

Case Study on GNU/UNIX Based Operating Systems

Samrat Nakarmi

November 2024

1 Introduction

The operating system is the most important part of any computer system, it is often called the mother of all system software because of how important it is to a functioning computer system. The most popular operating system in the world is Windows which is a proprietary operating system by the Microsoft Corporation. Latest data show that Windows has a majority market share by a large amount, 74%. Other operating systems have rest of the market.

Operating System	Market Share%
Windows	73.41%
OSX	15.49%
Unknown	4.66%
Linux	4.31%
Chrome OS	2.12%
FreeBSD	0%

Windows is the most common among common people who use their computer for daily activities like browsing the web, watching movies etc and Linux is primarily used by sophisticated people who work in the tech domain. Almost all servers in the world are Linux and other Systems like Smart TV, mobile devices use different flavors of Linux. Almost all serious Software Development work is done in Linux and all real software is running on Linux Servers around the world.

1.1 Usage

UNIX is a popular operating system, and is used heavily in a large variety of scientific, engineering, and mission critical applications. Interest in UNIX has grown substantially high in recent years because of the proliferation of the Linux (a Unix look-alike) operating system. The following are the uses of the UNIX operating system:

- Universally used for high-end number crunching applications.
- Wide use in CAD/CAM arena.
- Ideally suited for scientific visualization.
- DHCP, NetNews, Mail, etc., due to networking being in the kernel for a long time now.

UNIX was founded on what could be called a “small is good” philosophy. The idea is that each program is designed to do one job efficiently. Because UNIX was developed by different people with different needs it has grown to an operating system that is both flexible and easy to adapt for specific needs. UNIX was written in a machine independent language and C language. So UNIX and UNIX-like operating systems can run on a variety of hardware. These systems are available from many different sources, some of them at no cost. Because of this diversity and the ability to utilize the same “user-interface” on many different systems, UNIX is said to be an open system. At the time the first UNIX was written, most operating systems developers believed that an operating system must be written in an assembly language so that it could function effectively and gain access to the hardware. The UNIX Operating System is written in C. The C language itself operates at a level that is just high enough to be portable to a variety of computer hardware. Most publicly distributed UNIX software is written in C and must be compiled before use. In practical terms this means that an understanding of C can make the life of a UNIX system administrator significantly easier.

In the earlier units, we have studied various function of OS in general. In this unit we will study a case study on UNIX and how the UNIX handles various operating system functions. In the MCSL-045, section – 1, the practical sessions are given to provide you the hands-on experience.

1.2 Brief History

In the late 1960s, General Electric, MIT and Bell Labs commenced a joint project to develop an ambitious multi-user, multi-tasking OS for mainframe computers known as MULTICS. MULTICS was unsuccessful, but it did motivate Ken Thompson, who was a researcher at Bell Labs, to write a simple operating system himself. He wrote a simpler version of MULTICS and called his attempt UNICS (Uniplexed Information and Computing System). To save CPU power and memory, UNICS (finally shortened to UNIX) used short commands to lessen the space needed to store them and the time needed to decode them - hence the tradition of short UNIX commands we use today, e.g., ls, cp, rm, mv etc. Ken Thompson then teamed up with Dennis Ritchie, the creator of the first C compiler in 1973. They rewrote the UNIX kernel in C and released the Fifth Edition of UNIX to universities in 1974. The Seventh Edition, released in 1978, marked a split in UNIX development into two main branches: SYS V (System 5) and BSD (Berkeley Software Distribution). Linux is a free open source UNIX

OS. Linux is neither pure SYS V nor pure BSD. Instead, it incorporates some features from each (e.g. SYSV-style startup files but BSD-style file system layout) and plans to conform to a set of IEEE standards called POSIX (Portable Operating System Interface). You can refer to the MCSL-045 Section – 1 for more details in a tabular form. This tabular form will give you a clear picture about the developments.

2 Structure of Operating System

UNIX is a layered operating system. The innermost layer is the hardware that provides the services for the OS. The following are the components of the UNIX OS.

2.1 The Kernel

The operating system, referred to in UNIX as the kernel, interacts directly with the hardware and provides the services to the user programs. These user programs don't need to know anything about the hardware. They just need to know how to interact with the kernel and it's up to the kernel to provide the desired service. One of the big appeals of UNIX to programmers has been that most well written user programs are independent of the underlying hardware, making them readily portable to new systems. User programs interact with the kernel through a set of standard system calls.

