

Anatomy of the Linux Kernel

The Linux® kernel is the core of a large and complex operating system, and while it is huge, it is well organized in terms of subsystems and layers. In this article, you can explore the general structure of the Linux kernel and get to know its major subsystems and core interfaces. Where possible, you get links to other IBM articles to help you dig deeper.

Given that the goal of this article is to introduce you to the Linux kernel and explore its architecture and major components, let's start with a short tour of Linux kernel history, then look at the Linux kernel architecture from 30,000 feet, and, finally, examine its major subsystems. The Linux kernel is over six million lines of code, so this introduction is not exhaustive. Use the pointers to more content to dig in further.

A short tour of Linux history

While Linux is arguably the most popular open source operating system, its history is actually quite short considering the timeline of operating systems. In the early days of computing, programmers developed on the bare hardware in the hardware's language. The lack of an operating system meant that only one application (and one user) could use the large and expensive device at a time. Early operating systems were developed in the 1950s to provide a simpler development experience. Examples include the General Motors Operating System (GMOS) developed for the IBM 701 and the FORTRAN Monitor System (FMS) developed by North American Aviation for the IBM 709.

In the 1960s, the Massachusetts Institute of Technology (MIT) and a host of companies developed an experimental operating system called Multics (or Multiplexed Information and Computing Service) for the GE-645. One of the developers of this operating system, AT&T, dropped out of Multics and developed their own operating system in 1970 called Unics. Along with this operating system was the C language, for which C was developed and then rewritten to make operating system development portable.

Twenty years later, Andrew Tanenbaum created a microkernel version of UNIX®, called MINIX (for minimal UNIX), that ran on small personal computers. This open source operating system inspired Linus Torvalds' initial development of Linux in the early 1990s

Linux quickly evolved from a single-person project to a world-wide development project involving thousands of developers. One of the most important decisions for Linux was its adoption of the GNU General Public License (GPL). Under the GPL, the Linux kernel was protected from commercial exploitation, and it also benefited from

operating system

the software that manages the sharing of the resources of a computer and provides programmers with an interface used to access those resources

the user-space development of the GNU project (of Richard Stallman, whose source dwarfs that of the Linux kernel). This allowed useful applications such as the GNU Compiler Collection (GCC) and various shell support.

Introduction to the Linux kernel

Now on to a high-altitude look at the GNU/Linux operating system architecture. You can think about an operating system from two levels.

At the top is the user, or application, space. This is where the user applications are executed. Below the user space is the kernel space. Here, the Linux kernel exists.

There is also the GNU C Library (glibc). This provides the system call interface that connects to the kernel and provides the mechanism to transition between the user-space application and the kernel. This is important because the kernel and user application occupy different protected address spaces. And while each user-space process occupies its own virtual address space, the kernel occupies a single address space. For more information, see the links in the resources section.

The Linux kernel can be further divided into three gross levels. At the top is the system call interface, which implements the basic functions such as read and write. Below the system call interface is the kernel code, which can be more accurately defined as the architecture-independent kernel code. This code is common to all of the processor architectures supported by Linux. Below this is the architecture-dependent code, which forms what is more commonly called a BSP (Board Support Package). This code serves as the processor and platform-specific code for the given architecture.

Properties of the Linux kernel

When discussing the architecture of a large and complex system, you can view the system from many perspectives. One goal of an architectural decomposition is to provide a way to understand the source better and that's what we'll do here.

The Linux kernel implements a number of important architectural attributes. At a high level, and at lower levels, the kernel is layered into a number of distinct subsystems. Linux can also be considered monolithic because it lumps all of the basic services into the kernel. This differs from a microkernel architecture, where the kernel provides basic services such as communication, I/O, and memory and process management, and more specific services are plugged in to the microkernel layer. Each has its own advantages, but I'll steer clear of that debate.

