

Chapter 5: Process Management in Theory and x86

Introduction

In the previous chapters, we have discussed the topics that helped us to understand the basics that are needed to initialize the environment for a 32-bit protected mode kernel running on x86. Starting from this chapter we are going to discuss the topics that belong to the kernel itself, that is, the responsibilities of the kernel. We start with a quick look on the theories that are traditionally presented on academic textbooks, then we move to the practical part in order to implement these theories (or part of them) in 539kernel. A good place to start from is *process management*.

A kernel has multiple responsibilities, one of these responsibilities is to manage the resources and make sure they are managed well. One important resource of the computers is the time of the processor (AKA: CPU time) which is the component that executes the code of software that we would like to run on our machines. Process management is the the part that studies how a kernel should manage and distribute CPU time among a bunch of *processes*.

The Most Basic Work Unit: A Process

Process is the term which is used in operating systems literature to describe a running program. In the previous chapters of this book we have encountered the concept of the process multiple times and you may recall from these encounters that every user-space software that we use in our computers are soulless sequence of bytes that are stored somewhere in the hard disk. When we decide to use a specific software, for example, the web browser, the first thing we do is to open it either through double clicking on its icon in graphical user interfaces or through writing its command in the shell. When we do that, the kernel is needed to be called through a *system call* and take the responsibility of “opening” this software, we can consider system calls as functions which are provided by the kernel to expose its services for the user-space software, one way of implementing system calls is to use interrupts, exactly the same way that we have used with BIOS.

However, there are multiple steps that are needed to be performed to open the software, for example, reading its data from disk, but our current focus is on process-related parts, eventually, the kernel creates a new process for the software that we requested to open. The kernel maintains a table of current processes, each entry represents a process and contains the information which is needed by the kernel to manage the process, this data structure which stores a process information is known as *process control block* (PCB).

Of course, the most important part of the process is the code of the software that this process represents and its data, both data ¹ and code should be loaded into memory, after that, its code can be executed by the processor. We need to note that a process is an instance of a software, in other words, one software can be opened more than one time with a separated process for each one, for example, opening multiple windows of the web browser on the same time, the software is one which is the web browser and it is represented by the binary file which is stored in the hard disk, but each opened window is a separated process and each process' content is stored in the main memory. While the described concept is well-known by the term "process", specially in the literature, other terms can be used for the same concept, for example *task* and *job* are other words which are used to point to the same concept.

Each process is known to have a *state*, when it is loaded to the memory its state will be indicated by the kernel as *ready* and can be run anytime, when the CPU time is given to a process its state will be *running*. Also, the state of the process can be *wait*, an example of a situation where a process state is changed to wait is when it performs I/O request (e.g. read from the hard disk), its state will be *wait* since it's waiting for the I/O device to fulfill the request. The state information about a process is stored in the process control block which is, as mentioned earlier, an entry in the processes table.

Sometimes, a bunch of processes in a system need to communicate with each other to share some data or tell each other to perform a specific operation, this need led to a broad topic known as *inter-process communication* (IPC) which provides mechanisms to make this type of communication possible. The applications of IPC is not restricted to operating system kernels, they are used in distributed computing for instance. One well known mechanism of IPC is *shared memory*, that is, a region of the memory that is accessible by more than one process, they can read and write to this region in order to share data. The ability to write to same place by more than one process can cause a problem known as *race condition*, given a shared variable, the situation which two or more processes try to change the value of this variable at the same time is known as race condition. There are multiple solutions for this problem and this topic is studied in a branch known *concurrency control* which is a shared topic by many applications, one of them is database management systems (DBMS) which needs these mechanisms when two users try to update the same row at the same time.

Processes are not the only entities that need to communicate, there is another unit of work which is known as *thread* and it can be described as lightweight process. A process can have more than one thread and when a software uses more than one thread to do its job, it is described as *multithreaded*. Threads are everywhere in our usage of computers, and a world without them is unimaginable. For example, when you use a text editor, the main thread of the software lets

¹We mean static data here, which are contained in the binary file of the software. While the data that are generated by the running process are not loaded from the binary file, instead they are created while the code is running.

you write your text, but when you click on save icon a separated thread within the text editor's process is used to perform this operation. Furthermore, another thread can be used for the spell checker while you are typing. If all there functionalities were on one thread, you need to wait each one of them to finish in order to let you to perform the other functionality, that is, the software without threads is going to run sequentially while threads provide us concurrency within one process. Threads and processes have many similarities, for example, both of them are there to be executed, hence, they need to use the processor and both of them need to be scheduled to give every one of them time from the processor. In contrast to processes, threads run as a part of same process and share the same address space of the software which makes the communication between them much easier.

