

Computer Systems - CheatSheet

IN BA4 - Jean Cédric Chappelier

Notes by Ali EL AZDI

A Computer Systems Cheatsheet has been authorized for the upcoming exam, and I'm sharing a copy of mine for anyone interested. It provides a concise summary of the key concepts and techniques covered in the course. For any updates or suggestions, feel free to reach out to me on Telegram at [elazdi_al](https://t.me/elazdi_al) or via EPFL email at ali.elazdi@epfl.ch.

April 3rd, 2025

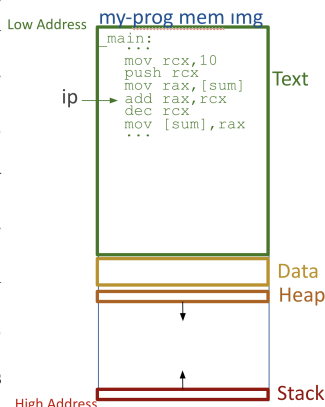
Program. Passive entity. A sequence of instructions stored in a file, not currently executing.

- **Storage:** Stored on persistent storage (e.g., disk). Loadable into main memory by the OS.
- **Access:** Static file. Contains instructions, data, and metadata used during execution.

Process. Active execution of a program. Managed by the OS. Each process has its own isolated virtual address space.

- **Storage:** Occupies main memory during execution.
- **Memory Image:** Text (code), data, heap, and one stack per thread.
- **Process ID (PID),** assigned and tracked by the OS, unique system-wide.
- **Status:** OS tracks current status (e.g., running, waiting, terminated).
- **Virtualization:** Each process has the illusion of exclusive access to memory.
- **OS-allocated Resources:** File descriptors, sockets, I/O handles, etc.
- **Pointer to Page Table.** The **Kernel** tracks the process-specific page table pointer.
- **Access:** Operates in isolated virtual memory. No direct access to other processes. Interacts with hardware via system calls. OS handles scheduling and resource management.

If a process has a single thread, we may refer to the thread's CPU context as the process's CPU context.



Limited Direct Execution. OS design allowing user programs to execute instructions directly on the CPU, with restrictions.

- **Goal:** Maximize performance while maintaining control and protection.
- **User Code Execution:** CPU runs user programs natively (not emulated) in **user mode**.
- **OS Control:** OS retains control over hardware via **privileged instructions** and mode switching.

User Mode vs Kernel Mode. Two CPU execution modes stored as a state bit in a protected CPU register (0 = user mode, 1 = kernel mode) controlling access to hardware and instructions.

- **User Mode:** Restricted. User code cannot execute privileged instructions or directly access hardware/memory management.
- **Kernel Mode:** Mode in which is ran a central part of the operating system, the **Kernel**.
- **Switching:** A transition from user mode to kernel mode is triggered by:
 - System Calls (e.g., file I/O, memory allocation)
 - Hardware Interrupts (keyboard, timer,...)
 - Software Traps (divide-by-zero, invalid memory access),
If an exception happens in kernel mode, an internal routine is called but no mode switch occurs.

After handling the event, the OS switches the CPU back to user mode to resume application execution.

Privileged Instructions. CPU instructions that can only execute in **kernel mode**. If attempted in user mode, the CPU raises a trap to the OS. Prevent user programs from interfering with other processes or the OS.

To do something that requires high privilege, a thread makes a syscall, invoking the kernel.

Exception/trap - synchronous signal. The CPU itself raises an exception/trap.

Interrupt - asynchronous signal. Some external entity raises an interrupt. If an external entity needs attention, it raises an **interrupt**, invokes the kernel.

Even if nothing external needs attention, the timer interrupt regularly invokes the kernel to run the OS scheduler.

fork() — Clone the current process

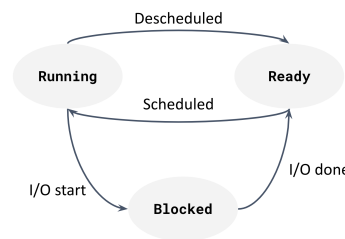
1. The operating system creates a new process (child) by duplicating the calling process.
2. The child initially shares all memory pages with the parent:
 - Pages are marked read-only and shared (Copy-On-Write).
 - If either process attempts to write to a page, it is copied privately for that process.
3. The child also receives:
 - Duplicated file descriptors (pointing to the same open files).
 - Identical program counter and stack pointer.
4. Returns:
 - In the parent process: the PID of the child.
 - In the child process: 0.
5. Both processes resume execution at the instruction following **fork()**.

exec() — Mutate the child into a new program

- Typically called by the child process immediately after **fork()**.
- Replaces the child's memory space with a new program image.
- Stack, heap, code, and data segments are replaced.
- If successful: the new program starts execution from **main()**.
- If failed: **exec()** returns -1, child continues old code.