Over time, the Linux kernel has become efficient in terms of both memory and CPU usage, as well as extremely stable. But the most interesting aspect of Linux, given its

buffer

a region of memory used to temporarily hold data while it is being moved from one place to another

VFS(Virtual File System)

an abstraction layer on top of a more concrete file system

size and complexity, is its portability. Linux can be compiled to run on a huge number of processors and platforms with different architectural constraints and needs. One example is the ability of Linux to run on a process with a memory management unit (MMU), as well as those that provide no MMU. The uClinux port of the Linux kernel provides for non-MMU support. See the resources section for more details.

Major subsystems of the Linux kernel

Now let's look at some of the major components of the Linux kernel using the breakdown.

System call interface

The SCI is a thin layer that provides the means to perform function calls from user space into the kernel. As discussed previously, this interface can be architecture dependent, even within the same processor family. The SCI is actually an interesting function-call multiplexing and demultiplexing service. You can find the SCI implementation in `./linux/kernel`, as well as architecture-dependent portions in `./linux/arch`. More details for this component are available in the resources section.

Process management

Process management is focused on the execution of processes. In the kernel, these are called threads and represent an individual virtualization of the processor (thread code, data, stack, and CPU registers). In user space, the term process is typically used, though the Linux implementation does not separate the two concepts (processes and threads). The kernel provides an application program interface (API) through the SCI to create a new process (`fork`, `exec`, or Portable Operating System Interface [POSIX] functions), stop a process (`kill`, `exit`), and communicate and synchronize between them (`signal`, or POSIX mechanisms).

Also in process management there is a need to share the CPU between the active threads. The kernel implements a novel scheduling algorithm that operates in constant time, regardless of the number of threads vying for the CPU. This is called the $O(1)$ scheduler, denoting that the same amount of time is taken to schedule one thread as it is to schedule many. The $O(1)$ scheduler also supports multiple processors (called Symmetric MultiProcessing, or SMP). You can find the process management sources in `./linux/kernel` and architecture-dependent sources in `./linux/arch`. You can learn more about this algorithm in the resources section.

Memory management

Another important resource that's managed by the kernel is memory. For efficiency, given the way that the hardware manages virtual memory, memory is managed in

Linux kernel

Unix-like

operating system kernel

kernel

the central component of most computer operating systems (OS). Its functions include managing the system's resources (the communication between hardware and software components

Minix

free/open source, Unix-like operating system (OS) based on a microkernel architecture

what are called pages (4KB in size for most architectures). Linux includes the means to manage the available memory, as well as the hardware mechanisms for physical and virtual mappings.

But memory management is much more than managing 4KB buffers. Linux provides abstractions over 4KB buffers, such as the slab allocator. This memory management scheme uses 4KB buffers as its base, but then allocates structures from within, keeping track of which pages are full, partially used, and empty. This allows the scheme to dynamically grow and shrink based on the needs of the greater system.

Supporting multiple users of memory, there are times when the available memory can be exhausted. For this reason, pages can be moved out of memory and onto the disk. This process is called swapping because the pages are swapped from memory onto the hard disk. You can find the memory management sources in `./linux/mm`.

Virtual file system

The virtual file system (VFS) is an interesting aspect of the Linux kernel because it provides a common interface abstraction for file systems. The VFS provides a switching layer between the SCI and the file systems supported by the kernel.

At the top of the VFS is a common API abstraction of functions such as open, close, read, and write. At the bottom of the VFS are the file system abstractions that define how the upper-layer functions are implemented. These are plug-ins for the given file system (of which over 50 exist). You can find the file system sources in `./linux/fs`.

GPL

a widely used free software license, originally written by Richard Stallman for the GNU project

Below the file system layer is the buffer cache, which provides a common set of functions to the file system layer (independent of any particular file system). This caching layer optimizes access to the physical devices by keeping data around for a short time (or speculatively read ahead so that the data is available when needed). Below the buffer cache are the device drivers, which implement the interface for the particular physical device.

Network stack

The network stack, by design, follows a layered architecture modeled after the protocols themselves. Recall that the Internet Protocol (IP) is the core network layer protocol that sits below the transport protocol (most commonly the Transmission Control Protocol, or TCP). Above TCP is the sockets layer, which is invoked through the SCI.

