3) Modify the clock interrupt service routine to increment the `zeos_ticks` variable.
- 4) Write the wrapper function. The identifier for this system call will be 10.
- 5) Update the system calls table
- 6) Write the service routine.
- 7) Return the result.

## 4 BASIC PROCESS MANAGEMENT

In this section we describe the code and necessary data structures to define and manage processes in ZeOS.

Take into account that the code related to the process management is key to achieve a good system performance. This means that the data structures have to take up minimum space and that the process management algorithms must be highly efficient.

As you know, a process has its own address space. A running process in user mode accesses its user stack, user data and user code. When running in system mode, it accesses the kernel data and code and uses its own system stack.

The process management implies knowledge of memory management. Here you will be given some basic concepts. A lot of work has already been done, but **you will have to understand the code provided to you**.

Chapters 3 and 11 of *Understanding the Linux Kernel* contain information about process management. You can also find useful information about the memory organization in Chapter 2.

This section is the most complex so you are advised to read the following sections carefully and to understand all the concepts explained.

### 4.1 Prior concepts

- Process. Data related to the process management: PCB (task\_struct), lists, etc.
- Process address space. Difference between logical and physical addresses. Paging.
- Context switch. Scheduling. Scheduling policies.

### 4.2 Definition of data structures

One of the goals of this section is to describe the data structures needed to manage processes in ZeOS. In particular:

- Process Control Block (PCB) implemented with the task\_struct structure in sched.h.
- A vector of task\_structs to hold the information of several processes.
- The system stack integrated into the task\_struct of every process.
- A free queue to link all the tasks that are available for new processes (to get a quick access to them).
- A ready queue to link all the tasks that are candidate to use the CPU which is already busy.

#### 4.2.1 Process data structures

##### PCB structure

The PCB should contain all necessary information needed to manage the processes in the system. A basic definition of the task\_struct is contained in the sched.h file:

```
struct task_struct {
    int PID;                                /* Process ID */
    page_table_entry * dir_pages_baseAddr; /* directory base address */
};
```

This structure has initially only two fields: an integer that will serve as the process identifier (PID), and a pointer to the directory pages of the processes (the directory pages points to the page table of the processes, see section 4.3 for more details). More fields will be added as you advance in the project.

For easy programming, rather than using two different structures: one for the `task_struct` and one for the stack for each process, both data structures will be overlapped in memory. To do so, we use the **union type of C**, which is a special kind of data. In C, if you declare a **union of three fields (a, b and c)**, the compiler reserves memory for the biggest field instead of the sum of their sizes.

To overlap the `task_struct` and the system stack, the `task_union` has been declared. The `task_union` shares the memory space between the process descriptor (`task_struct`) and the system stack of the process. Its declaration is:

```
union task_union {
    struct task_struct task;
    unsigned long stack[KERNEL_STACK_SIZE];
};
```

And you can see the result in Figure 18, where both structures have been depicted sharing the same memory area, a region of 4096bytes.

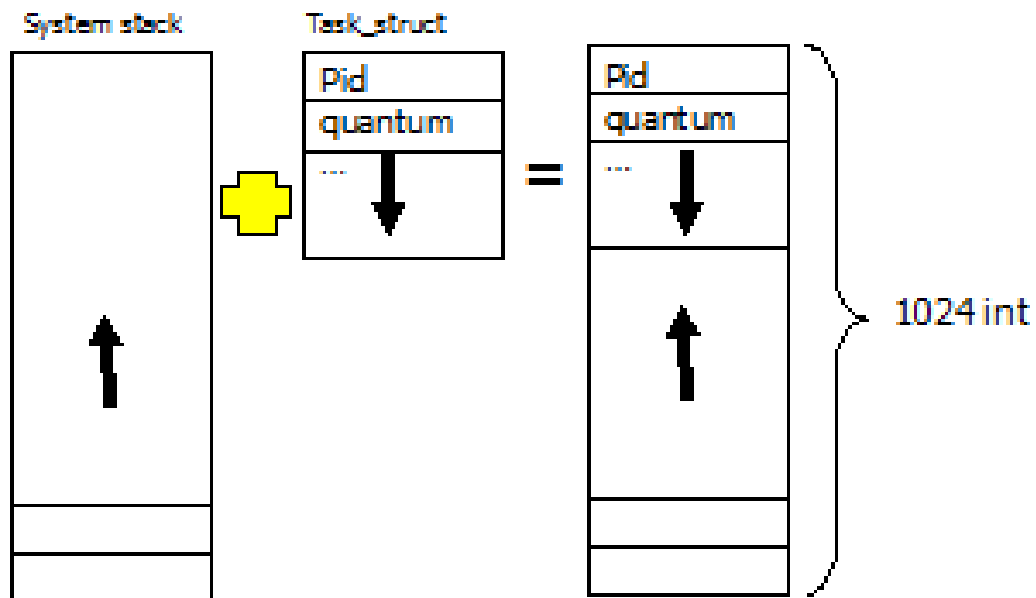


Fig. 18. The system stack and `task_struct` use the same page memory(4k)

The operating system will allow to create and destroy processes, meaning that a bunch of this structures will be needed. To store these processes in the system we will use a fixed size array<sup>15</sup>. The declaration of the task array with the `task_union` is also provided:

```
union task_union task[NR_TASKS]
```

15. A better option would be to create these `task_struct` structures dynamically, using some kind of dynamic memory allocation but we do not have this feature yet.

```
__attribute__((__section__(".data.task")));
```

The array is tagged with an attribute section to locate this array into the memory section named *data.task*. Figure 19 shows the placement of this task array in memory.

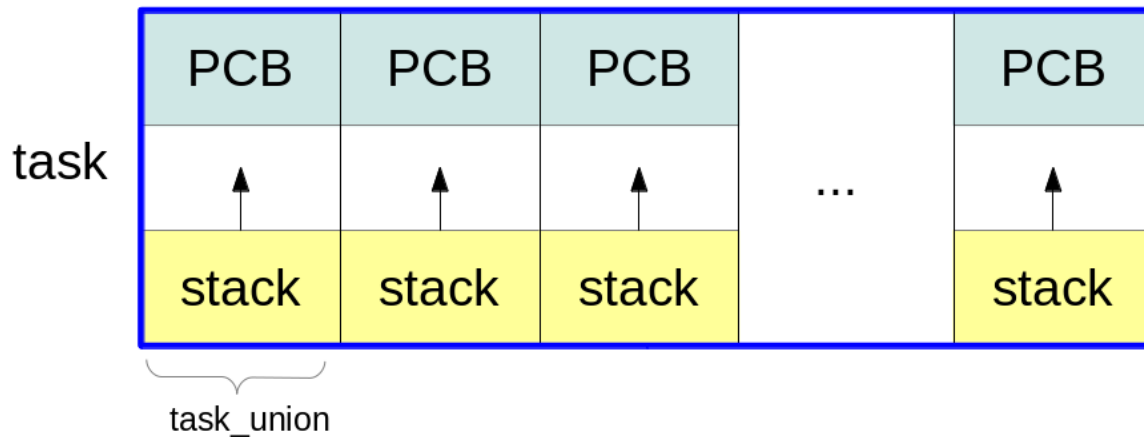


Fig. 19. Memory placement of the task array

The process allocated to entry 0 will be the idle process, referred to as process 0, idle process, or task 0. This process cannot be destroyed.

### Free queue

Finally, in order to get a quick access to the task\_structs not used in the task array a queue of the available PCBs is desired (freequeue). This queue must be initialized during the operating system initialization to link all the available tasks.

This queue is implemented using the `struct list_head` structure provided in the `list.h` file. In fact, this structure is the one used in Linux to implement whatever all the system queues.

The steps to include the free queue are:

- 1) **Declare and initialize** a variable named `freequeue` of type `struct list_head`.
- 2) Modify the declaration of the `task_struct` structure. Add one field named `list` that you will use to enqueue the structure into a queue.
- 3) Add all `task_structs` to this queue as, at this moment, all of them are available.

Due to the fact that the list functions are generic, they use a unique type: the `struct list_head`. So you will need a function to find the memory address of a `task_struct` given a particular `list_head` struct address. There is one function named `list_head_to_task_struct` which returns a pointer to a `task_struct` given a pointer to a `list_head` struct. Its header is:

```
struct task_struct *list_head_to_task_struct(struct list_head *l);
```

## Ready queue

Processes that are candidate to use the CPU but it is currently busy, remain queued in the ready queue. This queue holds the `task_struct` of all the process that are ready to start or continue the execution.

The steps to include the ready queue are:

- 1) **Declare and initialize** a variable named `readyqueue` of type `struct list_head`.

