

# **The OS-Compiler-CPU Triad: A Symphony of Computation**

The execution of any software, from the simplest script to the most complex enterprise application, is orchestrated through an intricate interplay between the Central Processing Unit (CPU), the Operating System (OS), and the Compiler. The CPU serves as the fundamental engine of computation, diligently executing instructions, performing arithmetic and logical operations, and manipulating data. The OS acts as the master conductor, managing the computer's hardware resources—including the CPU itself, memory, and input/output (I/O) devices—while providing a consistent and abstracted environment for applications to run. It is the OS's responsibility to ensure system stability, security, and the equitable sharing of resources among multiple competing processes. Bridging the gap between human-readable programming languages and the CPU's native tongue is the compiler. This sophisticated software tool translates high-level code into the precise sequence of machine instructions specific to a target CPU's Instruction Set Architecture (ISA). Beyond mere translation, compilers perform a host of optimizations, refining the generated code to enhance performance, reduce size, and improve energy efficiency. The symbiotic relationship between these three components is foundational to modern computing, enabling the vast and diverse world of software to operate efficiently and securely on an equally diverse range of hardware platforms. This report will delve into the critical mechanisms governing OS-CPU and Compiler-CPU interactions, elucidating concepts such as privilege modes, context switching, program execution mechanics, and code optimization, thereby providing a comprehensive understanding of how software leverages hardware capabilities.

## **I. The Operating System's Interface with the CPU: Orchestrating Hardware Access**

The operating system is pivotal in managing and interacting with the CPU to control program execution, oversee processes, and safeguard system resources. This control is established through fundamental CPU mechanisms that the OS leverages to create a structured and secure computing environment.

### **A. User Mode vs. Kernel Mode: The Foundation of Privileged Control**

Modern computing systems are built upon a foundational principle of privileged execution, primarily manifested through distinct CPU operating modes: user mode and kernel mode. This dual-mode operation is not merely a software convention but is deeply ingrained in the hardware architecture of contemporary CPUs.<sup>1</sup>

## Defining Privilege Levels and Protection Rings

CPUs implement at least two privilege levels, often conceptualized as protection rings, to segregate trusted system code from untrusted user applications.<sup>1</sup> Kernel mode, also known as supervisor mode, system mode, or Ring 0 in the x86 architecture, grants the executing code complete and unrestricted access to all hardware resources and memory.<sup>1</sup> This powerful mode is exclusively reserved for the core components of the operating system, often referred to as the kernel. In contrast, user mode, or Ring 3 in x86, imposes strict limitations. Applications running in user mode cannot directly access hardware or reference arbitrary memory locations; instead, they must request such services from the OS kernel via a well-defined interface, typically system calls.<sup>1</sup>

The rationale behind this separation is paramount for system integrity and security. It prevents user programs, which may be buggy or malicious, from interfering with the OS, accessing or modifying the memory of other programs, or directly manipulating sensitive hardware components.<sup>1</sup> Consequently, a crash or error occurring in a user-mode application is usually contained within that application and is recoverable, often by terminating the offending process. Conversely, a crash in kernel mode is typically catastrophic, leading to a system-wide halt or failure, as the code operating in this mode has the privileges to affect the entire system.<sup>1</sup> This hardware-enforced segregation is a cornerstone of secure and stable multi-user, multitasking operating systems. The inherent security provided by this dual-mode operation is a direct consequence of this hardware-software co-design, where the OS utilizes CPU-provided mechanisms. Without hardware enforcement, software-only protection schemes could be easily circumvented by malicious code that simply bypasses software checks to gain direct hardware access.

## The Mode Bit: Hardware Indication of Current Privilege

The CPU determines its current operational mode through a special hardware flag, often referred to as the 'mode bit', located within a control register such as the processor status register.<sup>3</sup> Conventionally, a mode bit value of 0 signifies kernel mode, while a value of 1 indicates user mode.<sup>3</sup> This bit is automatically checked by the CPU hardware whenever an attempt is made to execute a privileged instruction or access a protected memory region. If a user-mode process attempts such an operation, the hardware will typically generate an exception, transferring control to the OS to handle the violation.

## Transitions Between Modes: System Calls, Interrupts, and Exceptions

Processes transition between user and kernel modes through well-defined mechanisms:

1. **System Calls:** When a user program requires a service that necessitates privileged access—such as opening a file, creating a new process, or sending data over a network—it initiates a system call. This is not a simple function call but involves a special trap or software interrupt instruction (e.g., `syscall` on modern

Linux systems, or the historical `int 0x80` instruction).<sup>3</sup> Executing this instruction causes the CPU to switch from user mode to kernel mode and begin executing code at a predefined entry point within the OS kernel. The kernel then validates the request, performs the service, and subsequently executes another special instruction to return control to the user program, switching the CPU back to user mode.

2. **Interrupts:** External hardware devices, such as disk controllers, network interface cards, or timers, can signal the CPU that they require attention by generating an interrupt.<sup>1</sup> When an interrupt occurs, the CPU typically completes the current instruction, then suspends the currently executing process, saves its essential state, switches to kernel mode, and jumps to a specific interrupt handler routine within the OS. After the interrupt has been serviced, the OS may resume the interrupted process or switch to a different one.
3. **Exceptions (Traps):** Errors or unusual conditions arising during the execution of a program, such as division by zero, an attempt to access an invalid memory address, or execution of an illegal instruction, cause an exception or trap.<sup>1</sup> Similar to interrupts, traps cause the CPU to switch to kernel mode and transfer control to an OS exception handler. The OS then analyzes the nature of the exception and may terminate the process, deliver a signal to it, or attempt some form of recovery.

While these transitions are vital for system operation, they are not without cost. The act of switching modes involves saving the current execution context (registers, program counter), changing the mode bit, potentially flushing CPU pipelines to ensure instructions from the previous mode are not inadvertently executed with incorrect privileges, executing kernel code, and then restoring the context to return to user mode.<sup>13</sup> This process introduces a degree of performance overhead.<sup>8</sup> This overhead underscores a fundamental trade-off in system design: the enhanced security, stability, and control afforded by kernel-mediated access to resources come at the price of the performance cost associated with mode transitions. OS designers continually strive to minimize this overhead, for instance, by optimizing the code paths for frequent system calls or, in some specialized cases, by using mechanisms like vDSO (virtual dynamic shared object) in Linux to handle certain system calls entirely in user space, thus avoiding the mode switch.

Significance: Ensuring System Protection, Resource Management, and Stability

The dual-mode operation is of paramount significance for several reasons:

- **System Protection:** It forms the bedrock of system security by preventing user

programs, whether errant or malicious, from directly accessing or damaging critical OS data structures, other processes' memory, or hardware.<sup>1</sup> Privileged instructions, which could potentially compromise system integrity (e.g., halting the CPU, modifying memory management registers), are restricted to execution in kernel mode only.<sup>5</sup>

- **Resource Management:** By operating in kernel mode, the OS has the ultimate authority to manage and allocate all system resources, including CPU time, physical memory, disk space, and I/O devices.<sup>3</sup> This allows the OS to enforce policies for fair sharing, prevent resource starvation, and optimize overall system throughput.
- **Stability:** The isolation of user processes from each other and from the kernel greatly enhances system stability. If a user application encounters a critical error and crashes, the damage is typically confined to that application alone; the OS and other running applications can continue to function unaffected.<sup>1</sup> This is in stark contrast to older, single-mode operating systems where a faulty application could easily bring down the entire system.

