

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.

Each application has its own **context** — thread of execution.

Hardware needs to convert raw hardware into a usable computer system by keeping drivers hidden from user-level software.

Must support several simultaneous parallel activities, and switch activities at arbitrary times, so must ensure **safety**.

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 OS are time-oriented (not performance-oriented).

SENSOR NETWORKS (OS Tiny) 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.

Kernel update requires full system shutdown, but modular design is possible, 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 overhead 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 the kernel. 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, but better resource utilization simplifies scheduling.

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 **process creation** or **termination** later on. So no information needed about who is sleeping 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 context switching — move to a higher process state. **Indirect context switching** — context switch cost is proportional to perturbation of memory caches.

Flush TLB & pipeline flushing.

UNIX allows processes to create hierarchies. Windows has no notion of child processes, to be created and destroyed by making exec/clone system call of parent process image. Returns child PID in parent process and 0 in child process. No child created, and I returned to error.

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

How do processes **Synchronise** their operations to a task?

Critical section — processes access a shared resource. **Synchronisation mechanism** required at entry/exit of critical section.

Race condition occurs when multiple threads read and write shared data within critical section => their execution interleaved non-deterministically.

Assume **sequential consistency** throughout — operations of each thread executed in program order, operations of all threads executed in sequential order automatically.

Operations in the calling process. Never terminates.

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

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

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

Termination — normally process completes execution of body, *call(RT)*. After exit, it runs into error, or handled 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(pid, int sig) — returns two file descriptors: read end fd[0], write 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 data. Processes are too heavyweight. It is difficult to coordinate between address spaces, and expensive to context switch between them & create/destruction.

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

PTHREADS implemented by most Unix systems.

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

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

int pthread_cancel(...) — Cancels unstarted threads.

OS-level threads — OS kernel is 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) stops **all threads** in process.

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

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, Memory, overcommit. From context switching, scheduling decisions.

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

For Real-time Systems => Response time: Time between request issued and first response.

Non-preemptive scheduling lets processes run until it blocks / releases CPU. Preemptive scheduling lets processes 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 — processes轮流 blocks or time slices. Good for fairness, but gives 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, but still provide decent response times. Larger quantum => better performance. RR quantum & RR quantum + cts. switch time).

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 the kernel. 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, but better resource utilization simplifies scheduling.

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 **process creation** or **termination** later on. So no information needed about who is sleeping 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 context switching — move to a higher process state. **Indirect context switching** — context switch cost is proportional to perturbation of memory caches.

Flush TLB & pipeline flushing.

UNIX allows processes to create hierarchies. Windows has no notion of child processes, to be created and destroyed by making exec/clone system call of parent process image. Returns child PID in parent process and 0 in child process. No child created, and I returned to error.

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

How do processes **Synchronise** their operations to a task?

Critical section — processes access a shared resource. **Synchronisation mechanism** required at entry/exit of critical section.

Race condition occurs when multiple threads read and write shared data within critical section => their execution interleaved non-deterministically.

Assume **sequential consistency** throughout — operations of each thread executed in program order, operations of all threads executed in sequential order automatically.

Operations in the calling process. Never terminates.

int kill(pid, int sig) — returns two file descriptors: read end fd[0], write 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 data. Processes are too heavyweight. It is difficult to coordinate between address spaces, and expensive to context switch between them & create/destruction.

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

PTHREADS implemented by most Unix systems.

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

*void pthread_exit(void *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.

Mutual exclusion ensures if a process executes its critical section, no other process can be executing it. Two processes may be simultaneously inside a critical section.

2. No process running outside critical section.

3. No process requiring access to its critical section forever.

4. Operations made about the relative speed of processes.

A commonly desired feature of synchronisation mechanisms.

Example implementations => disabling interrupts, Peterson's solution, spin locks, binary semaphores w/ lock counter.

LOCK/MUTEX => mutual exclusion mechanism owned by actively excluding thread; lock is usually user/process object while mutex is same thing but OS/kernel object.

SMONOPHORE => the amount of data a lock is protecting.

Lock overhead => cost of using locks, e.g. memory space, initialisation, time to acquire/release.

Lock contention => number of processes waiting for lock; more contention => less parallelism.

Read/Write Locks can be used to minimize lock contention & maximize concurrency => exclusive access mode & non-blocking mode so multiple threads can acquire in read mode.

SEMAPHORES — processes cooperate by means of signals. Special variables, accessible via (atomic) *down(s)*, *up(s)*, *init(s)*, *d*. Two private components = a counter and a queue of processes waiting on the semaphore. In a binary semaphore, counter initialised to -1.

External fragmentation => total memory exists to satisfy request, but not contiguous. Reduced by **compaction** (place all free memory in one large block), **defragmentation** (relocate pages in memory).

Internal fragmentation = allocated memory larger than requested memory.

Swapping process memory temporarily out of memory to disk => transfer time is major part of swap time.

VIRTUAL MEMORY — Separation of user logical memory from physical memory. Each part of memory for execution, logical address space much larger than physical, can be shared by several processes.

PAGING — logical address space of process can be non-contiguous. **Frames** = fixed-size blocks of physical memory. **Pages** = fixed-size blocks of logical address space.

First-fit => best-fit: allocate first hole big enough => **Best-fit**: allocate smallest hole big enough => **Worst-fit**: allocate largest hole.

