

An Introduction to Operating Systems

What is an operating system?

A piece of *systems software* that provides a convenient, efficient environment for the execution of user programs.

It's probably the largest and most complex program you'll ever run!

Why do we need an operating system?

- ✦ **The user's viewpoint:**
to provide the user interface, command interpreter, and directory structure, and to execute application programs (word processor, email client, web browser, MP3 player).
- ✦ **The programming environment viewpoint:**
to enhance the bare machine, to provide utility programs (such as compilers, editors, filters), to provide high-level input and output (I/O), to structure information into files, and to improve access to memory (size, protection, sharing).
- ✦ **The efficiency viewpoint:**
to replace the (long departed) human operator, to schedule tasks, to efficiently store and retrieve data, and to invoke and share programs.
- ✦ **The economic viewpoint:**
to allow simultaneous use and scheduling of resources, including disk-bound data and expensive peripherals.

Traditionally, we would summarize an operating system's goals as making "the system" *convenient* to use and scheduling its resources *efficiently* and *fairly*.

In addition, it must support hardware and software not yet developed.

Operating System ≠ User/Computer Interface

An operating system is often simply seen and described as the user/computer interface.

We often (mistakenly) claim to understand, and like or dislike, an "operating system" based on its interface.

Such an interface provides us with:

- ✦ program creation (editors, compilers, debuggers, linkers)
- ✦ program execution (character and graphical)
- ✦ access to I/O devices (both fixed and removable)
- ✦ constrained access to files of media
- ✦ constrained access to "internal" resources
- ✦ error detection, response, reporting, and
- ✦ accounting and monitoring.

Whether or not a certain interface runs on a particular hardware or operating system platform is usually dictated by economics, marketing, and politics - not technology.

Operating System \equiv Resource Manager

An operating system is better considered as being in control of its hardware and software resources.

Better still, because the "controls" are often temporal or external to the operating system itself, let's consider the operating system as a *resource manager*.

An operating system is *just another program* running on the available hardware.

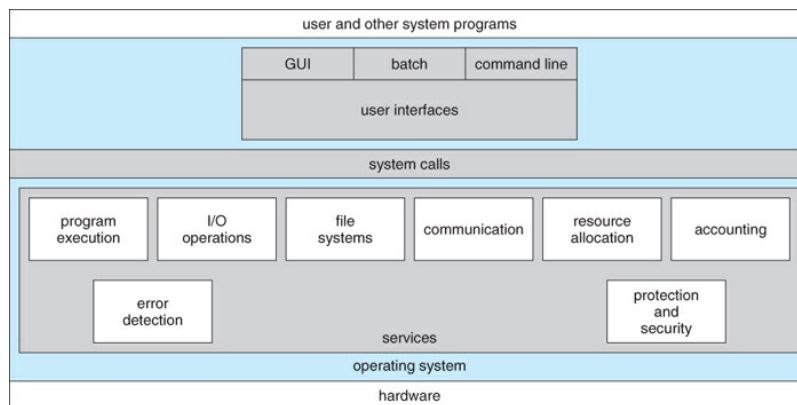
Most of the time, the operating system relinquishes control to the *user processes* until the hardware again dispatches the control.

Operating Systems Must Be Extensible

Of importance is an operating system's ability to *evolve* to meet new hardware and software demands:

- ✦ New hardware is constantly introduced - adding more memory presents little difficulty; new types of disks, video cards, etc, are more problematic.
- ✦ New application programs, tools, and system services are added.
- ✦ Fixes and patches are released to correct operating system *deficiencies*.

All of this suggests that the operating system, as a *program*, needs to be extensible - a *modular design* seems essential. Consider the following figure:



A view of operating system services

Of course, the above diagram provides a very simplified representation of operating systems and their services. In practice, the relationships between the modular components become very complex - [Linux](#) and [Windows-10](#)

and the [Interactive map of Linux kernel](#).

Traditional Operating System Services

CPU scheduling:

distribute or apportion computing time among several processes (or tasks) which appear to execute simultaneously.

Memory management:

divide and share physical memory among several processes.

Swapping:

move processes and their data between main memory and disk to present the illusion of a bigger machine.

I/O device support:

provide specialized code to optimally support device requirements.

File system:

organize mass storage (on disk) into files and directories.

Traditional Operating System Services, *continued*

Utility programs:

accounting, setting/constraining system resource access, manipulating the file system.

A command interface:

textual or graphical, to enable interactive interrogation and manipulation of operating system features.

System calls:

allow constrained access to the interior of the running operating system (as a program).

Protection:

keep processes from interfering with each other, their data, and "the system", whilst permitting sharing when requested.

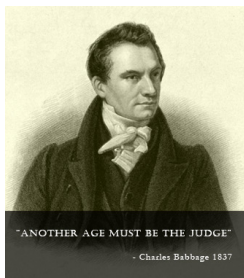
Communication:

allow users and processes to communicate within a single machine (inter-process communication), and across networks.

A Whirlwind History of Operating Systems

To understand the way modern operating systems are the way they are, it is useful to examine their evolution over the last almost eighty years.