**Table: Comparison of User Mode and Kernel Mode**

To crystallize the distinctions, the following table summarizes the key characteristics of user mode and kernel mode:

Feature	User Mode	Kernel Mode
Privilege Level	Low (Restricted, Unprivileged, Ring 3)	High (Privileged, Master, System, Ring 0)
Hardware Access	Indirect (via System Calls)	Direct and Unrestricted
Memory Access	Restricted to own process's address space	Access to all memory (kernel & user spaces)
Privileged Instructions	Cannot execute	Can execute
Mode Bit Value (Typical)	1	0

Typical Code	Applications, User-level Libraries	OS Kernel, Device Drivers
Impact of Crash	Usually affects only the individual process	Can halt the entire system
Transition Trigger	System Calls, Interrupts, Exceptions (to Kernel)	Return from SC/Interrupt/Exception (to User)
Virtual Address Space	Separate per process	Kernel space often shared across all processes

Sources: <sup>1</sup>

## B. System Calls: User Programs Requesting Kernel Services

System calls are the fundamental mechanism through which user-level programs request services from the operating system kernel.<sup>7</sup> These services encompass a wide range of operations that user programs are not privileged to perform directly, such as file manipulation, process control, network communication, and device management.<sup>7</sup>

### The Mechanism: Software Interrupts/Traps into the Kernel

When a user process needs to invoke an OS service, it executes a special machine instruction, often referred to as a software interrupt or a trap instruction (e.g., the syscall instruction on modern Linux x86-64 systems, or the older int 0x80 interrupt).<sup>3</sup> This instruction causes an immediate transfer of control from the user program to a predefined entry point within the kernel. Simultaneously, the CPU's execution mode is switched from user mode to kernel mode, granting the kernel the necessary privileges to perform the requested operation.<sup>7</sup>

To identify the specific service being requested, the user program typically loads a unique system call number into a designated CPU register (e.g., a7 in RISC-V for Linux, rax in Linux x86-64) before executing the trap instruction.<sup>7</sup> Arguments to the system call are also passed from the user program to the kernel, usually via other CPU registers or, for a larger number of arguments, through a block of memory on the user process's stack.<sup>7</sup> Once the kernel completes the requested service, it places a return value (indicating success or failure, and any result data) into a designated register and executes a special return-from-trap instruction, which switches the CPU mode back

to user mode and transfers control back to the user program, typically to the instruction immediately following the system call invocation.

This mandatory pathway through system calls for all privileged operations creates a controlled bottleneck. This "bottleneck" is a deliberate and essential design choice, as it allows the kernel to act as a vigilant gatekeeper. It validates all requests, enforces security policies and access controls (e.g., file permissions <sup>7</sup>), and serializes access to shared system resources, thereby preventing the chaos that would ensue from direct, unmediated resource access by multiple processes. While this mediation introduces some performance overhead <sup>8</sup>, it is a necessary trade-off for maintaining system order, security, and stability. The design of the system call interface itself—often kept relatively small and rigorously secured, as noted in <sup>7</sup>—is thus critical to the overall robustness of the operating system.

#### Categorization of System Calls

System calls can be broadly categorized based on the type of service they provide <sup>8</sup>:

- **Process Control:** These calls manage the lifecycle of processes, including creation (e.g., fork, CreateProcess), termination (e.g., exit, TerminateProcess), loading and executing new programs (e.g., exec family, CreateProcess), waiting for process completion (e.g., waitpid, WaitForSingleObject), and managing process attributes (e.g., getpid, setpriority). Examples include fork, exec, exit, and waitpid.<sup>7</sup>
- **File Management:** These calls handle operations on files and file systems, such as creating, deleting, opening, closing, reading, writing, and repositioning file pointers (e.g., open, read, write, close, lseek).<sup>9</sup> They also manage file attributes (e.g., stat, chmod).
- **Device Management:** These calls are used to request and release access to hardware devices, read from and write to devices, and manage device attributes.
- **Information Maintenance:** These calls allow programs to retrieve or set system information (e.g., current time, system configuration, process or file metadata).<sup>9</sup>
- **Communication:** These calls facilitate inter-process communication (IPC) and network communication, including creating and deleting communication connections (e.g., pipes, sockets), sending and receiving messages, and transferring status information.<sup>9</sup>
- **Protection:** These calls manage access control and permissions, such as getting and setting file permissions.<sup>9</sup>

System calls can also be classified based on their blocking behavior <sup>16</sup>:

- **Fully Synchronous (Fast) System Calls:** These calls are expected to return very quickly, typically within microseconds. They are generally CPU-bound, do not require the calling process to block or wait for external events, and complete their work without yielding the CPU to other processes. Examples include `getpid()` (get process ID), `getuid()` (get user ID), and `gettimeofday()` (get current time).<sup>16</sup>
- **Asynchronous System Calls:** These calls initiate an operation (often I/O-related) but return to the caller immediately, allowing the calling process to continue its execution while the operation proceeds in the background. The completion of the operation is typically signaled later, perhaps via a callback or an event. Examples include POSIX Asynchronous I/O functions like `aio_read`.<sup>16</sup>
- **Blocking with Timeout System Calls:** These calls may block the calling process while waiting for an event or resource, but they have a predefined maximum duration for waiting. If the event does not occur within the specified timeout period, the system call returns, often with an error indicating the timeout. Examples include `select()` or `poll()` with a timeout value.<sup>16</sup>
- **Indefinitely Blocking (Slow) System Calls:** These calls may block the calling process for an unpredictable or potentially indefinite amount of time, typically while waiting for I/O operations to complete or for specific external events to occur. When a process makes such a call and the required resource or event is not immediately available, the kernel usually suspends the calling process (places it on a wait queue) and schedules another process to run on the CPU.<sup>16</sup> Common examples include `read()` from a terminal, pipe, or socket when no data is available; `write()` to a full pipe or blocked socket; `accept()` waiting for an incoming network connection; and `wait()` waiting for a child process to terminate.<sup>16</sup>

The existence of this spectrum of system call "weights"—from fast, non-blocking calls to slow, indefinitely blocking calls—has significant implications for application design and responsiveness. While fast system calls have minimal impact on process scheduling, slow system calls can cause the calling process to be suspended by the OS, which then schedules other processes to run.<sup>16</sup> Applications that perform many synchronous, slow, blocking system calls without appropriate architectural considerations (like using separate threads for blocking operations or employing asynchronous programming patterns) can become unresponsive, leading to a poor user experience or low throughput in server applications. This reality drives the adoption of asynchronous programming models, event-driven architectures, and multithreading to maintain responsiveness and efficiency.