Thread. The smallest unit of CPU execution within a process. Multiple threads can exist within one process.

- **Memory:** Shares the parent process's virtual address space:
 - Shared: Code (text), heap, global/static data, open files.
 - Private: Each thread has its own stack (for local variables and function calls).
- **Execution context (per thread), values of all the CPU registers at the last moment the thread was running:**
 - **Instruction Pointer (IP) / Program Counter (PC):** Physical register pointing to the next instruction. Private per thread.
 - **Stack Pointer (SP):** Points to the top of the thread's private stack.
 - **General-purpose Registers:** Include temporary registers (e.g., RAX, RBX), status, and flags. **All private per thread.**
 - **Thread ID (TID).** The Thread ID is unique within a process, but not necessarily system-wide, assigned and tracked by the OS.
 - **Status:** OS tracks each thread's status (Running, Ready, Blocked).
 - **Virtualization:** Each thread has the illusion it exclusively occupies the CPU.
 - **Context switching:** OS saves/restores full register set (IP, SP, general-purpose regs). The CPU switches from one thread to another. Thread switches are faster than process switches. It saves to memory the CPU context of the current thread; it restores from memory the CPU context of the new thread.
 - **Thread management:** A main thread may create and manage others. In some models, a dedicated manager thread exists.
 - **Managed by OS:** The OS kernel schedules threads individually and tracks them.



OS Components.

- **Kernel (core of the OS):** Loaded at boot, runs in **kernel mode**, and manages hardware, memory, processes, and system calls. It is the only part of the OS that runs in kernel mode.
Three important routines inside of the kernel:
 - **Interrupt Handler or Interrupt Service Routine (ISR):** Handles interrupts from hardware devices.
 - **Exception / Trap Handler:** Handles exceptions and traps from user programs.
 - **System Call Handler:** Handles system calls from user programs.The three routines use two tables (**Trap/Interrupt Table** and **System Call Table**) to resolve the event type to a specific memory address of the handler routine managing the event.
- **Loader (part of the OS):** Prepares executables to run:
 - Loads program code/data into memory.
 - Maps required libraries (e.g., libc).
 - Sets up the process stack and heap.
 - Places command-line arguments and environment variables on the stack.
 - Sets the **%rsp** (stack pointer) and **%rdi** (argc), **%rsi** (argv) registers so the program can access arguments.
 - *Because the kernel manages processes and initiates execution, it is responsible for placing arguments in registers so the loader (running in user mode) can access them.*
 - Jumps to the program's entry point to begin execution in **user mode**.
- **User-level Programs:** Applications like shells, editors, browsers, etc. These are **part of the OS** but run entirely in **user mode**, relying on system calls to request kernel services.
- **System Libraries:** Shared libraries (e.g., libc) used by user programs. Loaded and linked by the loader, but executed in user mode. Provide wrappers around system calls.

fork() — Clone the current process

1. The operating system creates a new process (child) by duplicating the calling process.
2. The child initially shares all memory pages with the parent:
 - Pages are marked read-only and shared (Copy-On-Write).
 - If either process attempts to write to a page, it is copied privately for that process.
3. The child also receives:
 - Duplicated file descriptors (pointing to the same open files).
 - Identical program counter and stack pointer.
4. Returns:
 - In the parent process: the PID of the child.
 - In the child process: 0.
5. Both processes resume execution at the instruction following **fork()**.

exec() — Mutate the child into a new program

- Typically called by the child process immediately after **fork()**.
- Replaces the child's memory space with a new program image.
- Stack, heap, code, and data segments are replaced.
- If successful: the new program starts execution from **main()**.
- If failed: **exec()** returns -1, child continues old code.

wait() — Wait for a child process to terminate

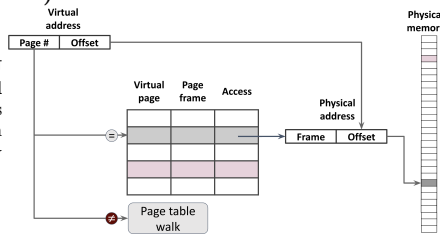
- Frees system resources allocated to the process.
- Sends termination status to the parent.
- If main function of a program returns, **exit()** is implicitly called.
- Blocks the calling process until one of its child processes terminates.
- Returns the PID of the terminated child.
- If no child processes exist, it waits indefinitely.