This queue must be initialized during system initialization and it is empty at the beginning.

### 4.2.2 Process identification

Once you have the process table, an important issue is how to identify a process. When entering the system (through a system call or a clock interrupt, for example) you need to know which process is on execution in order to access its data structures. **Where can you obtain this information? From the system stack pointer (remember that the `task_struct` of a process and its system stack share the same memory page).**

When entering the kernel, you know that:

- The system stack and the `task_struct` are physically in the same space
- The stack size is fixed (its size is 4KB).
- The `task_struct` is located at the beginning of the page that it occupies.
- The `esp` and `ss` registers are changed by the hardware automatically when changing from user mode to system mode.

Thus, if you apply a bitmask to the `esp` register, setting the latest `X` bits to 0 (where `X` is the number of bits used by the stack size), you will have the address of the beginning of the `task_struct`.

A function (or macro) named "current" that returns the address of the current `task_struct`, based on this idea, is provided. This macro gives the pointer to the task that is currently running.

```
struct task_struct * current()
```

## 4.3 Memory management

This section introduces some concepts and data structures related to memory management, necessary for process management.

In ZeOS, you are provided with the segmentation already initialized (and fixed) and the paging initialized for the initial process. All processes executing on ZeOS access the same logical pages and a logical-to-physical page translation is needed for each process.

In order to understand the creation of the processes, you need some knowledge about the memory organization in the 80x86 architectures.

Two types of ZeOS addresses should be differentiated: logical (or virtual) and physical.

- Logical address<sup>16</sup>. A 32-bit unsigned integer that represents up to 4GB of memory. The address range is from 0x00000000 to 0xffffffff. All the addresses generated by the cpu are logical addresses and translated to physical addresses by the MMU.
- Physical address. They are represented by 32-bit unsigned integers. Used to address memory from the chip.

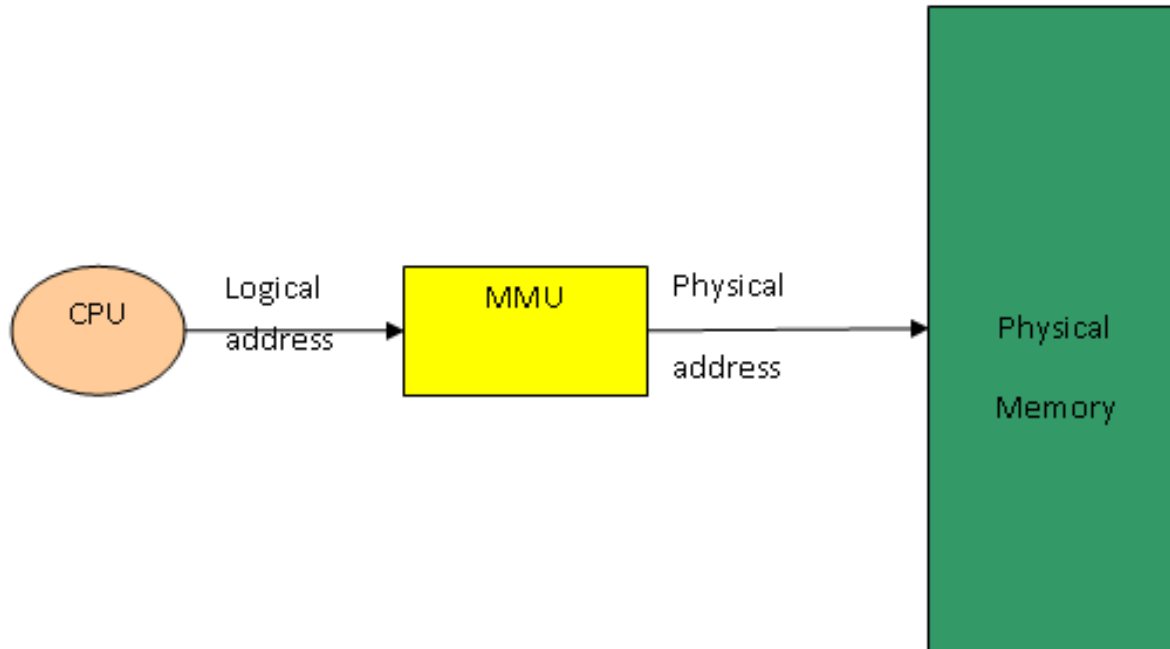


Fig. 20. Translation from logical addresses to physical addresses

Figure 20 shows the relationship between the different addresses. A brief description of segmentation and paging concepts will now be given.

#### 4.3.1 MMU: Paging mechanism

Paging translates logical addresses into physical addresses. The memory is arranged in page frames of equal size. The data structure that translates linear addresses into physical addresses is the page table.

To activate the paging mechanism in the 80x86 architectures, it is necessary **to activate the PG bit of the cr0 register**. PG=0 means that the linear addresses are the same as the physical addresses. This initialization is already done in the code provided in ZeOS.

#### 4.3.2 The directory and page table

The addressing mode in ZeOS has to use the two indexing levels offered by the 80x86 paging architecture: *page directory* and *page table*. Each page directory contains the base addresses of the page tables defined in the system for one process. The page tables contain the actual translations

16. A difference is established between logical and linear in "Understanding the Linux Kernel". What we refer to as "logical" here is called "linear" in the book

from logical to physical pages of the processes. `cr3` is a special register that keeps the physical base address<sup>17</sup> addresses for all kernel code and data, included page directories, are absolute. This means that logical addresses are the same than physical addresses of the page directory used by the MMU to do the current translations.

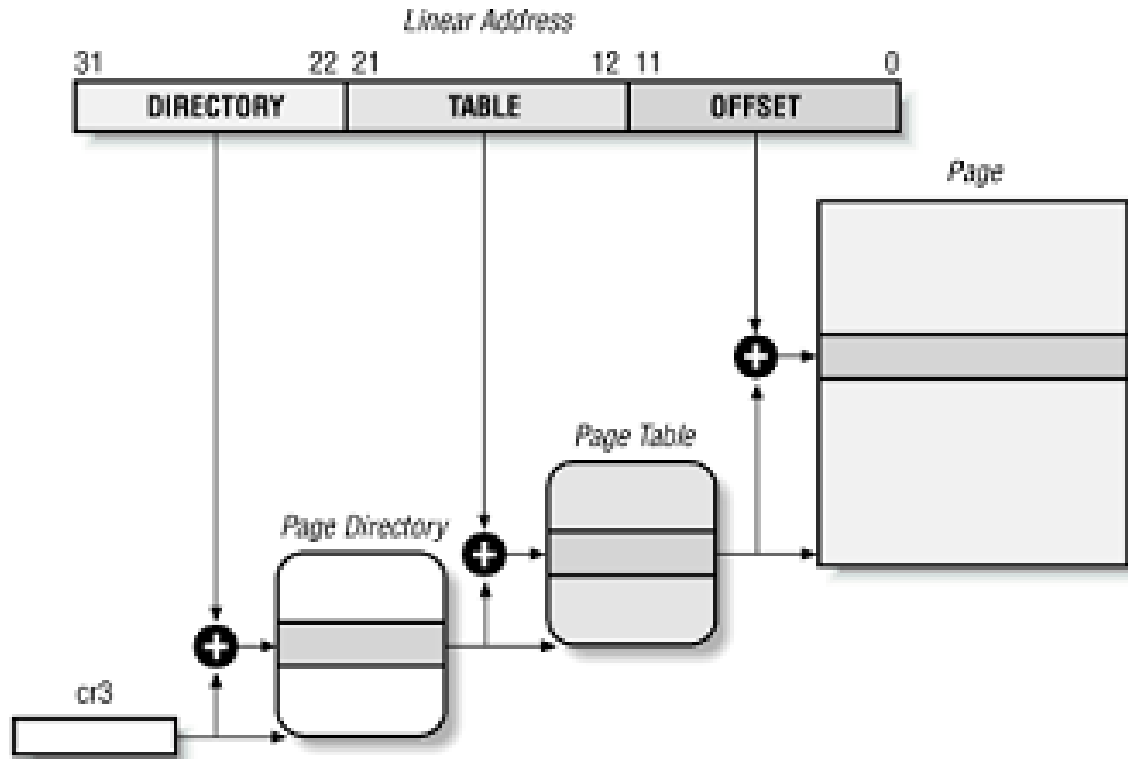


Fig. 21. Paging in the 80x86 architecture

Figure 21 shows the basic address translation mechanism based on paging used in ZeOS. In the Intel architecture, the 10 most significant bits (bits 31-22) of a 32-bit address define the entry to the *page directory*. The 10 bits in the middle (bits 21-12) define the *logical page* number and indexes the page table. The 12 least significant bits (bits 11-0) define the *offset* within the page (up to 4KB). Therefore, to obtain the logical page number for a given address, it is only necessary to shift it to the right:  $\text{page\_number} = (\text{address} \gg 12)$ .