The System Call Interface and Wrapper Libraries (e.g., `libc`)



Operating systems expose their system call functionalities to application programs through an Application Programming Interface (API). Direct invocation of system calls using trap instructions and register manipulation is complex, architecture-dependent, and error-prone. Therefore, systems typically provide a standard library, such as the C library (e.g., glibc on Unix-like systems, msvcrt.dll providing wrappers around the Native API in ntdll.dll on Windows), that sits between user applications and the OS kernel.<sup>9</sup>

This library provides wrapper functions for most system calls. These wrappers often have names identical or similar to the underlying system calls they invoke (e.g., the C function `open()` is a wrapper for the `open` system call).<sup>7</sup> The wrapper function takes care of the low-level details: marshalling arguments into the correct registers or onto the stack according to the specific ABI (Application Binary Interface), loading the appropriate system call number into the designated register, executing the trap instruction to enter the kernel, and then checking the return value from the kernel to set `errno` in case of an error and return an appropriate value to the application. This abstraction provided by library wrappers significantly enhances the portability of application code across different operating systems (that adhere to a standard like POSIX) and simplifies the task of programming for application developers.

### **C. Interrupts and Signals: Handling Asynchronous Events and Inter-Process Communication**

Beyond the deliberate requests for kernel services via system calls, the CPU's execution flow can also be altered by asynchronous events, managed through hardware interrupts and software signals. These mechanisms are crucial for the OS to respond to external stimuli and for processes to communicate or react to exceptional conditions.

#### **Hardware Interrupts: Device Attention and CPU Response**

A hardware interrupt is an asynchronous signal sent from a hardware device (e.g., disk controller, network card, timer chip, keyboard) to the CPU, indicating that the device requires immediate attention.<sup>11</sup> For instance, a disk controller might interrupt the CPU when a requested data block has been read from the disk and is available in memory, or a network card might interrupt when a new packet has arrived.

If the CPU is not currently executing a task that has a higher priority than the interrupt itself, it will suspend its current activity upon receiving an interrupt signal.<sup>11</sup> The CPU saves the essential state of the currently running process (at a minimum, the program counter and status register), switches to kernel mode if it wasn't already in it, and then transfers control to a specific piece of OS code known as an Interrupt Service Routine (ISR) or interrupt handler, which is designated to handle that particular type of



interrupt.<sup>11</sup>

Interrupts are fundamental for several reasons <sup>7</sup>:

- **Efficient I/O Management:** They allow the CPU to perform other tasks while I/O devices operate concurrently. Instead of the CPU repeatedly polling a device to check its status (busy-waiting), the device can notify the CPU via an interrupt only when it needs service. This is significantly more efficient.<sup>7</sup>
- **Real-time Responsiveness:** Interrupts enable the system to react promptly to real-world events, such as mouse movements or keyboard input, making interactive systems feel responsive.<sup>7</sup>
- **Multitasking and Scheduling:** Timer interrupts are a cornerstone of preemptive multitasking. The OS programs a hardware timer to generate interrupts at regular intervals. When a timer interrupt occurs, the OS regains control from the currently running process and can make scheduling decisions, such as whether to continue running the current process or to switch to a different one (context switch).<sup>7</sup> This ensures fair sharing of CPU time among multiple processes.

The very ability of an operating system to multitask and react to its environment in a timely fashion is built upon the interrupt mechanism. Without interrupts, the OS would either be forced into inefficient polling loops for device status or would effectively lose control once a user process began execution, thereby preventing preemption and true multitasking. The latency involved in servicing an interrupt—the time from when an interrupt is signaled to when the ISR begins execution and completes its critical tasks—is therefore a key performance metric for any OS.

#### The Interrupt Handling Process

The OS maintains a table (often called an interrupt vector table) that maps each type of interrupt to the memory address of its corresponding ISR. When an interrupt occurs, the CPU (or an interrupt controller) uses the interrupt's identifier to look up the address of the appropriate ISR in this table and then transfers execution to it.<sup>12</sup>

An ISR typically performs the following tasks <sup>12</sup>:

1. Saves any additional CPU state (registers) that it might modify and which were not automatically saved by the hardware.
2. Determines the cause of the interrupt, especially if multiple devices can trigger the same interrupt line (interrupt sharing). For polled interrupts, the ISR may need to query several devices to find the one that signaled. Vectored interrupts directly provide information identifying the source.
3. Services the device: This might involve reading status information, transferring

data between the device and memory (often with the help of Direct Memory Access - DMA), or acknowledging the interrupt to the device to stop it from further signaling.

4. Performs any necessary I/O request-related processing, such as waking up a process that was waiting for this I/O to complete.
5. Restores the saved CPU state.
6. Executes a special return-from-interrupt instruction, which causes the CPU to switch back to the mode and state of the interrupted process, allowing it to resume execution as if nothing had happened (unless a context switch is decided upon).

Modern systems support various types of interrupts, including legacy (fixed-pin) interrupts, Message-Signaled Interrupts (MSI), and Extended Message-Signaled Interrupts (MSI-X), which offer more flexibility and scalability, especially for PCI Express devices.<sup>12</sup>

Signals: Kernel-to-Process and Process-to-Process Communication

Signals are a software-based inter-process communication (IPC) mechanism used in Unix-like operating systems and other POSIX-compliant systems. They serve as a way for the kernel to notify a user-space process of an asynchronous event or to instruct it to take a specific action, and also for processes to send notifications to each other.<sup>7</sup>

Examples of events that can trigger signals include:

- Hardware exceptions detected by the CPU (e.g., illegal instruction, division by zero, segmentation fault), which the kernel translates into signals (e.g., SIGILL, SIGFPE, SIGSEGV).
- User actions, such as pressing Ctrl-C in a terminal, which typically sends a SIGINT (interrupt) signal to the foreground process.
- Software events, such as a timer expiring (SIGALRM), or one process explicitly sending a signal to another using the kill() system call.<sup>7</sup> A process can send a signal to itself using raise().<sup>7</sup>

Each process can define how it wishes to handle specific signals. For a given signal, a process can:

1. **Perform the default action:** Each signal has a default disposition (e.g., terminate the process, ignore the signal, stop the process, dump core).
2. **Ignore the signal:** The process requests that the signal be discarded. (Some signals, like SIGKILL and SIGSTOP, cannot be ignored or caught).
3. **Catch the signal:** The process registers a custom function, called a signal

handler. When the signal is delivered to the process, its normal flow of execution is temporarily interrupted, and the registered signal handler is executed.<sup>7</sup> After the handler finishes, the process typically resumes execution from where it was interrupted, unless the handler itself modified the execution flow (e.g., by calling `exit()` or `longjmp()`).

