

Processor controls hardware and executes instructions.
Memory stores data and code.
I/O components read and write from I/O devices.
System bus interconnects different components and provides communication between them.
OS must make **efficient** use of available resources and must **share resources** (CPU cores, caches, RAM, GPUs, clocks, timers, persistent storage) among multiple users.
OS must support **multiple parallel activities** to convert raw hardware into a usable computer system by keeping details hidden from user-level software.
Must support several simultaneous parallel activities, and may switch activities at arbitrary times, so must ensure **safe concurrency**.

MULTIPROCESSORS (Windows, macOS, Linux) have many CPU cores and CPUs.

SERVER OS (Solaris, FreeBSD, Linux) shares hardware/

MAINFRAMES have bespoke hardware, limited workload.

SMARTPHONES have power-efficient CPUs.

EMBEDDED OS (QNX, VxWorks) are used for home utilities, and only run trusted software.

REAL-TIME OSs are time-oriented (not performance/I/O).

SENSOR NETWORK OS (Tiny OS) is resource efficient.

The **OS kernel** is a privileged program that implements OS functionality.

Monolithic kernels are a single black box with all functionality. Only executable with its own address space. **Microkernels** have a layered design, with protection between them.

Microkernels have as little functionality in the kernel as possible, functionality passed to a user-level server. Kernel does IPC between servers - separate ones for I/O, file access, memory management, etc. leads to overhead. Less error-prone and servers have clean interfaces.

Hybrid kernels combine features of both.

Linux is a variant of Unix. **Interrupt handlers** are the primary means to interact with devices. I/O scheduler orders disk operations. Supports static in-kernel components and dynamically loadable modules.

NTOS practices **system calls** in **WINDOWS**. Programs run on top of **Windows API**. NTOS loaded from **ntoskrnl.exe** at boot, consists of **executive** and **kernel** layers. **Device drivers** loaded into kernel memory. NT structured like a **microkernel**, most components run in same address space. **Hardware Abstraction Layer** abstracts out BIOS config, CPU types, etc...

PROCESSES allow a single processor to run multiple programs "simultaneously". Provides concurrency, isolation between programs, allows better resource utilization, simplifies programming.

Pseudo concurrency single physical processor is switched between processes by interleaving.

Real concurrency utilises multiple physical processors.

On a **context switch**, the processor switches from executing process A to executing process B. May occur in response to events/interrupts, or as a result of periodic rescheduling. Process can **switch later** to do some other information needed should be stored in **process descriptor** / control block, kept in the **process table**.

Process Control Block — process has its own VM, should store PC, page table register, stack pointer, ..., process management info, file management info.

Context switches are expensive, so avoid unnecessary switching by changing process state. **Indirect calls** — perturbation of memory caches.

^ Flushing TLB & pipeline flushing

UNIX allows processes to create hierarchies. **Windows** has no notion of this.

int fork(void) Creates new child process, to be controlled by parent by **fork() + exec()**. Same resources of parent process image. Returns child **PID** in parent process and 0 in child process. No child created, and -1 returned to parent on error.

execve(char *path, char **argv, char **envp) — Changes process image and runs new process. Lots of useful wrappers.

int waitpid(int pid, int *status, int options) — Suspends execution of calling process until process **pid** terminates normally or signal received. Can wait for more than one child.

Simple design philosophy, so process API made up of basic blocks that can be easily combined.

CreateProcess() (equivalent of **fork()** & **exec()**) has 10 parameters!

Termination — normally process completes execution of body, **call exit()**. **Abnormal exit** — runs into error or received exception. **Aborted** — another process has overruled execution. Some run in an endless loop.

void exit(int status) — Terminates a process. Never returns in the calling process.

int kill(int pid, int sig) — sends signal to process **pid**.

SIGNALS are an IPC mechanism

SIGABRT Abort signal from assert

SIGFPE Floating point exception

SIGILL Kill signal similar to delivery of hardware interrupts.