These system calls request services to be provided by the kernel. Such services would include accessing a file: open close, read, write, link, or execute a file; starting or updating accounting records; changing ownership of a file or directory; changing to a new directory; creating, suspending, or killing a process; enabling access to hardware devices; and setting limits on system resources. UNIX is a multi-user, multi-tasking operating system. You can have many users logged into a system simultaneously, each running many programs. It is the kernel's job to keep each process and user separate and to regulate access to system hardware, including CPU, memory, disk and other I/O devices.

2.2 The Shell

The shell is often called a command line shell, since it presents a single prompt for the user. The user types a command; the shell invokes that command, and then presents the prompt again when the command has finished. This is done on a line-by-line basis, hence the term "command line". The shell program provides a method for adapting each user's setup requirements and storing this information for re-use. The user interacts with /bin/sh, which interprets each command typed. Internal commands are handled within the shell (set, unset), external commands are cited as programs (ls, grep, sort, ps). There are a number of different command line shells (user interfaces).

- Bourne (sh)

- C Shell (csh)
- Bourne Again Shell (bash)
- Z-Shell (zsh)

2.3 System Utilities

The system utilities are intended to be controlling tools that do a single task exceptionally well (e.g., `grep` finds text inside files while `wc` counts the number of words, lines and bytes inside a file). Users can solve problems by integrating these tools instead of writing a large monolithic application program. Like other UNIX flavours, Linux's system utilities also embrace server programs called daemons that offer remote network and administration services (e.g. `telnetd` provides remote login facilities, `httpd` serves web pages).

2.4 Application Programs

Some application programs include the `emacs` editor, `gcc` (a C compiler), `g++` (a C++ compiler), `xfig` (a drawing package), `latex` (a powerful typesetting language). UNIX works very differently. Rather than having kernel tasks examine the requests of a process, the process itself enters kernel space. This means that rather than the process waiting “outside” the kernel, it enters the kernel itself (i.e. the process will start executing kernel code for itself). This may sound like a formula for failure, but the ability of a process to enter kernel space is strictly prohibited (requiring hardware support). For example, on x86, a process enters kernel space by means of system calls - well known points that a process must invoke in order to enter the kernel. When a process invokes a system call, the hardware is switched to the kernel settings. At this point, the process will be executing code from the kernel image. It has full powers to wreak disaster at this point, unlike when it was in user space. Furthermore, the process is no longer pre-emptible

3 Process Management

When the computer is switched on, the first thing it does is to activate resident on the system board in a ROM (read-only memory) chip. The operating system is not available at this stage so that the computer must “pull itself up by its own boot- straps”. This procedure is thus often referred to as bootstrapping, also known as cold boot. Then the system initialization takes place. The system initialization usually involves the following steps. The kernel,

- Tests to check the amount of memory available.
- Probes and configures hardware devices. Some devices are usually compiled into the kernel and the kernel has the ability to autoprobe the hard-

ware and load the appropriate drivers or create the appropriate entries in the `/dev` directory.

- Sets up a number of lists or internal tables in RAM. These are used to keep track of running processes, memory allocation, open files, etc.

Depending on the UNIX version, the kernel now creates the first UNIX processes. A number of dummy processes (processes which cannot be killed) are created first to handle crucial system functions. A `ps -ef` listing on each OS shows will show you the existing processes. `init` is the last process created at boot time. It always has a process ID (PID) of 1. `init` is responsible for starting all subsequent processes. Consequently, it is the parent process of all (non-dummy) UNIX processes. Don't confuse the `init` process with the system `init` command. The `init` command (usually found in `/sbin` or `/usr/sbin`) is used by `root` to put the system into a specific run level. All subsequent processes are created by `init`. For example, one of the processes started by `init` is `inetd`, the internet super daemon. (`inetd`, in turn, creates many other processes, such as `telnetd`, on demand.) In the Unix process hierarchy, `init` is called the parent process and `inetd` the child of `init`. Any process can have any number of children (up to the kernel parameter `nproc`, the maximum allowed number of processes). If you kill the parent of a child process, it automatically becomes the child of `init`. Each running process has associated with it a process ID or PID. In addition, each process is characterized by its parent's PID or PPID. Finally, each process runs at a default system priority (PRI). The smaller the numerical value of the PRI, the higher the priority and vice versa.

3.1 Management of Process By Kernel

For each new process created, the kernel sets up an address space in the memory. This address space consists of the following logical segments:

- *text* - contains the program's instructions
- *data* - contains initialized program variables
- *bss* - contains uninitialized program variables
- *stack* - a dynamically growable segment, it contains variables allocated locally and parameters passed to functions in the program

Each process has two stacks: a user stack and a kernel stack. These stacks are used when the process executes in the user or kernel mode (described below).