A particularly complex aspect of signals is their interaction with system calls, especially blocking (slow) system calls. If a signal is delivered to a process while it is blocked in a system call (e.g., waiting for input in a `read()` call), the system call may be interrupted and return prematurely, typically with an error code of `EINTR` (Interrupted system call).<sup>16</sup> The application program is then responsible for checking for `EINTR` and, if appropriate, restarting the interrupted system call. Some systems offer an option (e.g., the `SA_RESTART` flag when registering a signal handler with `sigaction()`) to have certain system calls automatically restarted by the kernel after the signal handler returns.<sup>16</sup>

Signals provide a user-space mechanism for handling asynchronous events that is somewhat analogous to how hardware interrupts are handled at the kernel level. Both disrupt the normal flow of execution to invoke a dedicated handler routine.<sup>7</sup> However, the `EINTR` behavior highlights a significant complexity in signal handling. Robust application code must be prepared to deal with system calls that are prematurely terminated by signals, often by implementing retry logic. While `SA_RESTART` aims to simplify this, it is not universally applicable to all system calls or all situations, meaning the underlying challenge of managing these interactions persists. This inherent complexity in traditional signal handling, especially for I/O-bound applications, has been one of the motivating factors for the development of alternative asynchronous event notification mechanisms in modern operating systems, such as `epoll` in Linux, `kqueue` in BSD systems, and I/O completion ports in Windows, which can offer more predictable and manageable semantics for handling multiple I/O events concurrently.

#### **D. Context Switching: Enabling Concurrency and Multitasking**

Context switching is a fundamental operation performed by the operating system's scheduler to switch the CPU's focus from one executing process (or thread) to another.<sup>17</sup> This mechanism is the cornerstone of multitasking operating systems, allowing multiple processes to share a single CPU core, thereby creating the illusion that they are running simultaneously.<sup>15</sup>

Definition, Necessity, and Triggers

The primary necessity for context switching lies in its ability to facilitate CPU sharing among various tasks, manage the execution of high-priority processes, handle situations where processes must wait for I/O operations to complete, and respond effectively to hardware interrupts.<sup>17</sup>

Context switching can be triggered by several events:

- **Multitasking/Time Slicing:** In preemptive multitasking systems, each process is allocated a small unit of CPU time called a time slice or quantum. When a process's time slice expires, the OS scheduler may decide to switch to another ready process.<sup>15</sup>
- **I/O Wait / Resource Wait:** If a running process initiates an operation that requires waiting for a resource (e.g., reading data from a disk, waiting for network input, or acquiring a locked mutex), the OS will typically suspend that process and switch to another process that is ready to run. This prevents the CPU from idling while the first process waits.<sup>17</sup>
- **Interrupts:** The arrival of a hardware interrupt (e.g., from an I/O device or a timer) causes the current process to be suspended and an interrupt handler to execute. After the interrupt is handled, the OS scheduler may decide that a different process (perhaps one that was waiting for the event signaled by the interrupt, or a higher-priority process) should run next.<sup>17</sup>
- **System Calls:** Certain system calls can explicitly cause the calling process to yield the CPU or block. For example, a call to `sleep()` will suspend the process for a specified duration, or a `wait()` call will block the process until a child terminates.<sup>19</sup>
- **Higher Priority Task Arrival (Preemption):** If a process with a higher priority than the currently running process becomes ready to execute (e.g., it was unblocked or its status changed), a preemptive scheduler will interrupt the current process and switch to the higher-priority one.<sup>17</sup>

It is important to distinguish a context switch between processes from a mode switch (user to kernel mode or vice-versa) within the *same* process. While events like system calls or interrupts trigger a mode switch, they only lead to a context switch if the kernel subsequently decides to schedule a *different* process to run.

#### The Process: Saving and Restoring CPU State

The mechanics of a context switch involve two primary actions: saving the state of the currently running ("old") process and restoring the state of the next process to be run (the "new" process).<sup>17</sup> This entire operation is orchestrated by the OS.

The detailed steps are typically as follows <sup>19</sup>:

1. **Save Old Process State:** The OS saves the complete execution context of the currently running process. This context is stored in a per-process data structure maintained by the kernel, known as the Process Control Block (PCB).<sup>17</sup> The saved information includes:
  - **CPU Registers:** The contents of all general-purpose registers, the program counter (PC) (which holds the address of the next instruction to be executed), the stack pointer (SP), and any status registers.
  - **Process State:** Information about the current state of the process (e.g., running, ready, waiting).
  - **Memory Management Information:** Pointers to page tables or segment tables that define the process's address space.
  - **Scheduling Information:** Process priority, time slice remaining, etc.
  - **Accounting Information:** CPU time used, resource consumption.
  - **I/O Status Information:** List of open files, I/O devices allocated to the process.
2. **Update Old Process's PCB:** The PCB of the outgoing process is updated with its new state (e.g., ready, waiting for I/O) and placed into the appropriate queue (e.g., ready queue, I/O wait queue).
3. **Select New Process:** The OS scheduler selects the next process to run from the ready queue, based on its scheduling algorithm (e.g., priority-based, round-robin).
4. **Load New Process State:** The OS loads the saved context of the selected "new" process from its PCB into the CPU's registers. This includes restoring the program counter, stack pointer, general-purpose registers, and updating memory management unit (MMU) registers with the new process's address space information.
5. **Resume New Process:** The new process then resumes execution from the point where it was last stopped, as if it had never been interrupted.<sup>17</sup>

#### Role of the Process Control Block (PCB)

The Process Control Block (PCB), sometimes called a Task Control Block (TCB), is a critical kernel data structure that stores all the information the OS needs about a particular process.<sup>17</sup> It serves as the repository for the process's context when it is not running. During a context switch, the PCB of the outgoing process is where its state is saved, and the PCB of the incoming process is the source from which its state is restored. Without the PCB, the OS would have no way to manage multiple processes or to switch between them seamlessly.

#### Performance Implications and Overhead

Context switching, while essential for concurrency, is not a free operation. It represents pure

overhead from the perspective of user applications, as the CPU is dedicated to the administrative task of switching contexts rather than executing user-level code.<sup>17</sup> The time taken to perform a context switch can be significant, typically in the range of a few microseconds to milliseconds, depending on the OS, the CPU architecture, and the amount of state that needs to be saved and restored.<sup>18</sup> This "cost of concurrency" directly impacts system throughput and responsiveness. If context switches occur too frequently (e.g., due to very short time slices), a substantial portion of CPU time can be consumed by the switching overhead itself, reducing the time available for actual computation. This leads to a critical trade-off for OS designers: shorter time slices can improve system responsiveness to interactive events, but they also increase the frequency of context switches and thus the overhead. Conversely, longer time slices reduce overhead but can make the system feel less responsive. This tension often motivates OS designs and application architectures that aim to minimize unnecessary context switches, such as employing asynchronous I/O or carefully managing thread pools.

Context switches between threads within the same process are generally faster than context switches between different processes.<sup>18</sup> This is because threads of the same process share the same memory address space and many other resources (like open file descriptors).<sup>15</sup> Therefore, a thread context switch primarily involves saving and restoring thread-specific data, such as CPU registers, the stack pointer, and thread-local storage, without the need to change memory management structures like page tables. This differentiated cost strongly encourages the use of multithreading for concurrent tasks that operate on shared data within a single application, as it offers a more lightweight mechanism for concurrency compared to using multiple, separate processes. This performance advantage is a major reason for the widespread adoption of multithreaded programming models in applications requiring high levels of concurrency, such as web servers, database management systems, and graphical user interfaces.

While some CPU architectures have included hardware support for context switching (e.g., the Task State Segment - TSS in x86 architecture<sup>18</sup>), the dominant method in general-purpose operating systems is software-based context switching, managed by the OS scheduler.<sup>18</sup> Hardware mechanisms, while seemingly efficient, can be rigid and may save more state than necessary (e.g., almost all registers, as noted for TSS in<sup>18</sup>), potentially making them slower than a more selective software approach.