The purpose of the two indexing levels mechanism is to allocate space for the page tables as it is needed, this is as the address space of a process grows. In ZeOS, the address space of all processes is small and only one page table per process is needed. Therefore each page directory has only one valid entry (entry 0).

Each process has its own page directory and page table, and the operating system is the responsible to load the `cr3` register with the base address of the page directory for the process that is going to use the cpu. For this reason, ZeOS keeps the base address of the directory pages of each process in its PCB (field `dir_pages_baseAddr` in the `task_struct` structure).



Variables *dir\_pages* and *pagusr\_table* are two fixed sized vectors that hold the page directories and page tables respectively for each process in ZeOS (see file mm.c). These vectors have NR\_TASKS entries, and each of these entries is linked to a *task\_struct* in the *tasks* vector, through the *dir\_pages\_baseAddr* field.

The initialization of all the page directories is given as part of the ZeOS code (see function *init\_dir\_pages* in file mm.c). This initialization consists on setting the entry 0 of each page directory to point to its corresponding page table.

The page directory and the page table entries have the same structure. Each entry has the following flags (check the ZeOS code):

```
typedef union
{
    unsigned int entry;
    struct {
        unsigned int present : 1;
        unsigned int rw      : 1;
        unsigned int user    : 1;
        unsigned int write_t : 1;
        unsigned int cache_d : 1;
        unsigned int accessed : 1;
        unsigned int dirty   : 1;
        unsigned int ps_pat  : 1;
        unsigned int global  : 1;
        unsigned int avail   : 3;
        unsigned int pbase_addr : 20;
    } bits;
} page_table_entry;
```

Each entry has been defined as an union to initialize it easily. Setting the field *entry* to 0, sets the whole structure to 0. A new routine to enable this operation has been defined:

```
void del_ss_pag(page_table_entry* PT, unsigned logical_page);
```

The *pbase\_addr* field contains the page frame number that will be used for the current logical page. Therefore, if this structure is an entry of a page table<sup>18</sup>, it will translate from a **logical page number to a physical page number**. Remember that the content of the *pbase\_addr* field is the **physical page number, not the address**. To change the association between a logical page and a physical frame, you can use this function (mm.c):

```
void set_ss_pag(page_table_entry* PT, unsigned logical_page, unsigned
physical_frame);
```

### 4.3.3 Translation Lookahead Buffer (TLB)

The TLB is a table that the architecture uses to optimize the memory access. The page table entries used by the current process are stored in this table. The TLB is like a cache of the page table entries. This helps us to save one access to the directory when translating the linear addresses

18. instead of an entry of a page directory

into physical ones. Moreover, the TLB access is very fast. **Bear in mind that the TLB is not synchronized automatically with the page table.** Therefore, when the page table is modified, the associated TLB entries can become incorrect, so these TLB entries must be invalidated. In particular, every time a page table entry is deleted or modified, it is necessary to invalidate the TLB. This action is called TLB flush. In the 80x86 paging architecture, the hardware flushes the TLB each time the value of the cr3 register changes. **To invalidate the TLB, rewrite the cr3 register using the routine `set_cr3(page_table_entry* dir)` (see file `mm.c`).**

#### 4.3.4 Specific organization for ZeOS

##### *Logical space*

The logical space of a process is the memory space addressable for user programs. It is composed of segments and pages. The ZeOS segmentation has been reduced to a minimum, like in Linux. If you look at the `bootsect.S` file, you will see how the GDT entries have been initialized and that the segments are defined to start in the address 0 with an infinite size. This means that a segment:offset address will be translated basically to a linear address containing the same value of the offset. To make everything work correctly and due to the lack of a loader, the linker will generate addresses for user programs starting at the `0x100000` (`L_USER_START` in `mm_address.h`). In order to have a global understanding of memory management, it is necessary to understand the organization of segments. However, it is not necessary to modify the one provided.

##### *Logical space in ZeOS*

The linear address space is composed of N code pages and M data+stack pages, both consecutive in memory.

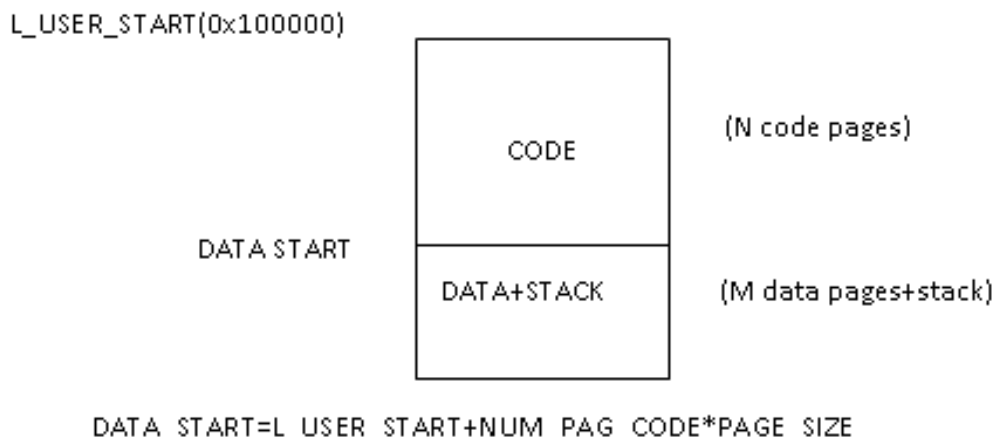


Fig. 22. Linear address space

Figure 22 shows the user addresses of a process. The left side of the figure shows the logical addresses of any process. The logical space of all processes starts at the address `L_USER_START` and has `NUM_PAGE_CODE` code pages and `NUM_PAG_DATA` pages of data.

In the case of ZeOS, the first 255 entries of the page table (addresses `0x0` to `0xfffff`) correspond to the system space (see the initialization code in file `mm.c`). Addresses of the system space are absolute, i.e. they do not require translation because logical addresses match physical address.

Each entry represents 4KB (it is the page size).  $L\_USER\_START \gg 12$  is the first logical page of the user code (you must define a constant to access this page in the file `mm_address.h`). Therefore, the entire process has  $255 + NUM\_PAGE\_CODE + NUM\_PAGE\_DATA$  valid entries in the page table: 0-255 for the system, only accessible when in kernel mode,  $NUM\_PAGE\_CODE$  pages for the user code and  $NUM\_PAGE\_DATA$  pages for the data and user stack. In the initial ZeOS configuration,  $NUM\_PAGE\_CODE=8$  and  $NUM\_PAGE\_DATA=20$  (see file `mm_address.h`).

Both logical and physical pages for system and user code are shared. This means that the logical entries in the page table of all processes point to the same physical pages.

**The pages for user data+stack are private and therefore the translation from logical addresses to physical addresses (or frame) will be different for each process.**

For each new process, it is necessary to initialize its whole page table with the suitable translation for each logical page.

As few pages will be used, a single entry of each page directory will be used (`ENTRY_DIR_PAGES`). This entry has the address of the corresponding page table. And register `cr3` points to the base address of the current page directory.

For each context switch it is necessary to change the value of the `cr3` register to point to the process that it is going to use the cpu.

Remember that the hardware uses the information in the page table to translate all the accesses to memory: it checks that the type of access is enable for the target logical address and translate the address to complete the access. If the hardware is not able to complete the access it generates an exception.

In order to avoid page faults inside the operating system, ZeOS has to check that all the addresses passed by the user as parameters are valid. For this purpose, a function named *access\_ok* has been provided. This function, given a user memory address, a block size and a type of access, checks in the page table of the process if the access is valid.

```
int access_ok(int type, const void *addr, unsigned long size)
type READ or WRITE
addr user address
size block size to check
returns 1 if correct and 0 otherwise
```

### *Physical space*

As mentioned above, **the logical space is consecutive but the physical space is not**. The physical memory will be organized as follows: the code pages will be shared by all the processes and they will start at `PH_USER_START` (defined in `segment.h`). Starting at this address, the system will place memory pages for user data and user stacks of created processes (see Figure 23).

Notice that the only difference between the various processes will be the entries of the page table corresponding to the user data and user stack.