3.1.1 Mode Switching

At least two different modes of operation are used by the UNIX kernel - a more privileged kernel mode, and a less privileged user mode. This is done to protect some parts of the address space from user mode access.

User Mode Processes, created directly by the users, whose instructions are currently executing in the CPU are considered to be operating in the user-mode. Processes running in the user mode do not have access to code and data for other users or to other areas of address space protected by the kernel from user mode access.

Kernel Mode Processes carrying out kernel instructions are said to be running in the kernel-mode. A user process can be in the kernel-mode while making a system call, while generating an exception/fault, or in case on an interrupt. Essentially, a mode switch occurs and control is transferred to the kernel when a user program makes a system call. The kernel then executes the instructions on the user's behalf. While in the kernel-mode, a process has full privileges and may access the code and data of any process (in other words, the kernel can see the entire address space of any process).

3.1.2 The Context of a Process and Context Switching

The context of a process is essentially a snapshot of its current runtime environment, including its address space, stack space, etc. At any given time, a process can be in user-mode, kernel-mode, sleeping, waiting on I/O, and so on. The process scheduling subsystem within the kernel uses a time slice of typically 20ms to rotate among currently running processes. Each process is given its share of the CPU for 20ms, then left to sleep until its turn again at the CPU. This process of moving processes in and out of the CPU is called context switching. The kernel makes the operating system appear to be multi-tasking (i.e. running processes concurrently) via the use of efficient context-switching.

At each context switch, the context of the process to be swapped out of the CPU is saved to RAM. It is restored when the process is scheduled its share of the CPU again. All this happens very fast, in microseconds.

To be more precise, context switching may occur for a user process when

- a system call is made, thus causing a switch to the kernel-mode,
- a hardware interrupt, bus error, segmentation fault, floating point exception, etc. occurs,
- a process voluntarily goes to sleep waiting for a resource or for some other reason, and
- the kernel preempts the currently running process (i.e. a normal process scheduler event).

Context switching for a user process may occur also between threads of the same process. Extensive context switching is an indication of a CPU bottleneck.

3.1.3 Communication between the Running Processes

UNIX provides a way for a user to communicate with a running process. This is accomplished via signals, a facility which enables a running process to be notified about the occurrence of a) an error event generated by the executing process, or b) an asynchronous event generated by a process outside the executing process. Signals are sent to the process ultimately by the kernel. The receiving process has to be programmed such that it can catch a signal and take a certain action depending on which signal was sent.

4 Memory Management

One of the numerous tasks the UNIX kernel performs while the machine is up is to manage memory. In this section, we explore relevant terms (such as physical vs. virtual memory) as well as some of the basic concepts behind memory management.

4.1 Physical vs Virtual Memory

UNIX, like other advanced operating systems, allows you to use all of the physical memory installed in your system as well as area(s) of the disk (called swap space) which have been designated for use by the kernel in case the physical memory is insufficient for the tasks at hand. Virtual memory is simply the sum of the physical memory (RAM) and the total swap space assigned by the system administrator at the system installation time.

$$VirtualMemory(VM) = PhysicalRAM + SwapSpace$$

4.2 Dividing Memory into Pages

The UNIX kernel divides the memory into manageable chunks called pages. A single page of memory is usually 4096 or 8192 bytes (4 or 8KB). Memory pages are laid down contiguously across the physical and the virtual memory

4.3 Cache Memory

With increasing clock speeds for modern CPUs, the disparity between the CPU speed and the access speed for RAM has grown substantially. Consider the following:

- Typical CPU speed today: 250-500MHz (which translates into 4-2ns clock tick)
- Typical memory access speed (for regular DRAM): 60ns
- Typical disk access speed: 13ms

In other words, to get a piece of information from RAM, the CPU has to wait for 15-30 clock cycles, a considerable waste of time. Fortunately, cache RAM has come to the rescue. The RAM cache is simply a small amount of very fast (and thus expensive) memory which is placed between the CPU and the (slower) RAM. When the kernel loads a page from RAM for use by the CPU, it also prefetches a number of adjacent pages and stores them in the cache. Since programs typically use sequential memory access, the next page needed by the CPU can now be supplied very rapidly from the cache. Updates of the cache are performed using an efficient algorithm which can enable cache hit rates of nearly 100% (with a 100% hit ratio being the ideal case). CPUs today typically have hierarchical caches. The on-chip cache (usually called the L1 cache) is small but fast (being on-chip). The secondary cache (usually called the L2 cache) is often not on-chip (thus a bit slower) and can be quite large; sometimes as big as 16MB for high-end CPUs (obviously, you have to pay a hefty premium for a cache that size).