Software-based switching provides greater flexibility, allowing the OS to integrate context switching with complex scheduling algorithms and tailor the process to specific system needs and policies. This preference for software flexibility over fixed hardware solutions is a recurring theme in system design, highlighting the adaptability required to manage evolving hardware and software landscapes.



## II. Program Execution Mechanics: The Runtime Choreography of Software

The execution of a program, particularly the invocation of functions and the management of their data, involves a precise and standardized set of low-level operations. These mechanics are centered around the call stack and agreed-upon calling conventions, which together ensure orderly execution and data integrity.

### A. Stack Frames (Activation Records): Managing Function Invocations

At the heart of function call management is the concept of the stack frame, also known as an activation record. This structure is fundamental to how procedural and functional programming languages implement subroutines.

Purpose: Storing State for Active Subroutines

Each time a function is called during program execution, a new block of memory is allocated on the call stack. This block is the stack frame for that particular invocation of the function.<sup>19</sup> The primary purpose of the stack frame is to hold all the necessary information associated with that function call, including its arguments, local variables, and the information needed to return control to the caller. The call stack itself operates in a Last-In-First-Out (LIFO) manner: the stack frame for the most recently called function resides at the top of the stack, and when a function returns, its frame is removed (popped) from the stack.<sup>19</sup> This LIFO behavior elegantly handles nested function calls and recursion; each new call simply pushes a new frame, and each return pops the current frame, automatically restoring the context of the calling function without complex bookkeeping.<sup>23</sup>

Typical Contents: Function Arguments, Local Variables, Return Address, Saved Frame Pointer  
A typical stack frame contains several key pieces of information <sup>23</sup>:

- **Function Arguments (Parameters):** The values passed to the function by its caller. Depending on the calling convention, these might be placed on the stack by the caller before the call, or some might be passed in registers and then potentially saved into the stack frame by the callee.
- **Local Variables:** Space allocated for variables that are declared within the function and are local to its scope. This ensures that each invocation of a function (especially in recursion) has its own private set of local variables.
- **Return Address:** The memory address of the instruction in the calling function's code to which execution should return after the current function completes.<sup>19</sup> This is typically pushed onto the stack by the CALL instruction itself.
- **Saved Frame Pointer (Old Base Pointer):** The value of the frame pointer register from the *calling* function's stack frame. This is saved so that when the current function returns, the caller's frame pointer (and thus its stack frame context) can



be restored.

- **Saved Registers:** Some calling conventions require the called function (callee) to preserve the values of certain CPU registers for the caller. If the callee uses these registers, it must save their original values (often in its stack frame) before modifying them and restore them before returning.

#### Stack Pointer (SP) and Frame Pointer (FP/BP) Registers

Two special-purpose CPU registers are crucial for managing the call stack and its frames:

- **Stack Pointer (SP):** This register always points to the current "top" of the stack—the memory address of the last item pushed onto it or the next free location, depending on the architecture and convention.<sup>23</sup> The SP is automatically decremented when data is pushed onto the stack (assuming the stack grows downwards in memory, which is common) and incremented when data is popped off.
- **Frame Pointer (FP) or Base Pointer (BP):** This register, when used, points to a fixed location within the current active stack frame.<sup>23</sup> Typically, it might point to the location where the old frame pointer was saved or to the beginning of the local variable area. The FP provides a stable base address from which function arguments and local variables can be accessed using fixed offsets. This is useful because the SP might change during the function's execution as temporary values are pushed and popped.

However, the use of a dedicated frame pointer represents a trade-off. While it simplifies debugging (as GDB and other debuggers can easily traverse the stack by following the chain of saved FPs<sup>26</sup>) and makes accessing stack variables via fixed offsets straightforward, it consumes a general-purpose register. In architectures with a limited number of registers, or for highly optimized code, compilers might choose to omit the frame pointer (e.g., through the `-fomit-frame-pointer` option in GCC<sup>26</sup>). In such cases, local variables and arguments are accessed via offsets relative to the SP. This can make debugging more complex, requiring the debugger to analyze the code to determine frame layouts, but it frees up a register for general computation, potentially improving performance. Thus, debug builds often retain frame pointers, while release builds might omit them.

#### Creation and Destruction of Stack Frames During Call and Return

The lifecycle of a stack frame is tightly coupled with the function call and return sequence:

- **Function Call (Simplified):**
  1. The *caller* prepares for the call. This may involve pushing arguments onto the

stack (if the calling convention dictates stack-based argument passing for some or all arguments).

2. The caller executes a CALL instruction. This instruction typically pushes the *return address* (the address of the instruction immediately following the CALL) onto the stack. Control is then transferred to the first instruction of the *callee* (the function being called).
3. The *callee* then establishes its own stack frame (prologue):
  - It saves the caller's frame pointer by pushing the current value of FP/BP onto the stack.
  - It sets its own frame pointer: FP/BP is set to the current value of SP. This new FP/BP now points to a fixed location in the callee's frame.
  - It allocates space for its local variables by decrementing the SP by the required amount. (Other registers that need to be preserved might also be saved on the stack here).

- **Function Return (Simplified):**

1. The *callee* prepares to return (epilogue):
  - It places its return value in a designated register (e.g., EAX/RAX for integers) or on the stack if it's a large value, according to the calling convention.
  - It deallocates its local variables, typically by setting SP to the value of FP/BP. This effectively discards the local variable space.
  - It restores the caller's frame pointer by popping the saved FP/BP value from the stack back into the FP/BP register.
  - (Other saved registers are restored here).
2. The callee executes a RETURN instruction. This instruction typically pops the *return address* from the top of the stack into the program counter (PC).
3. Execution resumes in the *caller* at the instruction following the original CALL.
4. The *caller* may then need to clean up any arguments it pushed onto the stack (if the calling convention specifies caller cleanup).

This meticulous process ensures that each function invocation operates within its own isolated environment on the stack, with its own arguments and local variables, and that control can be correctly passed back and forth.<sup>23</sup>

While this stack-based mechanism is efficient and fundamental, the contiguous allocation of stack frames and the storage of control information like return addresses within them also present a potential security risk. If a function writes data beyond the allocated buffer for a local variable (a buffer overflow), it can overwrite adjacent data

in the stack frame, including the saved frame pointer and, critically, the return address.<sup>23</sup> An attacker could craft an input that causes such an overflow and deliberately overwrite the return address to point to malicious code injected elsewhere in memory or to existing library functions in a way that subverts the program's intended behavior (a technique known as stack smashing). This vulnerability has spurred the development of numerous defense mechanisms, such as stack canaries (random values placed on the stack that are checked before a function returns), non-executable stack memory (NX bit or Data Execution Prevention - DEP), and Address Space Layout Randomization (ASLR), all aimed at making such exploits more difficult.