The physical pages or frames are computed from physical addresses. For example, the physical pages for the code of process `P0` will start at the frame  $PH\_USER\_START \gg 12$ . A macro name `PH_PAGE` to compute the page number starting from a memory address has been provided:

```
#define PH_PAGE(x) (x >> 12)
```

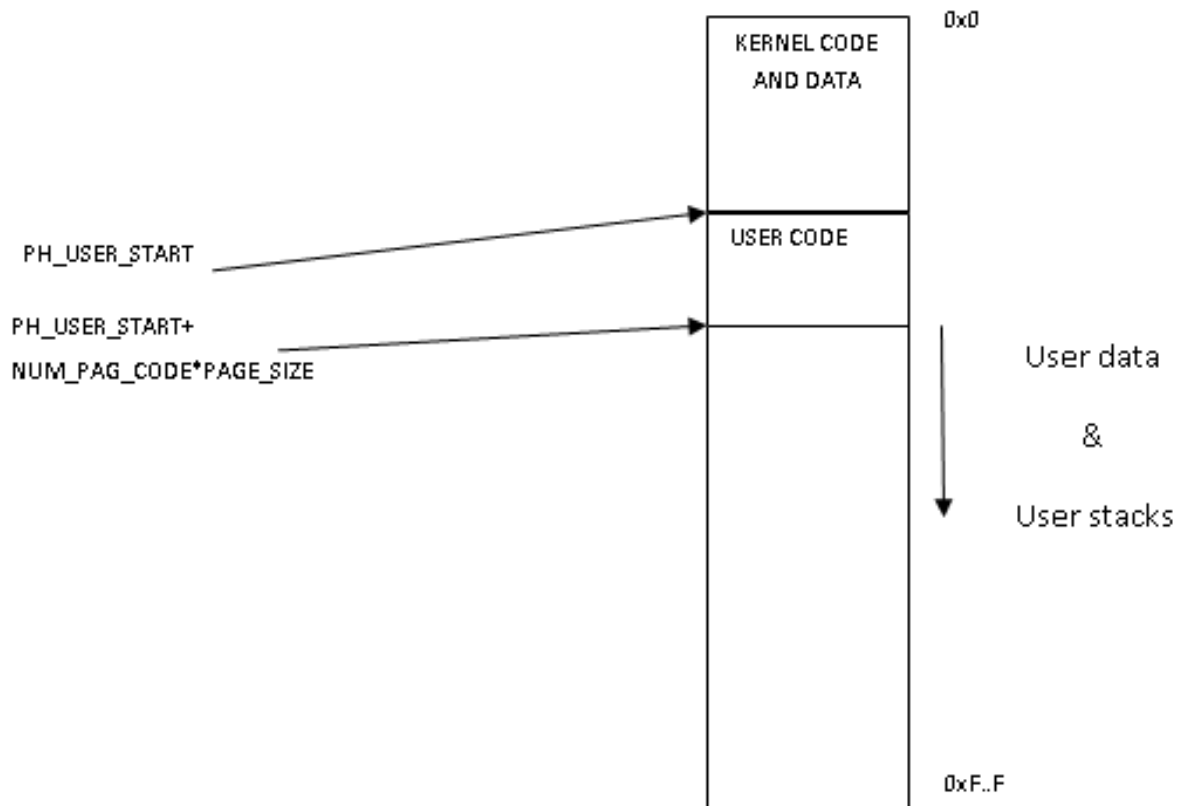


Fig. 23. Physical memory in Zeos

Every time a process is created, it will be necessary to look for available frames in the physical memory for mapping data and stack to the process. And every time a process ends an execution, it will mark as available the data and stack frames that were assigned to the process.

The `phys_mem` table (`mm.h` file) has been defined to have information in the system about busy frames (assigned to a logical page of a process) and free frames. It contains an entry for each frame with its status. The structure is initialized as a part of the `init_mm` routine (in `mm.c`) and it is based on marking all system frames and the frames used by the user code as busy (see Figure 24).

To complete the `phys_mem` management interface, you are provided with the `alloc_frame` function (to allocate one frame) and the `free_frame` (to deallocate one frame) functions in the `mm.c` file. The interface of these two functions is described below.

```
int alloc_frame(void)
```

The `alloc_frame` function searches one available frame. If it finds one, then it is marked as busy in the `phys_mem` table and returned. If there are no free frames, it returns -1.

```
void free_frame(unsigned int frame)
```

The `free_frame` function marks the frame passed as a parameter as free.

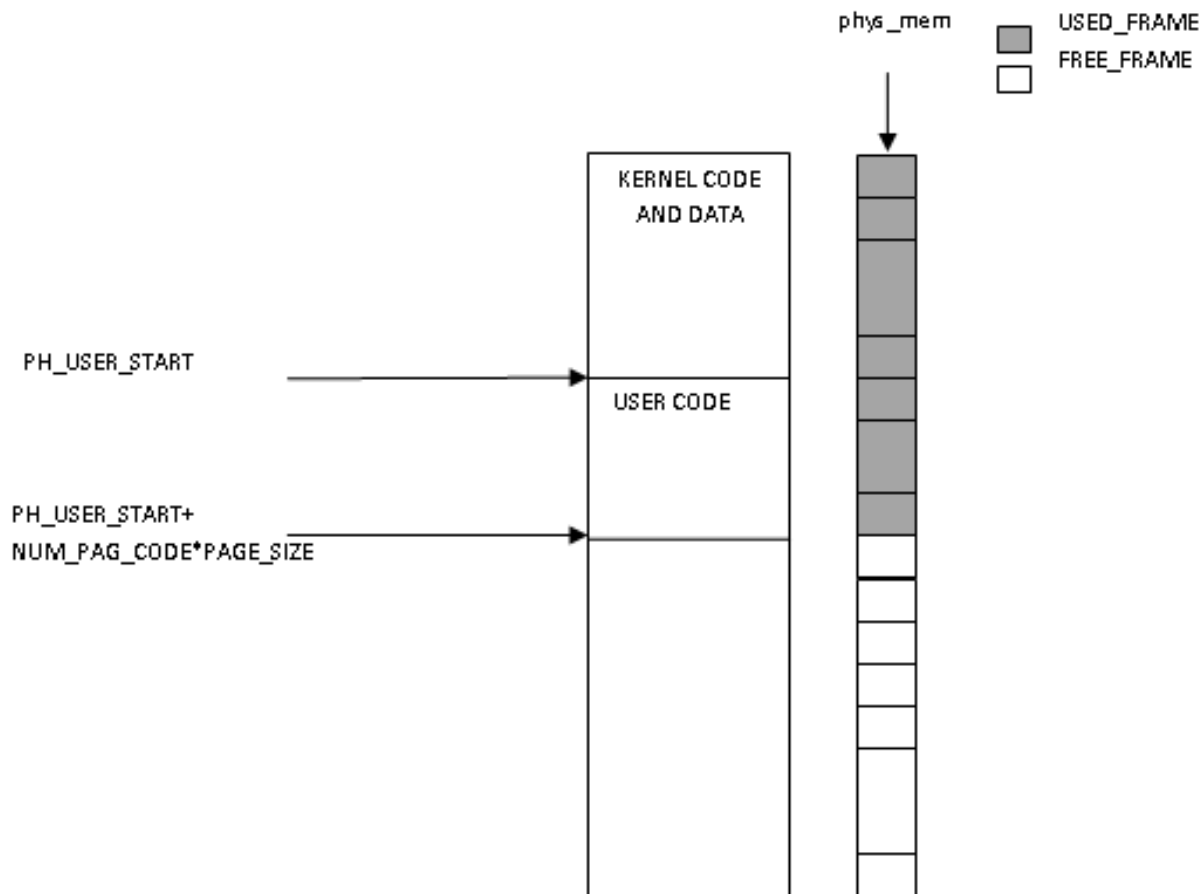


Fig. 24. Initializing physical memory and `phys_mem` table

## 4.4 Initial processes initialization

There are two special processes in ZeOS: the idle process and the init process. Both processes are created and initialized during the initialization code of the system (see file `system.c`).

### 4.4.1 Idle process

The **idle process** is a kernel process that never stops and that always execute the funtion `cpu_idle`, which executes the assembler instruction `sti` and implements an empty endless loop (see file `sched.c`).

The assembler instruction `sti` is required because in ZeOS, interrupts are disabled while executing in kernel mode (disabling interrupts is the default behavior when an interrupt happens in the 80x86 architecture). As the idle process executes in kernel mode, it is necessary to explicitly enable interrupts when the idle process starts using the `cpu`, in order to be able to interrupt it and to give the `cpu` to any user process that becomes ready to use it.

The PID of this process is 0, always executes on kernel mode and only should use the CPU if there are not any other user process to execute (its purpose is just to keep the CPU always busy).