The Basics of Multitasking

When we write a kernel, multiple design questions should be answered ² and the part of process management is not an exception of that. There are multiple well-known answers for some basic design questions, each one of those answers tries to solve a problem that faced the person who gave us this answer, for example, one of well-known features of the modern operating systems is *multitasking* which is the successor of *monotasking* and both of them can be considered as answers for a design question in operating systems. In multitasking environment, the system can run multiple processes at a time even if there is only one processor available, while monotasking environment, the system can run only one process at a time until this process finishes its work or the user closes it, only after that, another process can be run.

Mutliprogramming & Time-Sharing

In the days of monotasking, we were facing a serious problem that led to the birth of multitasking, it has been noticed that the processes tend to have idle time, for example, when the process is waiting for the hard disk to fetch some stored data, the process will be idle while it is taking over the processor, that is, the processor is under its control but it doesn't use it for something useful, that means, in this case, we are wasting the valuable resource of CPU time in waiting for some action to finish, we need to utilize the processor as much as possible, and here came the solution of this problem, by letting the kernel to have a list of processes that are *ready* to run, assuming the machine has just one processor with one core, the CPU time will be given to, say process A, for some time, at some point of running time, process A requests from the disk some data and due to that it becomes idle waiting for disk to respond, instead of keep the control of the processor under the process A, which is doing nothing but waiting right now, the kernel suspends process A and give the CPU time to another process, say

²Remember the job of a kernelist!

process B, this switching between two processes is known as *context switch*. The process B is going to use the processor while process A is waiting for the disk to respond. At some point, process B will perform some action that makes it idle which means that the kernel can switch to another ready process and so on. This solution is known as *multiprogramming*. To sum it up, we have a list of ready processes, choose one, give it the CPU time and wait for it until it becomes idle, since it's waiting for something switch to another process which is not idle.

Better yet, multiprogramming has been extended to utilize the processor more efficiently. Instead of waiting for the currently running process to perform something which makes it idle, why don't we suspend it after some period of running time whether it is idle or not and switch to another process? This solution is known as *time sharing* which is with multiprogramming represent the scheme that modern operating systems use for multitasking. In time sharing, a list of ready processes is available for the kernel, in each unit of time, say for example, every 1 second ³, the currently running process will be suspended by the kernel and another process is given the CPU time and so on.

Process Scheduling

You may recall from the previous chapter the system timer which emits an interrupt every unit of time, this interrupt can be used to implement time sharing to switch between the processes of the system, of course the kernel needs an algorithm to choose which process to run next, this kind of algorithms are known as *scheduling algorithms* and in general this part of the topic is known as *scheduling* in which we try to find the best way of choosing the next process to run in order to satisfy our requirements. The *scheduler* is the part of the kernel that schedules the next process by using some specific scheduling algorithm. There are many scheduling algorithms to deal with different requirements and one of them is known as *round-robin*. In this algorithm, each process is given a fixed amount of CPU time known as *quantum*, when the running process finishes its quantum the scheduler will be called and the CPU time will be given to the next process in the list until its quantum finishes and so on until the schedule reaches to the end of the process list where it starts with the first process again. The value of the quantum is decided by the kernelist, for example 50 milliseconds, which means each process will run 50 milliseconds then suspended to run the next one on the list and so on.

Process Context

As you know, when a process is executing, it can use the registers of the processor (e.g. `eax`) to store its own data. Also, it may change the values of segment registers in case it is running under a system that employs segmented-memory

³In practice, it is shorter.

instead of flat memory model. Furthermore, the value of the instruction pointer EIP will contain an address which belong to the process' address space. All these values that are related to a process and stored inside the registers of the processor are known as *process context*, another term which is *process state* may also be used in some places to mean the same thing, but to not confuse this concept with the one that we have defined the term "state" with previously, it is better to use the term process context.

When a scheduler decides to change the currently running process, let's call it A, through a context switch, a copy of the context of the process A should be taken and be stored somewhere, that is, a snapshot of the last context of A is taken before switching to another process. By taking this snapshot, it will be easy later to resume process A by just loading its context to the processor's register and jump to the the value of EIP which has been just loaded.

Preemptive & Cooperative Multitasking

Both multiprogramming and time-sharing solutions give us a type of multitasking known as *preemptive multitasking*, the processes are forced by the kernel to give the CPU time to another process and no process can take over the processor for the whole time. Another type of multitasking is known as *cooperative multitasking* (or *non-preemptive multitasking*), in this type the context switching is not perform forcibly, instead, the currently running process should cooperate and voluntarily tells the kernel when it should be suspended and a context switch should be performed. One of the obvious problem of this way, at least for the well-known workloads (e.g. servers or desktops), that a process, which runs a code that has been written by someone we don't know, cannot be trusted. It simply may take over the CPU time and never cooperate and give it up due to an error in the code or even in purpose ⁴.