Memory Management Unit (MMU) Address Translation.

Hardware component managed by the operating system that translates virtual addresses into physical addresses using per-process page tables. Each process has its own page table, and the currently active one is referenced in a special CPU register.

Translation Lookaside Buffer (TLB).

A small, fast cache inside the MMU that stores recent **virtual-to-physical page mappings**. The TLB reduces the need to access the page table on every memory reference, significantly speeding up address translation.



Virtual Address Translation:

1. The CPU executes an instruction (e.g., load, store, or fetch) that references a memory location using a **virtual address**.
2. The MMU extracts the **virtual page number (VPN)** and **page offset** from the virtual address.
3. The MMU first checks the **TLB** for a cached translation of the VPN.
4. **If the translation is found in the TLB (TLB hit):**
 - (a) The corresponding **physical frame number** is retrieved directly.
5. **If the translation is not found (TLB miss → page table walk):**
 - (a) The MMU uses the active **page table** to look up the VPN and obtain the PFN.
 - (b) The new mapping is inserted into the TLB for future accesses.
6. **If the page table entry is invalid or present = 0 (page fault):**
 - (a) The CPU triggers a **trap into the OS kernel**.
 - (b) The **page fault handler** in the OS's memory management subsystem is invoked.
 - (c) The handler checks whether the page is:
 - i. **Never allocated before** → allocate a new physical page and zero-initialize it.
 - ii. **Swapped out to disk** →
 - A. Check the **swap cache** for the page.
 - B. If found in the swap cache, use the cached page directly.
 - C. If not in the cache, read the page from the **swap space** on disk into a free physical frame, and insert it into the swap cache.
 - (d) The newly loaded or allocated page is mapped in the page table, and **present** is set to 1.
 - (e) The TLB entry for the VPN is updated if needed.
 - (f) The faulting instruction is retried.
7. The final **physical address** is then used to **access RAM**.

Swapping

Transfer of memory pages between physical RAM and disk-based swap space to free up RAM.

- **Swap Space** Disk region reserved for evicted pages; holds non-resident memory to extend usable RAM.

- **Swap Cache** In-memory buffer of recently swapped-out pages; enables fast lookup and avoids redundant disk I/O during swap-in.

-Swap-Outs:

1. The OS detects that the number of free physical pages has fallen below a predefined threshold.
2. A background kernel thread is triggered to reclaim memory.
3. The kernel scans memory to identify candidate pages for eviction, typically using an aging or Least Recently Used algorithm.
4. A candidate page is selected if:
 - (a) It is not currently in use (i.e., not recently accessed),
 - (b) It is not locked, pinned, or shared with kernel-critical structures.
5. If the selected page is dirty (i.e., has been modified), its contents are written to the swap space on disk.
6. The page is inserted into the **swap cache** so it can be quickly retrieved if needed again.
7. The page table entry is updated: **present** = 0, and the swap location is recorded.
8. The physical frame is freed and returned to the pool of available memory.

CPU Cache.

The **CPU cache** is a small, fast memory located on or near the processor that stores copies of frequently accessed data from main memory (RAM). Its purpose is to reduce memory access latency and improve overall performance by exploiting **temporal** and **spatial locality**.

Cache Levels:

- **L1 Cache (Level 1):** The smallest and fastest cache, located right on the CPU core. 64 KB, but takes $< 1nsec$ to access.
- **L2 Cache (Level 2):** Larger and slightly slower than L1, still located on the CPU core. 256-512 KB, takes $< 4nsec$ to access.
- **L3 Cache (Level 3):** Shared among multiple cores, larger but slower than L1 and L2. Located on the CPU die. 6-32 MB, takes $10s\ of\ nsec$ to access.

Hierarchy Behavior:

Caches form a **hierarchy**: the CPU checks L1 first, then L2, then L3, and finally main memory if needed. This optimizes for latency and hit rate.

Paging.

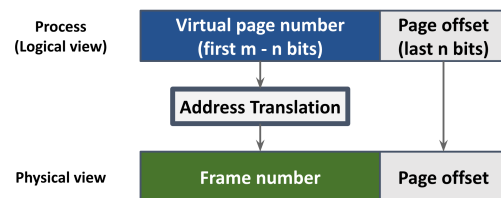
Paging divides virtual memory into fixed-size **pages**, mapped to physical frames via a **page table**. It avoids **ext. fragmentation**, but may cause **internal fragmentation** if a page is only partially used.