4.4 Memory Organisation by the Kernel

When the kernel is first loaded into memory at the boot time, it sets aside a certain amount of RAM for itself as well as for all system and user processes. Main categories in which RAM is divided are

- *Text*: to hold the text segments of running processes.
- *Data*: to hold the data segments of running processes.
- *Stack*: to hold the stack segments of running processes.
- *Shared Memory*: This is an area of memory which is available to running programs if they need it. Consider a common use of shared memory: Let us assume you have a program which has been compiled using a shared library (libraries that look like libxxx.so; the C-library is a good example - all programs need it). Assume that five of these programs are running simultaneously. At run-time, the code they seek is made resident in the shared memory area. This way, a single copy of the library needs to be in memory, resulting in increased efficiency and major cost savings.
- *Buffer Cache*: All reads and writes to the file system are cached here first. You may have experienced situations where a program that is writing to a file doesn't seem to work (nothing is written to the file). You wait a while, then a sync occurs, and the buffer cache is dumped to disk and you see the file size increase.

4.5 The System and User Areas

When the kernel loads, it uses RAM to keep itself memory resident. Consequently, it has to ensure that user programs do not overwrite/corrupt the kernel data structures (or overwrite/corrupt other users' data structures). It

does so by designating part of RAM as kernel or system pages (which hold kernel text and data segments) and user pages (which hold user stacks, data, and text segments). Strong memory protection is implemented in the kernel memory management code to keep the users from corrupting the system area. For example, only the kernel is allowed to switch from the user to the system area. During the normal execution of a Unix process, both system and user areas are used. A common system call when memory protection is violated is SIGSEGV.

4.6 Paging vs. Swapping

Paging When a process starts in UNIX, not all its memory pages are read in from the disk at once. Instead, the kernel loads into RAM only a few pages at a time. After the CPU digests these, the next page is requested. If it is not found in RAM, a page fault occurs, signaling the kernel to load the next few pages from disk into RAM. This is called demand paging and is a perfectly normal system activity in UNIX. (Just so you know, it is possible for you, as a programmer, to read in entire processes if there is enough memory available to do so.) The UNIX daemon which performs the paging out operation is called pageout. It is a long running daemon and is created at boot time. The pageout process cannot be killed.

Swapping Let's say you start ten heavyweight processes (for example, five xterms, a couple netscapes, a sendmail, and a couple pines) on an old 486 box running Linux with 16MB of RAM. Basically, you do not have enough physical RAM to accommodate the text, data, and stack segments of all these processes at once. Since the kernel cannot find enough RAM to fit things in, it makes use of the available virtual memory by a process known as swapping. It selects the least busy process and moves it in its entirety (meaning the program's in-RAM text, stack, and data segments) to disk. As more RAM becomes available, it swaps the process back in from disk into RAM. While this use of the virtual memory system makes it possible for you to continue to use the machine, it comes at a very heavy price. Remember, disks are relatively slower (by the factor of a million) than CPUs and you can feel this disparity rather severely when the machine is swapping. Swapping is not considered a normal system activity. It is basically a sign that you need to buy more RAM.

The process handling swapping is called sched (in other UNIX variants, it is sometimes called swapper). It always runs as process 0. When the free memory falls so far below minfree that pageout is not able to recover memory by page stealing, sched invokes the syscall sched(). Syscall swapout is then called to free all the memory pages associated with the process chosen for being swapping out. On a later invocation of sched(), the process may be swapped back in from disk if there is enough memory.

4.7 Demand Paging

Berkeley introduced demand paging to UNIX with BSD (Berkeley System) which transferred memory pages instead of processes to and from a secondary device; recent releases of UNIX system also support demand paging. Demand paging is done in a straightforward manner. When a process needs a page and the page is not there, a page fault to the kernel occurs, a frame of main memory is allocated, and then the process is loaded into the frame by the kernel.

The advantage of demand paging policy is that it permits greater flexibility in mapping the virtual address of a process into the physical memory of a machine, usually allowing the size of a process to be greater than the amount of availability of physical memory and allowing more Processes to fit into main memory. The advantage of a swapping policy is that is easier to implement and results in less system overhead.