Multiple-partition Allocation:

Hole is a block of available memory. When new process arrives, allocate memory from large-enough hole. OS maintains info on allocated & free partitions (holes).

SEMANTICALLY — processes cooperate by means of signals. Special variables, accessible via (atomic) *down(s)*, *up(s)*, *init(s)*, *d*. Two private components = a counter and a queue of processes waiting on the semaphore. In a binary semaphore, counter initialised to -1.

External fragmentation => total memory exists to satisfy request, but not contiguous. Reduced by **compaction** (place all free memory in one large block), **defragmentation** (relocate pages in memory).

Internal fragmentation = allocated memory larger than requested memory.

Swapping process memory temporarily out of memory to disk => transfer time is major part of swap time.

VIRTUAL MEMORY — Separation of user logical memory from physical memory. Each part of memory for execution, logical address space much larger than physical, can be shared by several processes.

Counting algorithm keep a count of references to a page. If page is not used for a long time, it is freed.

LRU (Least Recently Used) — page replaced least recently used => never forgets heavy page usage.

LFU (Least Frequently Used) — replaces least count => accounts for small-count pages being new.

Locality of Reference — Programs tend to request same pages in space/time. If program's favoured subset of pages in space/time, we can use **locality** for **swapping** — excessive paging activity causing low CPU utilization.

Working set of pages $W[t]$ is the set of pages referenced by a process in time interval $[t-w, t]$.

WS Clock Algorithm adds a time of fast to use the second-chance algorithm => On each page fault => if age < w then if the page is clean, replace, otherwise trigger a write back and move to next page] => if age >= w then if the page is clean, replace, otherwise trigger a write back and move to next page] => can choose w based on page fault frequency => if too large then working set won't fit into memory.

Processes transition between working sets, so temporarily maintains in-memory pages outside of working set size.

Local page replacement — each process gets fixed allocation of physical memory. Need to pick up changes in working set size.

Global page replacement — dynamically share memory between runnable processes. Initially allocate memory proportional to process size.

Linux uses variation of clock algorithm to approximate LRU page-replacement strategy. MMU uses two linked lists.

Address Translation: One-Level Paging

For a given logical address 2^m and page size 2^n

2^m bits of address space, where each address refers to a byte in memory.

2^n bytes of page size, spanned by 2^n bits of address space.

First m MSBs of logical address is **page number** p => used to **index into page table** $PT[p]$ => to get **frame address** f

entry procedure **insert(item)**

while (count == N) wait (not_full);

insert _item (item); count++;

signal (not_empty);

entry procedure **remove(item)**

while (count == 0) wait (not_empty);

remove _item (item); count--;

signal (not_full);

end monitor;

producer inserts & consumer removes items w/ max. capacity N => only insert if not-full & remove if full => only insert if non-full & remove if empty.

Set of processes is **DEADLOCKED** if each process is waiting for an event no other process can cause.

Resource deadlock the most common:

1. Mutual exclusion,

2. Hold and wait (process can request resources while it already holds others earlier),

3. Non-preemption,

4. Circular wait (set of processes in circular waiting chain).

How to deal with deadlocks?

Ignore it — If contention for resources is low, then deadlocks are infrequent.

Detection & Recovery:

After system deadlocked, dynamically build resource ownership graph and look for cycles => e.g. DFS search to look for cycles: when an edge has been inspected, it is marked and not visited again.

Recover from deadlock by preemption (temporarily take resource from owner and give to another) => rollback (periodic checkpoints are made for process context switch will random process ...))

SIGKILL or **SIGSTOP** — Default action for most signals, fails to meet rules of mutual exclusion => if T_0 's non-crit. section is long enough, it will prevent T_1 from entering crit. section. \uparrow \uparrow

T1

while (true) {

while (turn != 0) {

turn++;

critical_section();

turn = 0;

noncritical_section0(); noncritical_section1(); }

T0

while (true) {

while (turn != 1) {

turn++;

critical_section();

turn = 0;

noncritical_section0(); noncritical_section1(); }

PETERSON'S SOLUTION fixes issues w/ strict alternation: thread 1 calls enter/leave_critical() & thread 2 calls enter/leave_critical(). \uparrow \uparrow

alternation: thread 1 calls enter/leave_critical() & thread 2 calls enter/leave_critical(). \uparrow \uparrow

interested[2] = (0, 0)

if (interested[0] == 0) {

// thread is 0 or 1

interested[0] = 1;

interested[1] = 0;

interested[2] = (1, 0)

if (interested[0] == 1) {

int other = 1 - turn;

while (turn == other) {

turn++;

critical_section();

turn = 0;

noncritical_section0(); noncritical_section1(); }

PETERSON'S SOLUTION fixes issues w/ strict alternation: thread 1 calls enter/leave_critical() & thread 2 calls enter/leave_critical(). \uparrow \uparrow

alternation: thread 1 calls enter/leave_critical() & thread 2 calls enter/leave_critical(). \uparrow \uparrow

interested[2] = (0, 0)

if (interested[0] == 0) {

// thread is 0 or 1

interested[0] = 1;

interested[1] = 0;

interested[2] = (1, 0)

if (interested[0] == 1) {

int other = 1 - turn;

while (turn == other) {

turn++;