Address Representation.

Virtual Address Size is not always equal to Physical Address Size!

Let m be the number of bits in the **virtual address** (i.e., $m = \log_2(\text{virtual address space size})$).

This also corresponds to the number of bits required to uniquely address every byte in the virtual memory. If the virtual memory consists of 2^k pages and each page is 2^n bytes, then the total virtual address space is $2^k \times 2^n = 2^m$ bytes, and thus $m = k + n$.



A virtual address is split into:

- **Offset** (n bits): identifies a byte within a page of size 2^n bytes, where $n = \log_2(\text{page size})$.
- **Virtual Page Number (VPN)** ($m - n$ bits): selects the page entry from the page table.

Page Table.

A process's **page table** maps each **Virtual Page Number (VPN)** to a **Physical Frame Number (PFN)**. The page table is indexed by the VPN and stores **Page Table Entries (PTEs)**.

Linear Page Table.

A single-level table with 2^{m-n} entries. Each VPN directly indexes a PTE. Simple but potentially large: \Rightarrow Table size = $2^{m-n} \times \text{PTE size}$ (may span multiple memory pages).

PTE Content (Typical Metadata Bits):

- **Present bit (P):** Valid address translation exists
- **Protection bits (R/W/X):** Read, write, execute permissions
- **User/Supervisor (U/S):** Access control (user vs. kernel)
- **Dirty bit (D):** Set if the page has been modified
- **Access/Reference bit (A):** Set on access; used in replacement policies

Multi-Level Page Tables.

Used to reduce memory overhead. Split the VPN into multiple parts to form a tree-like hierarchy. The **offset** remains the last $n = \log_2(\text{page size})$ bits of the address.

General Breakdown (Multi-Level Paging):

$$\text{Virtual Address Size} = m = \sum k_i + n$$

where:

- $n = \log_2(\text{page size})$: offset bits
- k_i : number of bits used at each level of the page table

Two-Level Page Table:

$$\text{VPN} = \underbrace{\text{Level 1 index}}_{k_1} \parallel \underbrace{\text{Level 2 index}}_{k_2}, \quad k_1 + k_2 = m - n$$

- Level 1 index (k_1 bits): selects the **page directory entry (PDE)**
- Level 2 index (k_2 bits): selects the **page table entry (PTE)** from the second-level table

Three-Level Page Table:

$$\text{VPN} = \underbrace{\text{L1 index}}_{k_1} \parallel \underbrace{\text{L2 index}}_{k_2} \parallel \underbrace{\text{L3 index}}_{k_3}, \quad k_1 + k_2 + k_3 = m - n$$

- Level 1 index (k_1 bits): selects **first-level page directory**
- Level 2 index (k_2 bits): selects **second-level directory/table**
- Level 3 index (k_3 bits): selects the final PTE

Bit Allocation in Multi-Level Page Tables:

The VPN portion ($m - n$ bits) is divided into L parts ($k_1 + k_2 + \dots + k_L = m - n$), one per level.

- **Bit division is not necessarily even.** Systems may assign more bits to higher levels.
- If uneven, **last levels (closer to the leaf)** typically receive **fewer bits** (fewer entries).
- Each level i contains 2^{k_i} entries, indexing the next level or the final PTE.

Segmentation.

Segmentation divides memory into variable-sized **logical segments** (e.g., code, data, stack), each with a base and limit. It aligns with program structure and supports segment-level protection, but suffers from **ext. fragmentation** due to variable-sized allocations.

Performance Metrics.

- **Latency** — time per I/O ($\mu s-ms$)
- **Throughput** — data/sec (MB/s)
- **IOPS** — ops/sec

Caching.

Keep frequently accessed data in memory to reduce latency. Speeds up reads.

Batching.

Group multiple I/O operations to reduce syscall overhead and disk seeks.

Block Cache.

- In-memory cache of disk blocks
- **Good for reads** — avoids repeated disk access
- **Limited for writes** — consistency requires flushing

Write Caching.

- **Write-through** — data written to both cache and disk immediately
(slower, but consistent)
- **Write-back** — data stays in cache, flushed later
(faster, but risk of data loss on crash)
- `fsync()` — forces flush of dirty blocks to disk