SIGKILL Kill memory reference

SIGPIPE Broken pipe; write to pipe with no reader

SIGPOLL Timer signal from alarm

SIGPWR Termination signal sent to another process, if it has permission. Generated when an exception occurs, or when a process receives an event, when certain key combinations typed in terminal, or programmatically using **setitimer()**.

Default action for most signals is to terminate process, but receiving process may ignore/handle it instead (unless **SIGKILL** or **SIGSTOP**).

Pipes — method of connecting stdout of one process to stdin of another. **Named** or **unnamed**. Sender should close the read end, receiver should close the write end.

Named Pipes / **FileDescriptors** are persistent pipes that survive processes which created them. Stored on file system. Any process can open it like a regular file.

int pipe(int fd[2]) — returns two file descriptors: **read end(fd[0])**, **write end(fd[1])**.

THREADS are execution streams that share the same **address space**. **Multithreading** — each process can contain multiple threads.

Useful for applications which contain multiple parallel activities that share the same data. Processes are too heavyweight. It is difficult to coordinate between them & create/destroy activities.

Shared address space can lead to memory corruption and concurrency bugs. How to handle **forking** and **signals**?

PTHREADS implemented by most Unix systems.

int pthread_create(...) — Creates a new thread.

void pthread_exit(void *value_ptr) — Terminates the thread and makes **value_ptr** available to any successful join with terminating thread. Called implicitly when the

thread's start routine returns.

int pthread_yield(void) — Releases CPU to let another thread run. Returns 0 on success (always), or error code.

int pthread_join(...) — Blocks until thread terminates.

User-level threads — OS kernel not aware of threads, each thread manages its own threads. Better performance, and each application can have its own scheduling algorithm. But blocking system calls (e.g. during page fault) stop **all threads** in process.

Hybrid approaches use kernel threads and multiplex user-level threads onto some of these.

SCHEDULING

— If multiple processes are ready, which one should be run?



Goals of Scheduling Algorithms => Ensure fairness:

Comparable processes should get comparable services; Avoid indefinite postponement: No process should starve. **Enforce policy**: E.g. priorities; **Maximize resource utilization**: CPU, I/O devices; **Minimize overhead**: From context switches, decision-making decisions

For **Batch Systems** => **Throughput**: Jobs per unit of time; **Turnaround time**: Time Between job submission & completion

For **Interactive Systems** => **Response time**: Time between request issued and first response

Non-preemptive scheduling lets process run until it blocks / releases CPU. Preemptive scheduling lets process run for a maximum amount of fixed time.

FIRST-COME FIRST-SERVE — No indefinite postponement. Easy implementation. But suffers from "head-of-line blocking" if long job followed by many short jobs.

ROUND-ROBIN — Processor runs until it blocks or time quantum expires, and goes to ready queue. Poor turnaround time when jobs have similar run-times.

^ Good avg. turnaround time when different runtimes. RR quantum (time slice) => should be larger than context switch cost, but still provide decent response times. Turnaround = $\frac{\text{Total Job Execution Time}}{\text{RR quantum}} + \text{RR quantum} \cdot \text{ctx. switch time}$.

PRIORITY SCHEDULING — jobs run based on their priority. Priority can be defined based on some process-specific metrics. Can be static or dynamic.

Generally, favour short and I/O-bound jobs. Quickly determine the nature of job and adapt to changes.

MULTILEVEL FEEDBACK QUEUES — Form of priority scheduling. One queue for each priority level. Within each queue, use round robin scheduling (usually RR). Need to determine current nature of job. Need to prevent starvation of lower-priority jobs. **Feedback** comes from periodically recomputing priorities: e.g. based on how much CPU they used recently, increase job's priority. Problem: starvation. Solution: weighted moving average. Not very flexible, does not react quickly to changes, cheating is a concern.

LOTTERY SCHEDULING — Jobs receive lottery tickets for various resources — random ticket chosen at each scheduling decision. Number of lottery tickets is meaningful as a result. Highly responsive with no starvation. Jobs can exchange tickets. But unpredictable response time.