The main steps to initialize the idle process have to be implemented in the `init_idle` function and are the following:

- 1) Get an available `task_union` from the freequeue to contain the characteristics of this process
- 2) Assign PID 0 to the process.
- 3) Initialize field `dir_pages_baseAaddr` with a new directory to store the process address space using the `allocate_DIR` routine.
- 4) Initialize an execution context for the process to restore it when it gets assigned the cpu (see section 4.5) and executes `cpu_idle`.
- 5) Define a global variable `idle_task`

```
struct task_struct * idle_task;
```

- 6) Initialize the global variable `idle_task`, which will help to get easily the `task_struct` of the idle process.

The first time that the idle process will use the cpu will be due to a cpu scheduling decision, when there are not more candidates to use the cpu. Thus, the context switch routine will be in charge of putting this process on execution. This means that we need to initialize the context of the idle process as the context switch routine requires to restore the context of any process. This context switch routine restores the execution of any process in the same state that it had before invoking it (see subsection 4.5 for more details). This is achieved by undoing the dynamic link in the stack and then return using the code address in top of the stack (see figure 25). This means that we need to:

- 1) **Store in the stack of the idle process the address of the code that it will execute (address of the `cpu_idle` function).**
- 2) **Store in the stack the initial value that we want to assign to register `ebp` when undoing the dynamic link** (it can be 0),
- 3) Finally, we need to **keep (in a field of its `task_struct`) the position of the stack where we have stored the initial value for the `ebp` register**. This value will be loaded in the `esp` register when undoing the dynamic link.

Notice that as part of this initialization it is not necessary to modify the page table of this process. The reason is that this process is a kernel process and only needs to use those pages that contain kernel code and kernel data which are already initialized during the general memory initialization in the `init_mm` function (see file `mm.c`). This function allocates all the physical pages required for the kernel pages (both code and data pages) and initializes all the page tables in the system with the translation for these kernel pages, which are common for all the processes (see functions `init_frames` and `init_table_pages` in file `mm.c`).

#### 4.4.2 Init process

The **init process** is the first one executed after booting the OS. Its PID is 1 and it is a user process, executing the user code. Since the init process is the first one executed, it becomes the parent for the rest of processes in the system. The code for this process is implemented in the `user.c` file.

The steps to initialize the init process have to be implemented in the `init_task1` function which is called from the function `main` in `system.c`. Those steps are:

- 1) Assign PID 1 to the process
- 2) Initialize field `dir_pages_baseAaddr` with a new directory to store the process address space using the `allocate_DIR` routine.

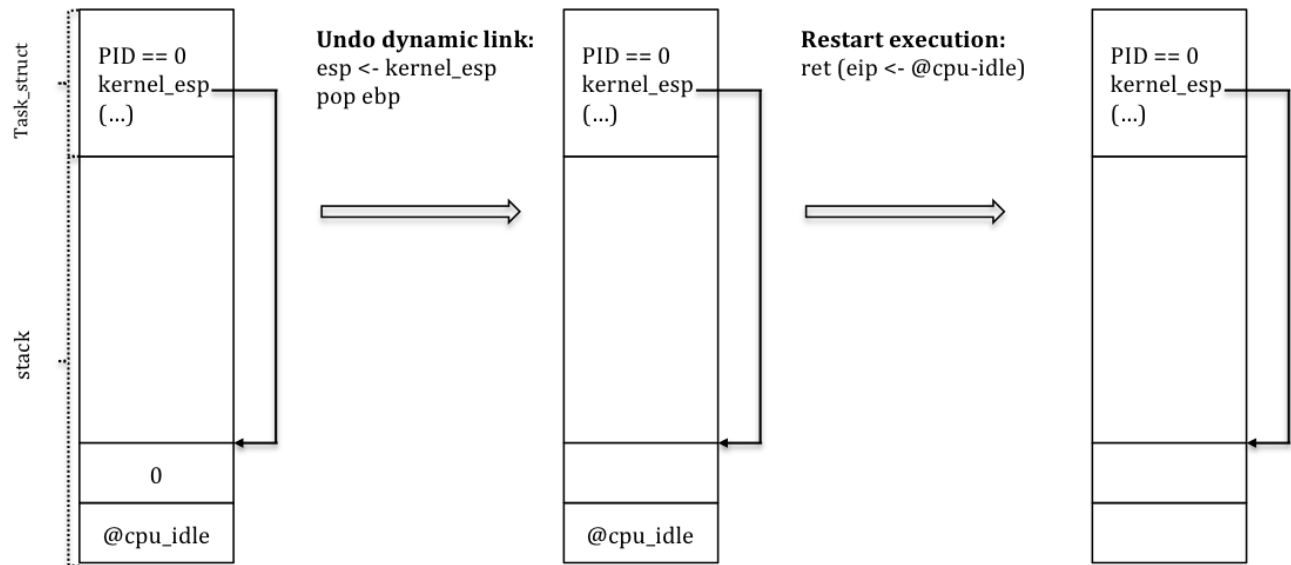


Fig. 25. Idle process: putting it on execution

- 3) Complete the initialization of its address space, by using the function `set_user_pages` (see file `mm.c`). This function allocates physical pages to hold the user address space (both code and data pages) and adds to the page table the logical-to-physical translation for these pages. Remind that the region that supports the kernel address space is already configure for all the possible processes by the function `init_mm`.
- 4) Update the TSS to make it point to the `new_task` system stack.
- 5) Set its page directory as the current page directory in the system, by using the `set_cr3` function (see file `mm.c`).

This process is the first process that will use the cpu after the system initialization. This means that the context switch routine is not involved with the first execution of this process. It is the function `return_gate` which prepares the stack of the process to enable the execution of the user code (see file `hardware.c`). At the end of `return_gate`, the `lret` assembler instruction is executed to downgrade the privilege level to user level and to execute the main function in the `user.c` file.

#### 4.4.3 Zeos implementation details

As part of the initialization code it is necessary to initialize both processes: the idle and the init. In file `sched.c` you have defined two empty functions for this purpose: `init_idle` and `init_task1`. These functions are invoked from the system `main` function (see file `system.c`), and you have to complete them with the necessary steps to initialize both processes. In addition, **you have to delete the sentence that invoke the function `monoprocess_init_addr_space`**. This function initializes the address space of the monoprocess version of ZeOS and it is not needed for the multitasking version that you are building now.

### 4.5 Process switch

Zeos provides a routine for making the context switch between two processes. This context switch is performed in kernel mode.

The context switch consists of exchanging the process that is being executed with another process. It changes the user address space (changing the page table), changes the kernel stack (to restore the new context), and restores the new hardware context. This routine will be named *task\_switch*.

Its header is:

```
void task_switch(union task_union *new)
new: pointer to the task_union of the process that will be executed
```

This routine 1) saves the registers ESI, EDI and EBX<sup>19</sup>, 2) calls the *inner\_task\_switch* routine, and 3) restores the previously saved registers.

The *inner\_task\_switch* routine has the same header as the previous *task\_switch*, and it restores the execution of the *new* process in the same state that it had before invoking this routine. This is achieved by changing to the new kernel stack, undoing the dynamic link in the stack (usually created by the C compiler when a function is called) and then continue the execution at the code address in the top of the stack. Figure 26 shows the stages of the stack of the current process during the invocation of the context switch. The last picture is the expected state of any process stack to be restored.