The sockets layer is the standard API to the networking subsystem and provides a user interface to a variety of networking protocols. From raw frame access to IP protocol data units (PDUs) and up to TCP and the User Datagram Protocol (UDP), the sockets

layer provides a standardized way to manage connections and move data between endpoints. You can find the networking sources in the kernel at `./linux/net`.

Device drivers

The vast majority of the source code in the Linux kernel exists in device drivers that make a particular hardware device usable. The Linux source tree provides a `drivers` subdirectory that is further divided by the various devices that are supported, such as Bluetooth, I2C, serial, and so on. You can find the device driver sources in `./linux/drivers`.

Architecture-dependent code

While much of Linux is independent of the architecture on which it runs, there are elements that must consider the architecture for normal operation and for efficiency. The `./linux/arch` subdirectory defines the architecture-dependent portion of the kernel source contained in a number of subdirectories that are specific to the architecture (collectively forming the BSP). For a typical desktop, the `i386` directory is used. Each architecture subdirectory contains a number of other subdirectories that focus on a particular aspect of the kernel, such as `boot`, `kernel`, `memory management`, and others. You can find the architecture-dependent code in `./linux/arch`.

Interesting features of the Linux kernel

If the portability and efficiency of the Linux kernel weren't enough, it provides some other features that could not be classified in the previous decomposition.

Linux, being a production operating system and open source, is a great test bed for new protocols and advancements of those protocols. Linux supports a large number of networking protocols, including the typical TCP/IP, and also extension for high-speed networking (greater than 1 Gigabit Ethernet [GbE] and 10 GbE). Linux also supports protocols such as the Stream Control Transmission Protocol (SCTP), which provides many advanced features above TCP (as a replacement transport level protocol).

Linux is also a dynamic kernel, supporting the addition and removal of software components on the fly. These are called dynamically loadable kernel modules, and they can be inserted at boot when they're needed (when a particular device is found requiring the module) or at any time by the user.

A recent advancement of Linux is its use as an operating system for other operating systems (called a hypervisor). Recently, a modification to the kernel was made called the Kernel-based Virtual Machine (KVM). This modification enabled a new interface to user space that allows other operating systems to run above the KVM-enabled kernel. In addition to running another instance of Linux, Microsoft® Windows® can

GNU

a computer operating system composed entirely of free software, initiated in 1984 by Richard Stallman

also be virtualized. The only constraint is that the underlying processor must support the new virtualization instructions. See the resource section for more information.

Going further

This article just scratched the surface of the Linux kernel architecture and its features and capabilities. You can check out the Documentation directory that is provided in every Linux distribution for detailed information about the contents of the kernel.

Resources

- The GNU site (<http://www.gnu.org/licenses>) describes the GNU GPL that covers the Linux kernel and most useful applications provided with it. Also described is a less restrictive form of the GPL called the Lesser GPL (LGPL).
- UNIX (<http://en.wikipedia.org/wiki/Unics>), MINIX (<http://en.wikipedia.org/wiki/Minix>) and Linux (<http://en.wikipedia.org/wiki/Linux>) are covered in Wikipedia, along with a detailed family tree of the operating systems.
- The GNU C Library (<http://www.gnu.org/software/libc/>), or glibc, is the implementation of the standard C library. It's used in the GNU/Linux operating system, as well as the GNU/Hurd (<http://directory.fsf.org/hurd.html>) microkernel operating system.
- uClinux (<http://www.uclinux.org/>) is a port of the Linux kernel that can execute on systems that lack an MMU. This allows the Linux kernel to run on very small embedded platforms, such as the Motorola DragonBall processor used in the PalmPilot Personal Digital Assistants (PDAs).
- "Kernel command using Linux system calls" (<http://www.ibm.com/developerworks/linux/library/l-system-calls/>) (developerWorks, March 2007) covers the SCI, which is an important layer in the Linux kernel, with user-space support from glibc that enables function calls between user space and the kernel.
- "Inside the Linux scheduler" (<http://www.ibm.com/developerworks/linux/library/l-scheduler/>) (developerWorks, June 2006) explores the new O(1) scheduler introduced in Linux 2.6 that is efficient, scales with a large number of processes (threads), and takes advantage of SMP systems.
- "Access the Linux kernel using the /proc filesystem" (<http://www.ibm.com/developerworks/linux/library/l-proc.html>) (developerWorks, March 2006) looks at the /proc file system, which is a virtual file system that provides a novel way for user-space applications to communicate with the kernel. This article demonstrates /proc, as well as loadable kernel modules.
- "Server clinic: Put virtual filesystems to work" (<http://www.ibm.com/developerworks/linux/library/l-sc12.html>) (developerWorks, April 2003) delves into the VFS layer that allows Linux to support a variety of different file systems through a common interface. This same interface is also used for other types of devices, such as sockets.