5 File System in Linux

Physical disks are partitioned into different file systems. Each file system has a maximum size, and a maximum number of files and directories that it can contain. The file systems can be seen with the `df` command. Different systems will have their file systems laid out differently. The `/` directory is called the root of the file system. The UNIX file system stores all the information that relates to the long-term state of the system. This state includes the operating system kernel, the executable files, configuration information, temporary work files, user data, and various special files that are used to give controlled access to system hardware and operating system functions. The constituents of UNIX file system can be one of the following types:

5.1 Ordinary Files

Ordinary files can contain text, data, or program information. Files cannot contain other files or directories. UNIX filenames are not broken into a name part and an extension part, instead they can contain any keyboard character except for `'/'` and be up to 256 characters long. However, characters such as `*,?,#` and `&` have special meaning in most shells and should not therefore be used in filenames. Putting spaces in filenames also makes them difficult to manipulate, so it is always preferred to use the underscore `'_'`.

5.2 Directories

Directories are folders that can contain other files and directories. A directory is a collection of files and/or other directories. Because a directory can contain other directories, we get a directory hierarchy. The “top level” of the hierarchy is the root directory.

The UNIX file system is a hierarchical tree structure with a top-level directory known as the root (designated by a slash `'/'`). Because of the tree structure,

a directory can have many child directories, but only one parent directory. The layout is shown in Figure 2. The path to a location can be defined by an absolute path from the root /, or as a relative path from the current working directory. To specify an absolute path, each directory from the source to the destination must be included in the path, with each directory in the sequence being separated by a slash. To specify a relative path, UNIX provides the shorthand “.” for the current directory and “..” for the parent directory e.g., The absolute path to the directory “play” is /sbin/will/play, while the relative path to this directory from “zeb” is ../will/play. Various UNIX directories and their contents are listed in the table given below.

Directory	Content
/	The "root" directory
/bin	Essential low-level system utilities
/usr/bin	Higher-level system utilities and application programs
/sbin	Superuser system utilities (for performing system administration tasks)
/lib	Program libraries for low-level system utilities
/usr/lib	Program libraries for higher-level user programs
/tmp	Temporary file storage space (can be used by any user)
/home	User home directories for personal file space
/etc	UNIX system configuration and information files
/dev	Hardware devices
/proc	Pseudo-filesystem for kernel interface and active processes

6 CPU SCHEDULING

CPU scheduling in UNIX is designed to benefit interactive processes. Processes are given small CPU time slices by a priority algorithm that reduces to round-robin scheduling for CPU-bound jobs. The scheduler on UNIX system belongs to the general class of operating system schedulers known as round robin with multilevel feedback which means that the kernel allocates the CPU time to a process for small time slice, pre-empts a process that exceeds its time slice and feed it back into one of several priority queues. A process may need much iteration through the “feedback loop” before it finishes. When kernel does a context switch and restores the context of a process, the process resumes execution from the point where it had been suspended. The scheduler on UNIX system belongs to the general class of operating system schedulers known as round robin with multilevel feedback which means that the kernel allocates the CPU time to a process for small time slice, pre-empts a process that exceeds its time slice and feed it back into one of several priority queues. A process may need much iteration through the “feedback loop” before it finishes. When kernel does a context switch and restores the context of a process, the process resumes execution from the point where it had been suspended. Each process table entry contains a priority field. There is a process table for each process

which contains a priority field for process scheduling. The priority of a process is lower if they have recently used the CPU and vice versa. The more CPU time a process accumulates, the lower (more positive) its priority becomes, and vice versa, so there is negative feedback in CPU scheduling and it is difficult for a single process to take all the CPU time. Process aging is employed to prevent starvation. Older UNIX systems used a one second quantum for the round-robin scheduling. 4.33SD reschedules processes every 0.1 second and recomputes priorities every second. The round-robin scheduling is accomplished by the time-out mechanism, which tells the clock interrupt driver to call a kernel subroutine after a specified interval; the subroutine to be called in this case causes the rescheduling and then resubmits a time-out to call itself again. The priority recomputation is also timed by a subroutine that resubmits a time-out for itself event. The kernel primitive used for this purpose is called sleep (not to be confused with the user-level library routine of the same name.) It takes an argument, which is by convention the address of a kernel data structure related to an event that the process wants to occur before that process is awakened. When the event occurs, the system process that knows about it calls wakeup with the address corresponding to the event, and all processes that had done a sleep on the same address are put in the ready queue to be run.

7 Summary

We discussed issues broadly related to features of UNIX OS, boot process, system initialization, process management, memory management, file system and CPU scheduling in UNIX operating system. In this unit we discussed several theoretical concepts of UNIX operating system in general, it is often useful to use them in your lab for practice. Refer to the Section – 1 of MCSL-045, in which we have covered the practical component of the UNIX operating system