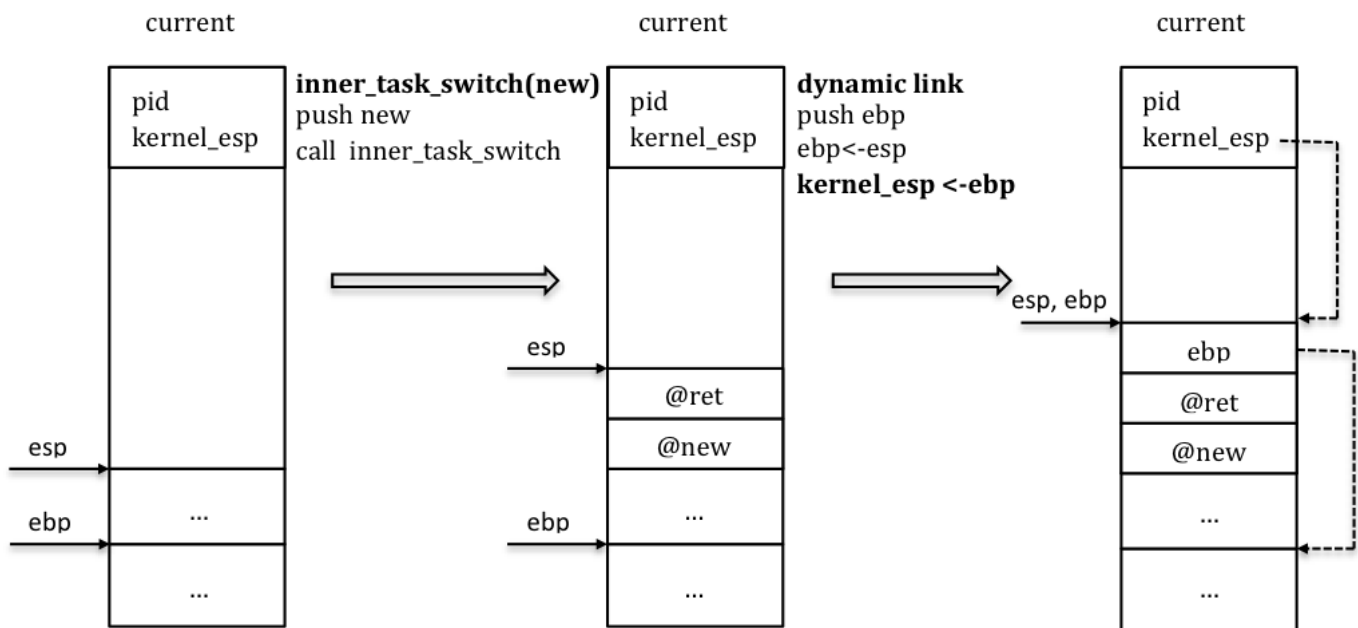


Fig. 26. Stack stages during the context switch invocation

The following steps are needed to perform a context switch:

- 1) Update the TSS to make it point to the `new_task` system stack.
- 2) Change the user address space by updating the current page directory: use the `set_cr3` function to set the `cr3` register to point to the page directory of the `new_task`.

19. Usually the compiler is smart enough to detect that these registers are modified inside the function code and it saves them automatically, but, in this case, the compiler is not aware of the whole code of the function (due to the trickeries that you will see soon) and potentially these registers may be modified, meaning that these registers would have a different value when returning from this function. The easy solution is to save and restore them manually.



- 3) Store the current value of the EBP register in the PCB. EBP has the address of the current system stack where the `inner_task_switch` routine begins (the dynamic link).
- 4) Change the current system stack by setting ESP register to point to the stored value in the *new* PCB.
- 5) Restore the EBP register from the stack.
- 6) Return to the routine that called this one using the instruction `RET` (usually `task_switch`, but...).

The context switch is used when the scheduling policy decides that a new process must be executed, or when a process blocks.

#### 4.5.1 Initial testing

You may use the idle and init processes to check that it works (for example, executing the init process at the beginning and after some ticks make a `task_switch` to the idle).

## 4.6 Process creation and identification

In this part, two system calls related to process management will be described in ZeOS: **getpid** and **fork**. `getpid` returns the pid of the process that calls it. `fork` creates a new process, a copy of the one that calls it (referred to as parent), inserts it into the processes list, ready to execute, and returns its pid.

#### 4.6.1 getpid system call implementation

The header of the call is:

```
int getpid(void)
```

It must return the pid of the process that called it. The identification of the `getpid` call is 20.

Since this is a very simple system call, its implementation is straightforward.

The wrapper routine in `libc.c` moves to `eax` the identifier of the `getpid` syscall (20), traps the kernel and returns the corresponding PID also in the `eax` register. The service routine returns the PID in the `task_struct` of the current process.

#### 4.6.2 fork system call implementation

The header of the call is:

```
int fork(void)
```

It returns a different value depending on whether it is the parent or the child. If it is the parent, it returns the pid of the created process; if it is the child, it returns a 0. If an error occurs it returns a -1.

Regarding the memory initialization, it is necessary to allocate free physical pages to hold the user data+stack region of the child process: system data and code are shared by all processes and since user code is read-only it is not necessary to replicate it. The data+stack region is

inherited from the parent and is private thereafter. This means that the new allocated pages will be initialized as a copy of the data+stack region of the parent.

The implementation of the fork system calls carries out the following steps:

- 1) Implement the wrapper of the system call. The identifier of the system call is 2.
- 2) Implement the `sys_fork` service routine. It must:
  - a) Get a free `task_struct` for the process. If there is no space for a new process, an error will be returned.
  - b) Inherit system data: copy the parent's `task_union` to the child. Determine whether it is necessary to modify the page table of the parent to access the child's system data. The `copy_data` function can be used to copy.
  - c) Initialize field `dir_pages_baseAddr` with a new directory to store the process address space using the `allocate_DIR` routine.
  - d) Search physical pages in which to map logical pages for data+stack of the child process (using the `alloc_frames` function). If there is no enough free pages, an error will be return.
  - e) Inherit user data:
    - i) Create new address space: Access page table of the child process through the directory field in the `task_struct` to initialize it (`get_PT` routine can be used):
      - A) Page table entries for the system code and data and for the user code can be a copy of the page table entries of the parent process (they will be shared)
      - B) Page table entries for the user data+stack have to point to new allocated pages which hold this region
    - ii) Copy the user data+stack pages from the parent process to the child process. The child's physical pages cannot be directly accessed because they are not mapped in the parent's page table. In addition, they cannot be mapped directly because the logical parent process pages are the same. They must therefore be mapped in new entries of the page table temporally (only for the copy). Thus, both pages can be accessed simultaneously as follows:
      - A) Use temporal free entries on the page table of the parent. Use the `set_ss_pag` and `del_ss_pag` functions.
      - B) Copy data+stack pages.
      - C) Free temporal entries in the page table and flush the TLB to really disable the parent process to access the child pages.
  - f) Assign a new PID to the process. The PID must be different from its position in the `task_array` table.
  - g) Initialize the fields of the `task_struct` that are not common to the child.
    - i) Think about the register or registers that will not be common in the returning of the child process and modify its content in the system stack so that each one receive its values when the context is restored.
  - h) Prepare the child stack with a content that emulates the result of a call to `task_switch`. This process will be executed at some point by issuing a call to `task_switch`. Therefore the child stack must have the same content as the `task_switch` expects to find, so it will be able to restore its context in a known position. The stack of this new process must be forged so it can be restored at some point in the future by a `task_switch`. In fact this new process has to a) restore its hardware context and b) continue the execution of the user process, so you must create a routine `ret_from_fork` which does exactly this. And use it as the restore point like in the idle process initialization 4.4.

- i) Insert the new process into the ready list: `readyqueue`. This list will contain all processes that are ready to execute but there is no processor to run them.
- j) Return the pid of the child process.

## 4.7 Process scheduling

In order to execute the different processes queued in the ready queue, a policy must be defined to select the next process. In this section, you will add a process scheduler that uses a round robin policy. However, the code of this scheduler must be generic enough to replace the round robin policy with a different one. For this reason we define an interface for the main scheduling functions.

### 4.7.1 Main scheduling functions

Following we describe the interface that you have to implement to provide ZeOS with a scheduling policy.

- Function to update the relevant information to take scheduling decisions. In the case of the round robin policy it should update the number of ticks that the process has executed since it got assigned the cpu.

```
void update_sched_data_rr (void)
```

- Function to decide if it is necessary to change the current process.

```
int needs_sched_rr (void)
returns: 1 if it is necessary to change the current process and 0
otherwise
```

- Function to update the state of a process. If the current state of the process is not *running*, then this function deletes the process from its current queue. If the new state of the process is not *running*, then this function inserts the process into a suitable queue (for example, the free queue or the ready queue). The parameters of this function are the *task\_struct* of the process and the queue according to the new state of the process. If the new state of the process is running, then the queue parameter should be NULL.

```
void update_process_state_rr (struct task_struct *t, struct list_head *dst_queue)
```

- Function to select the next process to execute, to extract it from the ready queue and to invoke the context switch process. This function should always be executed after updating the state of the current process (after calling function *update\_process\_state\_rr*).

```
void sched_next_rr (void)
```

You have to:

- 1) Add the necessary fields to the *task\_struct* to implement the scheduler.
- 2) Create a Round Robin scheduling policy. To do this, you should implement the previous interface and you should add the code necessary to invoke those functions.

### 4.7.2 Round robin policy implementation

The round robin policy consists in executing each process during a specific number of ticks (its *quantum*), doing a `task_switch` when the quantum finishes. This quantum can be different for each process and each process inherits it from its parent. Ticks are updated in the clock interrupt; the scheduling policy will then be invoked there to check if a new process must be executed.

Adding this policy requires implementing a new function, named *schedule*, and calling it whenever a scheduling decision is needed. The steps that this function has to perform are:

- Control how many ticks the current process has spent on the CPU. Use a global variable that is decreased with every tick.
- Decide if a context switch is required. A context switch is required when the quantum is over and there are some candidate process to use the CPU. Recall that the idle process should use the CPU **ONLY** if there are not any other user process that can advance with the execution.
- If a context switch is required:
  - Update the readyqueue, if current process is not the idle process, by inserting the current process at the end of the readyqueue.
  - Extract the first process of the readyqueue, which will become the current process;
  - And perform a context switch to this selected process. Remember that when a process returns to the execution stage after a context switch, its quantum ticks must be restored.