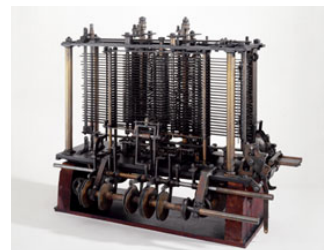
Advances in operating systems often accompanied advances in hardware, falling prices, and "exploding" capacities.



The first true digital computer was designed by English mathematician Charles Babbage (1792-1871).

Although Babbage spent most of his working life and fortune building his "analytical engine", its mechanical design and the wooden technology of the day could not provide the required precision.

Needless to say, the analytical engine did not have an operating system.



“ *Everything that can be invented has been invented.*

— Charles H. Duell, Commissioner, U.S. Office of Patents, 1899.

1945-55: Vacuum Tubes and Plugboards

Until World War II, little progress was made in constructing digital computers. Six significant groups can reasonably claim the first electrical computers:

- ▶ Tommy Flowers and Max Newman, Bletchley Park, England,
- ▶ Howard Aitken, Harvard,
- ▶ John von Neumann, Institute of Advance Studies, Princeton,
- ▶ Tom Kilburn and Freddie Williams, Manchester,
- ▶ J.P. Eckert and W. Mauchley, University of Pennsylvania, and
- ▶ Konran Zuse, Germany.

A single group of people designed, built, programmed, operated and maintained each machine. Although filling small warehouses, with tens of thousands of vacuum tubes, they were no match for today's cheapest home computers.

“ *I think there is a world market for maybe five computers.*

— Thomas Watson (1874-1956), Chairman of IBM, 1943.

Programs were loaded manually using console switches, or more likely by direct reconfiguration of wiring; indication of a program's execution and debugging returned through console lights.

Advantages:

- ▶ Interactive, and user received immediate response.

Disadvantages:

- ▶ Expensive machine was idle most of the time, because people's reactions (and thinking) were slow.
- ▶ Programming and debugging were tedious; hardware was very unreliable.
- ▶ Each program was self contained, including its own code for mathematical functions and I/O device support.

1955-65: Transistors and Batch Systems

Programming languages and operating systems (as we know them today) still unheard of. Collections of subroutines (procedures) to drive peripherals and to evaluate trigonometric functions were the first examples of operating systems services. The mid-1950s saw the user (still as the programmer) submitting a deck of punched Hollerith cards describing a job to be executed.

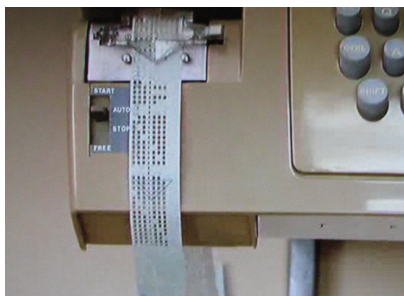


Given the high cost of computers, ways to increase their utility were quickly sought. The general solution was the *batch system*.

Similar/related programs, perhaps each requiring the FORTRAN (FORMula TRANslation) compiler, or a set of mathematical routines, were *batched* together so that the required routines need only be physically loaded once.



Programs were first written on paper tape, in the emerging FORTRAN language or in assembly language, and then copied to punched cards. Such decks of cards included *job control cards*, the program itself, and often the program's data.



Jobs submitted by different users were sequenced automatically by the operating system's *resident monitor*. Early peripherals, such as large magnetic tape drives, were used to *batch* input (jobs and data) and *spool* (from Simultaneous Peripheral Operation OnLine) output.

1955-65: Transistors and Batch Systems, *continued*

Generally, an inexpensive computer, such as an IBM 1401, was used for reading cards and printing from output tapes. The expensive machine, such as an IBM 7094, was used solely for the mathematical computations.

Advantages:

- ✦ The (true, computational) computer was kept busier.

Disadvantages:

- ✦ The computer was no longer interactive. Jobs experienced a longer *turnaround* time.
- ✦ The CPU was still idle much of the time for jobs. Other jobs remained *queued* for execution.

The significant operating system innovation at this time was the introduction of a *command interpreter* (a job control language - JCL) to describe, order, and commence execution of jobs.

The resident monitor was also protected from the user programs, and managed the automated loading of the programs after the monitor.

1965-1980: Integrated Circuits and Multiprogramming

In the early 1960s, computer manufacturers typically made two types of computers - word-oriented, large scale scientific computers (such as the IBM-7094), and character-oriented commercial computers (such as the IBM-1401), which were really better suited for I/O.

Incompatibility and a lack of an upgrade path were the problems of the day.

IBM attempted to address both problems with the release of their System/360, a *family* of software compatible machines differing only in capacity, price and performance. The machines had the same architecture and instruction set.

With heavy CPU-bound scientific calculations, I/O is infrequent, so the time spent (wasted) waiting was not significant. However, commercial processing programs in the emerging COBOL (Computer Oriented Business Organizational Language) often spent 80-90% of its time waiting for I/O to complete.

The advent of separate I/O processors made simultaneous I/O and CPU execution possible.