## B. Calling Conventions: Standardizing Inter-Function Communication

For functions to interact correctly, especially when they might be written by different programmers, compiled by different compilers, or even written in different programming languages, there must be a set of agreed-upon rules governing how this interaction occurs. These rules are known as calling conventions.<sup>27</sup>

Definition and Importance: Ensuring Interoperability and Predictability

A calling convention is a low-level, implementation-defined scheme that specifies precisely how subroutines (functions or methods) receive parameters from their caller and how they return a result. It dictates the mechanics of the function call interface at the machine code level.

The importance of calling conventions cannot be overstated<sup>27</sup>:

- **Consistency:** They provide a uniform and predictable way for functions to interact, regardless of where they originated or which compiler produced them.
- **Interoperability:** This is perhaps the most critical aspect. Calling conventions enable code compiled by different compilers (e.g., GCC and Clang) or written in different programming languages (e.g., C and Assembly, or C and Fortran) to be linked together and call each other successfully. This is fundamental for using pre-compiled libraries.
- **Optimization:** Well-defined calling conventions can aid compilers in generating more efficient code by minimizing the overhead associated with function calls (e.g., by favoring register usage over stack memory for arguments).
- **Maintainability:** A clear understanding of the calling convention being used makes it easier for developers to debug low-level code, understand assembly output, and write code that interfaces with other modules correctly.

Without such "social contracts" at the machine level, modular software development

as we know it—relying on libraries and separate compilation—would be practically impossible. If a calling function (caller) and a called function (callee) have different assumptions about where arguments are placed, how the stack should be managed, or where the return value will be, the program will likely crash, produce incorrect results, or exhibit other erratic behavior. This is why Application Binary Interfaces (ABIs), which standardize calling conventions among other things for a specific platform, are so crucial. Changes to established ABIs are rare and highly disruptive due to the widespread impact on software compatibility.

Key Elements: Argument Passing (Registers vs. Stack), Return Value Handling, Stack Cleanup (Caller vs. Callee)

Calling conventions typically define several key aspects of the function call interface:

- **Argument Passing:** This specifies how arguments are transmitted from the caller to the callee. Options include:
  - **Stack-based:** Arguments are pushed onto the call stack by the caller. The order (e.g., right-to-left or left-to-right) is also defined.
  - **Register-based:** Some or all arguments are placed in specific CPU registers. This is generally faster as register access is quicker than memory access.
  - **Mixed:** A combination, where some arguments are passed in registers and others on the stack (common in modern conventions).
- **Return Value Handling:** This defines how the function's result is returned to the caller. Small return values (like integers or pointers) are often returned in a specific register (e.g., EAX/RAX on x86/x64). Larger return values (like large structures) might be returned via a memory location pointed to by an implicit argument, or on the stack.
- **Stack Cleanup:** This determines who is responsible for removing the arguments from the stack after the function call returns:
  - **Caller Cleanup:** The calling function adjusts the stack pointer to deallocate the space used by the arguments.
  - **Callee Cleanup:** The called function adjusts the stack pointer before returning. The choice between caller and callee cleanup involves a trade-off. Caller cleanup allows for functions with a variable number of arguments (variadic functions, like printf) because only the caller knows how many arguments were actually passed. However, it means the cleanup code is duplicated at every call site, potentially increasing code size. Callee cleanup centralizes the cleanup code within the callee itself, which can lead to smaller overall code size, but it generally requires the function to have a fixed number of arguments.

- **Register Preservation:** The convention specifies which CPU registers are "caller-saved" (volatile) and which are "callee-saved" (non-volatile). A callee can freely modify caller-saved registers. If a callee needs to use a callee-saved register, it must first save its original value (e.g., on the stack) and restore it before returning to the caller.
- **Name Mangling/Decoration (Less directly part of the call, but related to linkage):** Compilers often modify function names in the object code (e.g., by adding prefixes or suffixes encoding argument types).<sup>29</sup> This can be influenced by the calling convention and is essential for type-safe linkage and supporting features like function overloading in C++.

Common Conventions (primarily for x86-32, but concepts apply broadly)

Several calling conventions became common, especially in the x86-32 bit world:

- **cdecl (C declaration):**<sup>27</sup>
  - **Argument Passing:** Arguments are pushed onto the stack in right-to-left order.
  - **Stack Cleanup:** The **caller** is responsible for cleaning the stack after the call returns. This makes cdecl suitable for variadic functions (e.g., printf, scanf) because the caller knows how many arguments it pushed.
  - **Return Value:** Typically returned in the EAX register (or EDX:EAX for 64-bit values).
  - **Name Decoration:** Often, an underscore (\_) is prefixed to the function name in the object code (e.g., \_myfunction). This is the default for many C and C++ compilers for x86-32.
- **stdcall (Standard call):**<sup>27</sup>
  - **Argument Passing:** Arguments are pushed onto the stack in right-to-left order, similar to cdecl.
  - **Stack Cleanup:** The **callee** is responsible for cleaning the stack before returning. This generally leads to slightly smaller code size because the cleanup code is in one place (the callee) rather than duplicated at each call site. However, it means stdcall functions cannot be variadic (unless special measures are taken).
  - **Return Value:** Typically returned in the EAX register.
  - **Name Decoration:** An underscore (\_) is prefixed, and a suffix @ followed by the total number of bytes of arguments on the stack is appended (e.g., \_myfunction@8 for a function taking two 4-byte integers). stdcall was famously used by the Win32 API.
- **fastcall:**<sup>27</sup>

- **Argument Passing:** Aims to improve performance by passing the first few (typically two) integer or pointer arguments in specific CPU registers (e.g., ECX and EDX on Microsoft compilers for x86-32). Any remaining arguments are pushed onto the stack, usually right-to-left. The exact registers used can vary between compilers.
- **Stack Cleanup:** The **callee** is responsible for cleaning the stack (for any arguments passed on the stack).
- **Return Value:** Typically returned in the EAX register.
- **Name Decoration:** Often, an @ symbol is prefixed, and a suffix @ followed by the argument byte count (for stack arguments) is appended (e.g., @myfunction@8).
- **thiscall (C++):**<sup>27</sup>
  - This convention is specifically used by C++ compilers for non-static member functions.
  - **Argument Passing:** The this pointer (which points to the object instance on which the member function is called) is passed in a register (commonly ECX on Microsoft compilers for x86-32). Other explicit arguments might follow stdcall or cdecl rules for stack passing (e.g., right-to-left on the stack).
  - **Stack Cleanup:** Usually, the **callee** cleans the stack for its explicit arguments.
  - This convention is typically implicit and managed by the C++ compiler.

The evolution of these conventions, particularly the introduction of fastcall and the design of modern 64-bit ABIs, reflects a continuous effort to enhance performance by reducing memory accesses. Passing arguments in registers is significantly faster than pushing them onto and popping them off the stack in memory.