Additionally, to complete the implementation it is necessary to add a function named *get\_quantum* that returns the quantum of the task passed as a parameter and a function *set\_quantum* that modifies the quantum of the process.

```
int get_quantum (struct task_struct *t);
void set_quantum (struct task_struct *t, int new_quantum):
```

Finally, check that the round robin policy works correctly by showing the pid of the running process.

## 4.8 Process destruction

This section describes the steps to destroy a process. This destruction is implemented in the *exit* system call, which destroys the process that calls it. Therefore, it must delete the process from the processes table and make a context switch.

Its header is:

```
void exit(void)
```

The steps to implement the *exit* system call are:

- 1) Implement the wrapper of the system call. Its identifier is 1.
- 2) Implement the service routine `sys_exit`. In this case, this system call does not return any value.
  - a) Free the data structures and resources of this process (physical memory, `task_struct`, and so). It uses the `free_frame` function to free physical pages.
  - b) Use the scheduler interface to select a new process to be executed and make a context switch.

## 4.9 Statistical information of processes

The objective of this section is to add statistical information about processes in ZeOS. In particular, for each process, it keeps information about the time spent in user mode, the time spent in the ready queue, and the time spent in the blocked state<sup>20</sup>. In addition it also keeps information about the number of context switches involving the process and the amount of ticks remaining in the current quantum of the process.

The stats.h file contains a data structure to be used in this new functionality. Remember that the definition of this structure will be used both by the user and by the kernel.

```
struct stats {
    unsigned long user_ticks; /* Total ticks executed by the process */
    unsigned long system_ticks; /* Total ticks executing system code */
    unsigned long blocked_ticks; /* Total ticks in the blocked state */
    unsigned long ready_ticks; /* Total ticks in the ready state */
    unsigned long elapsed_total_ticks; /* Ticks since the power on of the machine
    until the beginning of the current state */
    unsigned long total_trans; /* Total transitions ready --> run */
    unsigned long remaining_ticks; /* Remaining ticks to end the quantum */
};
```

Modify the task\_struct as required to keep the statistical data, implement the code to keep updated this data for each process and implement the new system call to access it.

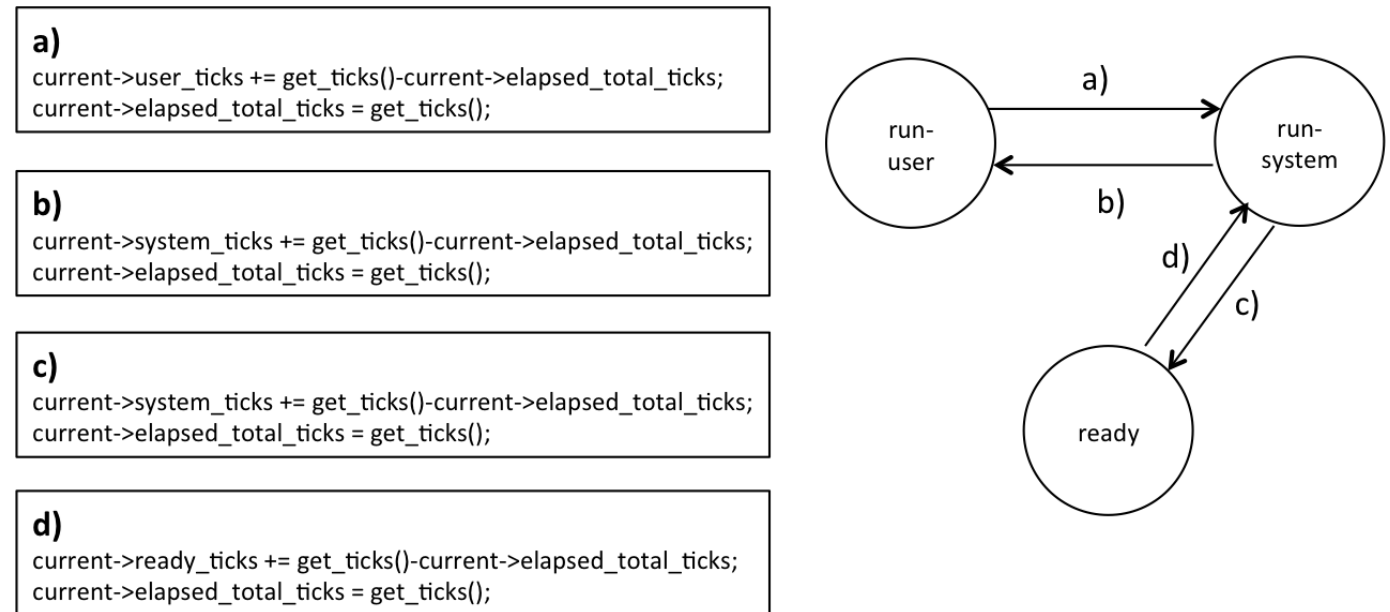


Fig. 27. States to consider when updating the execution time of a process

The already implemented function **get\_ticks** provides the elapsed ticks since the power on of the machine, and it can be used to calculate the different times. It has the following interface:

```
unsigned long get_ticks(void)
```

20. Notice that until now we have not implemented any blocking system call and thus, this time will be 0 for all the processes.

You must use this function each time a process changes the state both: 1) to update the amount of time spent in the state that the process is leaving and 2) to register into the field **elapsed\_total\_ticks** the initial time for the next state. For example, each time a process enters the system you have to update the **user\_ticks** field and to start counting the system time by performing the following operations:

```
current_ticks = get_ticks();
current->user_ticks += current_ticks - current-> elapsed_total_ticks;
current->elapsed_total_ticks = current_ticks;
```

Figure 27 represents the states and transitions that we are considering at this stage of the project, and the code involved to update the statistical information about the execution time. Think carefully about where to insert each of these lines of code inside the ZeOS implementation.

The identifier of the `get_stats` system call is 35. This system call returns -1 if an error is produced, 0 otherwise.

```
int get_stats(int pid, struct stats *st)
pid: process identifier whose statistics get_stats will get
st: pointer to the user buffer to store the statistics
```

The `pid` parameter contains the identifier of an alive process. The statistics of the process will be returned through the `st` parameter. Be aware that `st` is a pointer to the user address space. Think about error cases and how to deal with them.

## 4.10 Work to do

This section describes the steps to follow to make ZeOS become multitasking. You will complete the ZeOS code to create/identify/switch/destroy processes and add a new process scheduler:

- Adapt the `task_struct` definition.
- Initialize a free queue.
- Initialize a ready queue.
- Implement the initial processes initialization.
- Implement the `task_switch` function.
- Implement the `inner_task_switch` function.
- Implement the `getpid` system call.
- Implement the `fork` system call.
- Implement process scheduling.
- Implement the `exit` system call.
- Implement the `getstats` system call.

# 5 LIGHTWEIGHT THREADS AND PROCESS SYNCHRONIZATION

## 5.1 Work to do

In this delivery you will:

- Implement lightweight processes (threads)

- Adapt *allocate\_DIR*.
- Implement the *clone* system call.
- Implement semaphores to synchronize processes
  - Implement the *sem\_init* system call.
  - Implement the *sem\_wait* system call.
  - Implement the *sem\_signal* system call.
  - Implement the *sem\_destroy* system call.

### 5.1.1 Lightweight processes: threads

In this part, you have to implement lightweight processes: *threads*. You have to add a new *clone* system call with the following interface:

```
int clone (void (*function)(void), void *stack)
    function: starting address of the function to be executed by the new process
    stack    : starting address of a memory region to be used as a stack
returns: -1 if error or the pid of the new lightweight process ID if OK
```

It creates a new process that executes the code at *function*. The new process shares the whole address space with the current one, and *stack* is the base address of the memory region to be used as its user stack.

In order to finish the created thread, the user function **must** use the *exit* system call (otherwise it will generate an exception), because we will not implement a correct finalization of threads.

Use service number 19 as its system call identifier.

**NOTE:** Currently there is a page directory allocated for each process (variable *dir\_pages*). Each process initializes its *dir\_pages\_baseAddr* field with one of these page directories at the beginning of the system using the *allocate\_DIR* routine.

```
extern page_table_entry dir_pages[NR_TASKS][TOTAL_PAGES];

//Assign a page directory to a process.
void allocate_DIR (struct task_struct *p) {
    int pos = ((int) p - (int) task) / sizeof(union task_union);
    p->dir_pages_baseAddr = &dir_pages[pos];
}
```

This routine uses a very simple implementation which links each page directory to a single PCB structure (the first PCB uses the first page directory, the second the next one, and so on).