The CPU was *multiplexed* (shared), or employed *multiprogramming*, amongst a number of jobs - while one job was waiting for I/O from comparatively slow I/O devices (such as a keyboard or tape), another job could use the CPU.

Jobs would run until their completion or until they made an I/O request.

Advantages:

- Interactivity was restored.
- The CPU was kept busy if enough jobs were ready to run.

Disadvantages:

- The computer hardware and the operating system software became significantly more complex (and there has been no looking back since!).

1965-1980: Integrated Circuits and Multiprogramming, *continued*

Still, the desire for quicker response times inspired a variant of multiprogramming in which each user communicated directly with one of a multitude of I/O devices.

The introduction of *timesharing* introduced *pre-emptive scheduling*. Jobs would execute for at most pre-defined time interval, after which it was another job's turn to use the CPU.

The first serious timesharing system (CTSS, from MIT 1962) lacked adequate memory protection.

Most (modern) operating system complexity was first introduced with the support of multiprogramming - scheduling algorithms, deadlock prevention, memory protection, and memory management.

The world's first commercially available time-sharing computer, the DEC PDP-6, was installed in UWA's Physics Building in 1965 - cf.

[Cyberhistory](#), by Keith Falloon, UWA MSc thesis, 2001, and [pdp6-serials](#).

Early Operating System Security



Very early operating systems, on one-user-at-a-time computer systems, assisted the user to load their programs and commence their execution.

There was little to protect and, if an errant program modified the *executive* program, then only the current user was affected. As more of a courtesy, the executive might clear memory segments and check itself before accepting the next program.

As multi-tasking operating systems emerged, accountability of resource use became necessary, and operating system *monitors* oversaw the execution of programs.

As with modern operating systems, there was the need to:

- protect the operating system from the program,
- protect programs from themselves,
- protect programs from each other, and
- constrain data access to the correct program(s).

Until system resources became more plentiful (and cheaper) attempts were made to maximize *resource sharing* - security was a consequent, not an initial, goal.

e.g. process scheduling policies were dominated by already running processes requesting resources (such as libraries and tapes) that were already in use.

At this level, computer security is more concerned with reliability and correctness, than about deliberate attacks on programs, data, and the system.

1970s: Minicomputers and Microcomputers

“ *There are only two things to come out of Berkeley, Unix and LSD, and I don't think this is a coincidence.*

— *Jeremy S. Anderson.*

Another major development occurring in parallel was the phenomenal growth in minicomputers, starting with the DEC (Digital Equipment Corporation) PDP-1 (Programmed Data Processor) in 1961. The PDP-1, with 4K of 18-bit words cost only US\$120,000 - 5% of the IBM 7094.

The trend was towards many small mid-range personal computers, rather than a single mainframe.

Early minicomputers and microcomputers were simple in their hardware architectures, and so there was some regression to earlier operating system ideas (single user, no pre-emption, no multiprogramming).

For example, MS-DOS on an IBM-PC (circa. 1975) was essentially introduced as a batch system, similar to those of the 1960s, with a few modern additions, such as a hierarchical file system.

With some notable exceptions, the trend quickly moved towards support of all modern operating system facilities on microcomputers.

“ *There is no reason anyone would want a computer in their home.*

— *Ken Olsen, DEC Founder and Chairman, 1977.*

Perhaps most significant has been the evolution, and importance, of operating systems' user interfaces.

In particular, the graphical *desktop metaphor* has remained for some time.

1980-90s: Personal Computers and Networking

“ 640K ought to be enough for anybody.

— Bill Gates (1955-), in 1981.

The decentralization of computing resources, now *data* and not the hardware, required more support for inter-operating system communication - both physical support and application program support.

As minicomputers shrunk in size, but exploded in capacity, the powerful computer *workstation* was born. Companies such as Sun Microsystems (SUN) and Silicon Graphics (SGI) rode this wave of success.

Local-area networks (primarily Ethernet and token-ring) connected workstations, while wide-area networks connected minicomputers.

Operating system developments included the development of fast and efficient network communication protocols, data encryption (of networks and file systems), security, reliability, and consistency of distributed data.

2000s and Beyond: Speeds, Capacities, Mobility, and Ubiquity

- ✦ The *desktop metaphor* for computer interfaces becomes only an instance of the *widerweb metaphor*.
- ✦ Hardware prices drop, and capacities explode.
 - [Amazing Facts and Figures About the Evolution of Hard Disk Drives](#)
 - [An improved chart of SSD vs HDD historical and projected prices](#)
- ✦ High speed, long distance communication links encourage graphical and audio communication (surpassing text).
 - [TV and Video Will Triple Average Home Monthly Internet Usage to Beyond 1 TB By 2025](#)
- ✦ As CPU performance no longer keeps pace with [Moore's Law](#), and computing becomes increasingly mobile, focus shifts to battery life and metrics of performance/Watt.
 - [Moore's Law and Its Practical Implications](#)
 - [Performance per Watt Is the New Moore's Law](#)
 - [How Much Power Do Computers Consume?](#)

“ *For years, we thought that a million monkeys sitting at a million keyboards would produce the complete works of Shakespeare. Today, thanks to the Internet, we know that's not true.*

— Anon.