**Table: Overview of Common Calling Conventions (x86-32 focus for cdecl, stdcall, fastcall)**

Conventi on	Argument Passing (Order)	Argument Location	Stack Cleanup	Variadic Functions	Typical Use Case	Name Decoratio n Example (MSVC-lik e)
__cdecl	Right-to-L	Stack	Caller	Yes	C/C++ default,	_my_functi

	eft				libraries	on
__stdcall	Right-to-Left	Stack	Callee	No	Win32 API	_my_function@8
__fastcall	Right-to-Left (for stack)	First 2 in ECX/EDX, rest on stack	Callee	No	Performance-critical	@my_function@8
__thiscall	Right-to-Left (for stack)	this in ECX, rest on stack (usually)	Callee	No (usually)	C++ non-static member functions	(Compiler-specific, often mangled)

Sources: <sup>27</sup>

### The x64 ABI (Application Binary Interface): A Detailed Examination

For 64-bit architectures, such as x86-64 (also known as AMD64), standardized ABIs define the calling conventions. The two most prominent are the System V AMD64 ABI (used by Linux, macOS, and other Unix-like systems) and the Microsoft x64 calling convention (used by Windows). While they share many similarities due to the underlying architecture, they have some differences. The following details primarily draw from the Microsoft x64 convention as described in 30, but the general principles of register-based argument passing are common.

- **Register Usage for Arguments:**

- The x64 convention is a "fast-call" type, heavily utilizing registers for argument passing.
- **Integer/Pointer Arguments:** The first four integer or pointer arguments are passed in registers RCX, RDX, R8, and R9, respectively (for Microsoft x64; System V uses RDI, RSI, RDX, RCX, R8, R9 for the first six).
- **Floating-Point Arguments:** The first four floating-point arguments are passed in XMM0, XMM1, XMM2, and XMM3 registers.
- **Mixed Arguments:** If arguments are a mix of integer/pointer and floating-point types, they still consume their respective register slots in order. For example, if the first argument is an integer (RCX) and the second is a float (XMM1), the third integer argument would go into RDX.
- **Larger Arguments:** Arguments that are larger than 8 bytes (e.g., large



structs) or are not simple types (1, 2, 4, or 8 bytes) are typically passed by reference (a pointer to the data is passed in a register).<sup>30</sup> A single argument is never spread across multiple registers.

- **Stack Passing:** Any arguments beyond those that fit in the designated registers are passed on the stack. They are pushed right-to-left.
- **Shadow Store (Spill Space) (Microsoft x64):**
  - The caller is responsible for allocating a 32-byte area on the stack immediately above (at a higher address than) the return address. This area is known as the "shadow store" or "home space" for the first four register arguments (RCX, RDX, R8, R9 or XMM0-XMM3).<sup>30</sup>
  - The callee can use this space to save the values of these argument registers if it needs to reuse those registers for other purposes (e.g., before calling another function). The caller must always allocate this space, even if the callee takes fewer than four parameters. This simplifies support for unprototyped C functions and variadic functions.
- **Stack Alignment:**
  - The stack pointer (RSP) must be 16-byte aligned *before* a CALL instruction is executed.<sup>30</sup> The CALL instruction itself pushes an 8-byte return address onto the stack, which means that upon entry to the callee, RSP is (original RSP - 8), and thus aligned to an 8-byte boundary but not necessarily a 16-byte one.
  - The callee is responsible for re-establishing 16-byte alignment if it needs to (e.g., if it calls other functions, allocates significant stack space, or uses XMM operations that require 16-byte aligned memory operands).
- **Return Values:**
  - Integer-like values (up to 64 bits) are returned in the RAX register.
  - Floating-point values (float, double) are returned in the XMM0 register.
  - Larger structures or user-defined types are typically returned by having the caller allocate space for the return value and pass a pointer to this space as an implicit first argument (often in RCX for Microsoft x64, or RDI for System V).
- **Caller/Callee Saved Registers:**
  - **Volatile (Caller-Saved):** Registers like RAX, RCX, RDX, R8, R9, R10, R11, and XMM0-XMM5 (Microsoft x64) are considered volatile.<sup>30</sup> The callee can use these registers without saving their original values. If the caller needs the values in these registers preserved across a function call, the caller must save them before the call.
  - **Non-Volatile (Callee-Saved):** Other registers (e.g., RBX, RBP, RDI, RSI, RSP, R12-R15, and XMM6-XMM15 (lower 64 bits) for Microsoft x64) are non-volatile.

If a callee uses these registers, it must save their original values upon entry and restore them before returning.

- **Unwindability for Exception Handling:**

- To support structured exception handling, non-leaf functions (functions that call other functions or allocate their own stack frame and modify non-volatile registers) must provide metadata that describes how to "unwind" their stack frame.<sup>30</sup> This metadata (pdata and xdata on Windows) allows the system to walk the stack backwards during exception propagation, correctly restoring register values and finding exception handlers. This imposes restrictions on the structure of function prologs and epilogs to ensure they are describable by this metadata. Leaf functions that don't modify non-volatile registers or RSP significantly may not require this complex unwind information.

The increased number of general-purpose registers in x64 architectures compared to x86-32 has been a key driver in the design of these register-heavy calling conventions, significantly reducing the need for stack-based argument passing for common cases and thus improving performance.

**Table: x64 ABI Register Usage for Parameter Passing (Microsoft Convention Example)**

Parameter Type	1st Arg	2nd Arg	3rd Arg	4th Arg	5th+ Args
Integer/Pointer	RCX	RDX	R8	R9	Stack
Floating Point (float, double)	XMM0	XMM1	XMM2	XMM3	Stack
__m128, structs/unions > 8 bytes (passed by value), arrays,	Pointer in RCX/RDX/R8/R9 to caller-allocated memory	Pointer in RDX/R8/R9/Stack to caller-allocated memory	Pointer in R8/R9/Stack to caller-allocated memory	Pointer in R9/Stack to caller-allocated memory	Stack

strings					
Structs/Unions (8, 16, 32, or 64 bits), __m64 (passed as if integers)	RCX (as integer)	RDX (as integer)	R8 (as integer)	R9 (as integer)	Stack

Note: Arguments passed by reference (pointer) still consume a register slot for the pointer itself. The "Shadow Store" applies to the first four register slots regardless of type.

Sources: 30

### III. The Compiler's Crucial Role: From Source Code to Optimized CPU Instructions

The compiler is a sophisticated software tool that acts as the critical translator between human-written source code in high-level programming languages (like C++, Java, or Python) and the low-level machine code that a CPU can directly execute.<sup>31</sup> Its role extends far beyond simple translation; it also involves extensive analysis and optimization to ensure the generated code is efficient, correct, and tailored to the specific capabilities of the target hardware.

#### A. Generating Machine Code for Specific Instruction Set Architectures (ISAs)

The diversity of CPU designs necessitates that compilers produce code compatible with the particular architecture on which it is intended to run. This compatibility is defined by the Instruction Set Architecture (ISA).