With the addition of the new *clone* system call, the new thread shares his parent's page directory and therefore this routine to get a new page directory is no longer valid. A new allocator is needed in order to find a free directory to be used. This also implies that you should account when each directory is used (or reused) and when it is released.

### 5.1.2 Synchronization between processes: semaphores

In this part, you will add a new feature to help the programmer to synchronize his processes: the semaphores. Semaphores will be used to execute certain user code fragments in mutual exclusion.

A semaphore will be a resource of a process. When a process creates a new semaphore, it becomes its owner. When the process dies (exits), all the semaphores allocated to it will be destroyed.

Due to the lack of memory allocation, you will declare a static array of semaphores (a maximum of 20 semaphores will be used), and define the data structures needed for each semaphore, which are basically a counter and a queue for the blocked processes.

To manage the semaphores, you will implement some new system calls (`sem_init`, `sem_wait`, `sem_signal` and `sem_destroy`) following the steps of the system calls that you have already seen.

### *`sem_init` system call implementation*

The header of the call is:

```
int sem_init (int n_sem, unsigned int value)
    n_sem: identifier of the semaphore to be initialized
    value: initial value of the counter of the semaphore
returns: -1 if error, 0 if OK
```

Its identifier is 21. This system call initializes a semaphore identified by the number `n_sem`. It initializes the semaphore's counter to `value`. At the same time, it initializes the blocked processes queue in this semaphore and the data necessary for its correct use. The process that initializes the semaphore becomes its owner.

The return value 0 indicates a correct execution. If `n_sem` is not a valid identifier for a semaphore or it is already defined, the returned value will be -1.

### *`sem_wait` system call implementation*

Its header is:

```
int sem_wait (int n_sem)
    n_sem: identifier of the semaphore
return: -1 if error, 0 if OK
```

The identifier of this call is 22. If the counter for semaphore `n_sem` is less than or equal to zero, this call will block the process that has called it in this semaphore. If the counter is greater than zero, this call will decrement the value of the semaphore.

A return value 0 indicates a correct execution. If `n_sem` is not a valid identifier for a semaphore, or the semaphore is destroyed while the process is blocked, the returned value will be -1.

### *`sem_signal` system call implementation*

Its header is:

```
int sem_signal (int n_sem)
    n_sem: identifier of the semaphore
returns: -1 if error, 0 if OK
```



Its identifier is 23. If there is no blocked process in the semaphore `n_sem` then this call increments the counter of the semaphore `n_sem`. If there are one or more blocked processes in `n_sem`, this call will unblock the first process.

The return value 0 indicates a correct execution. If `n_sem` is not a valid identifier for a semaphore, the returned value will be -1.

### *sem\_destroy system call implementation*

Its header is:

```
int sem_destroy (int n_sem)
    n_sem: identifier of the semaphore to destroy
returns: -1 if error, 0 if OK
```

The identifier of this system call is 24. If the calling process is the owner of the semaphore, this call destroys the semaphore `n_sem`. The return value 0 indicates a correct execution. If there are blocked processes, these will be unblocked and the `sem_wait` will return -1.

If `n_sem` is not a valid identifier for a semaphore, the semaphore is not initialized, or the process is not the owner, the returned value will be -1.

Carry out the necessary tests to guarantee that the synchronization works properly (both mutual exclusion and process synchronization).

## **6 I/O MANAGEMENT AND DYNAMIC MEMORY**

### **6.1 MANAGEMENT OPERATIONS OF THE KEYBOARD DEVICE**

The main goal of this section is to describe the implementation of the keyboard-dependent *read* system call. The system has to be able to read from the keyboard in order to interact with the user.

You have to extend the keyboard interrupt functionality that you programmed during the first part of the project. An intermediate buffering system (a circular buffer) will be used to store characters read during keyboard interrupts. The *read* system call will get the characters from this buffer and *block* the process if there are other processes executing a previous *read* or if there are not enough characters.

As a consequence, the following tasks are required:

- Implement new data structures.
- Create the read system call.
- Extend the keyboard interrupt.

#### **6.1.1 New data structures**

- *Keyboardqueue*: A queue to store processes blocked by the keyboard.
  - Take into account that there are different points in the system where a process can be blocked/unblocked into a queue. It would be convenient to encapsulate these functionalities in two functions: *block* and *unblock* respectively.
- *Circular buffer*: Buffer to store the pressed keys. You can implement it using an array of fixed size. Define the necessary operations to access it in order to accomplish the requirements of a FIFO circular buffer. Define the new data type as well as the operations for handling it.

### 6.1.2 Keyboard read implementation

You must create a new system call with the following interface:

```
int read (int fd, char *buf, int count)
    read() attempts to read up to 'count' bytes from file descriptor 'fd' into
    the buffer starting at 'buf'.
```

Reads from different processes are processed in sequential order: if two processes call the read system call, the second process will not receive any character until the first one has finished reading its requested characters.

The behaviour of this system call is as follows:

- Check the parameters.
- Call the device-dependent read *sys\_read\_keyboard*. This function must copy the requested bytes from the circular buffer to the user buffer. Different cases appear:
  - If there are processes waiting for data (already blocked), then block the process at the end of the *keyboardqueue* and schedule the next process.
    - \* Block a process consists on moving the process to a list different from the *readyqueue*, getting the process out of the scheduling algorithms. In our case this corresponds to the *keyboardqueue*.
  - Otherwise, the function must copy the requested bytes from the circular buffer, and until this happens it must iterate through the following steps:
    - \* If the buffer contains all requested characters(*count*), copy them to the user buffer (*buf*) and return the total number of characters read.
    - \* If the buffer is full, then copy the whole content to the user buffer and block the process at the beginning of the *keyboardqueue*.
    - \* In any other case, block the process at the beginning of the *keyboardqueue*.

Take into account that, once a process gets blocked in the keyboard, it will be unblocked by the keyboard service routine and, therefore, the service routine has to detect when the blocked process has all the requested data.

**Note:** Remember the procedure to copy information between different address spaces.

### 6.1.3 Keyboard service routine

You must extend the functionality of the keyboard interrupt service routine to implement a buffering system. This routine must:

- Insert the new character inside the circular buffer.
- If the buffer is full and there are no processes waiting for data, then the new character is lost.
- Otherwise, 1) either the buffer is full and a process is waiting, or 2) there are enough characters in the buffer to satisfy the last read operation, this routine must unblock the first process of the *keyboardqueue*.
  - Unblock a process consists on moving the process to the *readyqueue* list, including the process in the scheduling algorithms.

## 6.2 DYNAMIC MEMORY

A mechanism to dynamically allocate/deallocate memory to user processes must be developed. A new memory region for the user address space, the *Heap*, must be defined per process. You have to add a new *sbrk* system call with the following interface:

```
void *sbrk(int increment);
```

*sbrk()* change the location of the program break, which defines the end of the process's heap segment (the program break is the first location after the end of the heap segment). Increasing the program break has the effect of allocating memory to the process; decreasing the break deallocates memory. *sbrk()* increments the program's data space by *increment* bytes. Calling *sbrk()* with an *increment* of 0 can be used to find the current location of the program break.

On success, *sbrk()* returns the previous program break. (If the break was increased, then this value is a pointer to the start of the newly allocated memory). On error, (void \*) -1 is returned, and *errno* is set to ENOMEM.

Figure 28 shows the resulting address space after using this new system call.

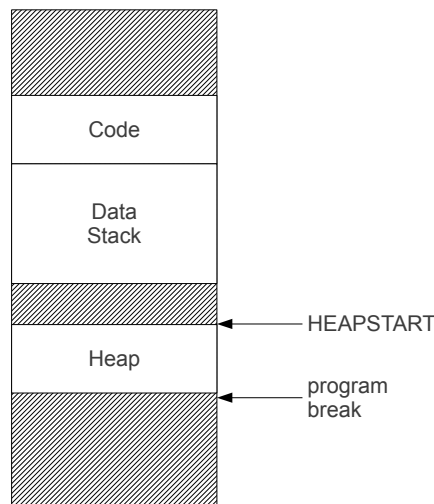


Fig. 28. Resulting user address space after using the *sbrk* system call.

**NOTE:** The memory is allocated per process and it is **not** shared. It must be *copied* when a new process is created and *freed* when the process is destroyed. Take into account that a page allocator must be created in the system code, but the user allocates bytes, and therefore, a fine grain accounting has to be maintained.