buffer cache

a collection of data duplicating original values stored elsewhere or computed earlier, where the original data is expensive to fetch (owing to longer access time) or to compute, compared to the cost of reading the cache

- "Inside the Linux boot process" (<http://www.ibm.com/developerworks/linux/library/l-linuxboot/index.html>) (developerWorks, May 2006) examines the Linux boot process, which takes care of bringing up a Linux system and is the same basic process whether you're booting from a hard disk, floppy, USB memory stick, or over the network.
- "Linux initial RAM disk (initrd) overview" (<http://www.ibm.com/developerworks/linux/library/l-initrd.html>) (developerWorks, July 2006) inspects the initial RAM disk, which isolates the boot process from the physical medium from which it's booting.
- "Better networking with SCTP" (<http://www.ibm.com/developerworks/linux/library/l-sctp/>) (developerWorks, February 2006) covers one of the most interesting networking protocols, Stream Control Transmission Protocol, which operates like TCP but adds a number of useful features such as messaging, multi-homing, and multi-streaming. Linux, like BSD, is a great operating system if you're interested in networking protocols.
- "Anatomy of the Linux slab allocator" (<http://www.ibm.com/developerworks/linux/library/l-linux-slab-allocator/>) (developerWorks, May 2007) covers one of the most interesting aspects of memory management in Linux, the slab allocator. This mechanism originated in SunOS, but it's found a friendly home inside the Linux kernel.
- "Virtual Linux" (<http://www.ibm.com/developerworks/linux/library/l-linuxvirt/>) (developerWorks, December 2006) shows how Linux can take advantage of processors with virtualization capabilities.
- "Linux and symmetric multiprocessing" (<http://www.ibm.com/developerworks/linux/library/l-linux-smp/>) (developerWorks, March 2007) discusses how Linux can also take advantage of processors that offer chip-level multiprocessing.
- "Discover the Linux Kernel Virtual Machine" (<http://www.ibm.com/developerworks/linux/library/l-linux-kvm/>) (developerWorks, April 2007) covers the recent introduction of virtualization into the kernel, which turns the Linux kernel into a hypervisor for other virtualized operating systems.
- Check out Tim's book GNU/Linux Application Programming (<http://www.charlesriver.com/Books/BookDetail.aspx?productID=91525>) for more information on programming Linux in user space.
- In the developerWorks Linux zone (<http://www.ibm.com/developerworks/linux/>), find more resources for Linux developers, including Linux tutorials (http://www.ibm.com/developerworks/views/linux/libraryview.jsp?type_by=Tutorials), as well as our readers' favorite Linux articles and tutorials (<http://www.ibm.com/developerworks/linux/library/l-top-10.html>) over the last month.
- Stay current with developerWorks technical events and Webcasts (http://www.ibm.com/developerworks/offers/techbriefings/?S_TACT=105AGX03&S_CMP=art).

Unix

a computer operating system originally developed in 1969 by a group of AT&T employees at Bell Labs including Ken Thompson, Dennis Ritchie and Douglas Ilroy