Multitasking in x86

With the assistance of system timer, multitasking can be realized fully by the kernel, that is, by the code. This type is known as *software multitasking*, that is, the kernel itself is responsible for storing the list of processes and their related information. Also, it's responsible for storing a snapshot of process context before performing context switching and resuming this snapshot when the process is resumed. On the other hand, In x86 some features are provided to handle these things with the assistance of the processor itself, this type is known as *hardware multitasking*.

⁴You may ask who would use cooperative multitasking and gives this big trust to the code of the software! In fact, the versions of Windows before 95 used this style of multitasking, also, Classic Mac OS used it. Why? You may ask, I don't know exactly, but what I know for sure is that humanity is in a learning process!

While hardware multitasking is available in x86, the modern operating system kernels don't use it, instead, multitasking is implemented by the kernel itself. One reason to take this decision is portability. Modern kernels tend to run on more than one architecture and not only x86, by using as little as possible of architecture's features it will be easier to port a kernel to other architectures.

In this section, we are going to cover the basics of hardware multitasking in x86 which are needed to initialize the environment to make it work correctly, in the same way as GDT and flat memory model. Furthermore, I think knowing the other available options is important, especially for kernelists. In 539kernel we are going to implement software multitasking as in modern kernels.

Task-State Segments

The most basic component of hardware multitasking in x86 is known as *task-state segment* (TSS) ⁵ which is a segment in the memory as any other code or data segment, it has what other segments have, a base address, a limit and properties. The difference from code and data segments is that TSS is a system segment ⁶, this segment stores the context of a specific process.

In hardware multitasking, each process should have its own TSS, and each TSS should have an entry in the GDT table, that is, *TSS descriptor*. A special register known as *task register* should contain the segment selector of the currently running process's TSS descriptor, the instruction `ltr` is used to store a value in this register.

Figure shows the structure of a task-state segment, as you can see, most of the fields are values of registers while the others are out of our topic's range except for previous task link which will be covered in a moment. You can see that stack segment register and stack pointer register have four entries instead of one, `SS`, `SS0`, `SS1` and `SS2` for stack segment register. `ESP`, `ESP0`, `ESP1` and `ESP2` for stack pointer register. These fields points to the stack that should be used when the process is in a specific privilege level, for example, `SS0:ESP0` will be used as the stack of the process when it switches to privilege level 0, when it switches back to privilege level 3 the stack `SS:ESP` will be used instead, and the same is applicable to the other similar fields. If we intend to implement software multitasking, the sole reason of defining at least one TSS is due to these fields, when a switch between privilege levels occurs, the processor needs a TSS to use these fields from it in order to switch between stacks. This is needed only when the system runs user-space code, that is, privilege level 3 code.

Figure shows the structure of TSS descriptor in GDT table. As you can see, all fields and flags are same as the ones that we have already explained in chapter . The differences are in the *type field* which has static values in TSS descriptor

⁵In x86, the term task is used instead of process.

⁶In chapter we have seen that there are two types of segments in x86, application segments such as code, data and stack segment. And system segments and they are LDT and TSS.

and a flag called B, or *busy flag* which should be 1 when the process that this TSS descriptor represents is active and 0 when it is inactive.

Context Switching in x86

One way of switching from a process to another ⁷ in x86 hardware multitasking is to call or jump to TSS descriptor in GDT, assume that the system timer caused the call of the scheduler which selects process A as the next process to run, the scheduler can cause context switch by using the instructions `call` or `jmp` and the operand should be the segment selector of A's TSS descriptor. In this way, the processor is going to take a copy of currently running process (call it B) and stores it in B's own TSS, then the values in A's TSS will be loaded into the processor registers and then execution of A begins.

Another way of context switching in x86 hardware multitasking is to call or jump to a task gate. In chapter , when we discussed the descriptors of IDT, we have said that one type of descriptor that can be defined is a task gate descriptor. This kind of descriptors is considered as a separated process by the process, when we jump or call a task gate, the previously explained mechanism of task switching will be performed. Task gates can also be defined in GDT and LDT. In the IDT table of 539kernel we have chosen to not define the interrupts as task gates, we don't want to perform a context switch with each interrupt.

When a process is called instead of jumped to, eventually, it should return to the caller process by using the instruction `iret`, for the processor, to be able to decide which task is the caller, the previous task link field of the callee's TSS will be updated to contain the segment selector of the caller process. In this way, when `iret` instruction is executed, it will be easy to know to which process the processor should switch back.

⁷In x86, context switch is known as task switch.