Understanding ISAs: The Software-Hardware Contract

An Instruction Set Architecture (ISA) is an abstract model that formally defines the interface between software and the CPU hardware.<sup>34</sup> It specifies the set of instructions the CPU can execute, the supported data types (e.g., integers, floating-point numbers of various sizes), the available CPU registers, how the hardware manages main memory (including addressing modes and features like virtual memory), and the input/output model.<sup>34</sup> Essentially, the ISA is the "programmer's manual" for the CPU at the lowest level, visible to assembly language programmers and compiler writers.<sup>36</sup>

A key importance of the ISA is that it ensures **binary compatibility**.<sup>34</sup> This means that machine code generated for a specific ISA (e.g., x86-64, ARMv8, RISC-V) can run on any CPU implementation that conforms to that ISA, regardless of the underlying

microarchitectural differences (like pipeline depth, cache sizes, or manufacturing process) between different CPU models from different vendors.<sup>34</sup> This allows, for example, a lower-performance, lower-cost machine to be replaced with a higher-cost, higher-performance machine running the same ISA, and the existing software will run without modification. This abstraction layer is fundamental to the computing ecosystem, allowing hardware and software to evolve somewhat independently. Hardware manufacturers can innovate on microarchitecture to improve performance, power efficiency, or cost, while software developers can write programs that target a stable ISA, confident that their applications will run on a wide range of current and future hardware. This "narrow waist" model, where a relatively stable ISA connects a vast array of software to diverse hardware implementations, has been a powerful enabler of innovation in the computing industry.

#### The Compiler Backend: Translating Intermediate Representation to Target-Specific Assembly/Machine Code

Compilers are typically designed with a modular, multi-phase architecture, commonly divided into three main parts: the front end, the middle end (optimizer), and the back end (code generator).<sup>33</sup>

- **Front End:** This part is responsible for processing the source code of a specific programming language. It performs:
  - **Lexical Analysis:** Breaking the source code into a stream of tokens (keywords, identifiers, operators, literals).<sup>38</sup>
  - **Syntax Analysis (Parsing):** Checking if the token stream conforms to the grammatical rules of the language and typically building an Abstract Syntax Tree (AST) to represent the program's structure.<sup>38</sup>
  - **Semantic Analysis:** Checking for semantic correctness (e.g., type compatibility, declaration of variables before use) and annotating the AST with type information and other semantic details.<sup>38</sup> The output of the front end is an **Intermediate Representation (IR)** of the program. The IR is a lower-level, language-independent (but often still machine-independent) form of the program.<sup>33</sup>
- **Middle End (Optimizer):** This part takes the IR from the front end and performs various transformations to improve its efficiency. These optimizations are generally independent of the specific target CPU architecture.<sup>33</sup> Examples include dead code elimination, constant propagation, loop optimizations like loop-invariant code motion, and common subexpression elimination.
- **Back End (Code Generator):** This is the part of the compiler that is highly dependent on the target CPU architecture (ISA).<sup>33</sup> It takes the (optimized) IR from

the middle end and translates it into the target machine's assembly language or directly into executable machine code. Key tasks performed by the back end include:

- **Instruction Selection:** Choosing the appropriate sequence of machine instructions from the target ISA to implement the operations specified in the IR.<sup>33</sup> This can be complex, as a single IR operation might have multiple possible machine code translations with different performance characteristics.
- **Register Allocation:** Assigning program variables and temporary values to the limited number of CPU registers available in the target architecture.<sup>33</sup> Efficient register allocation is crucial for performance, as register access is much faster than memory access. If there are not enough registers, some values must be "spilled" to memory (stack).
- **Instruction Scheduling:** Reordering machine instructions to optimize for the target CPU's pipeline structure, aiming to maximize instruction-level parallelism and minimize stalls caused by data dependencies or resource conflicts.<sup>33</sup>

The quality and sophistication of a compiler's backend are therefore critical for exploiting the full potential of the target hardware. It is in the backend that the compiler's detailed "knowledge" of the CPU's ISA and microarchitectural features (like pipeline depth, number of execution units, cache hierarchy, and specific instruction latencies) is applied to generate highly optimized code.

#### Ensuring Binary Compatibility and Portability

By generating code that adheres strictly to a specific ISA, the compiler ensures that the resulting executable program will run correctly on any CPU that implements that ISA.<sup>34</sup> The use of an IR is also a key factor in compiler design and portability. Because the front end (language-specific) and middle end (optimizations on IR) are largely independent of the target machine, a single front end (for a language like C++) and middle end can be paired with multiple different back ends, each targeting a different ISA (e.g., one for x86-64, one for ARMv8, one for RISC-V).<sup>33</sup> This modular design significantly reduces the effort required to support multiple programming languages across multiple hardware platforms.

#### **B. Handling CPU-Specific Optimizations: Unleashing Hardware Potential**

Beyond generating correct machine code for a given ISA, a primary goal of modern compilers is to optimize this code to enhance performance, reduce its size, and, increasingly, improve power efficiency.<sup>45</sup> Compilers employ a vast arsenal of sophisticated techniques to transform the initial, straightforward translation of source code into highly efficient machine instructions that take full advantage of the target

CPU's specific features and microarchitecture.

The Goal: Enhancing Performance (Speed, Code Size, Power Efficiency)

A compiler optimization is a transformation applied to a piece of code that results in a functionally equivalent piece of code but with improved characteristics, most commonly execution speed or reduced code size.<sup>45</sup> While it might be tempting for programmers to manually optimize critical sections of code by writing complex or low-level constructs, it is generally more effective to write clear, understandable, and maintainable high-level code and rely on the compiler's optimization capabilities.<sup>45</sup> Modern compilers are often better equipped to perform complex transformations and can consider interactions between different optimizations that a human programmer might overlook. In fact, premature manual optimization can sometimes obscure the code's intent from the compiler, hindering its ability to apply even more effective automatic optimizations.<sup>45</sup>

Key Optimization Categories (many are interrelated)

Compilers perform optimizations at various stages, often on the Intermediate Representation (IR) and during the final code generation phase. Some key categories include:

- **Strength Reduction:** Replacing computationally expensive operations with equivalent but cheaper ones. For example, replacing multiplication by a power of two with a bit shift (e.g.,  $x * 4$  becomes  $x \ll 2$ ), or transforming multiplication involving a loop counter into an additive update within the loop (e.g., if  $y = i * 1234$  is inside a loop where  $i$  increments by 1,  $y$  can be updated by adding 1234 in each iteration).<sup>47</sup> Integer division by a constant can also be replaced by a sequence of multiplications and shifts.<sup>47</sup>
- **Function Inlining:** Replacing a call to a function with the actual body of that function at the call site.<sup>38</sup> This eliminates the overhead of the function call (setting up a stack frame, passing parameters, branching) and, more importantly, often exposes the function's body to further optimizations in the context of the caller (e.g., constant propagation from call arguments into the inlined code). The trade-off is increased code size, which can sometimes negatively impact instruction cache performance if overused.
- **Constant Folding and Propagation:**
  - **Constant Folding:** Evaluating expressions whose operands are known constants at compile time and replacing the expression with its result (e.g.,  $x = 2 + 3$ ; becomes  $x = 5$ ;).<sup>38</sup>
  - **Constant Propagation:** If a variable is assigned a constant value, the compiler replaces subsequent uses of that variable with the constant itself, potentially enabling further constant folding or other optimizations.<sup>38</sup>
- **Common Subexpression Elimination (CSE):** Identifying identical subexpressions that are computed multiple times and whose operands do not

change between computations. The subexpression is computed once, its result saved, and then reused, avoiding redundant calculations.<sup>38</sup>

- **Dead Code Elimination:** Removing code that has no effect on the program's output or is determined to be unreachable.<sup>38</sup> This reduces code size and avoids executing unnecessary instructions. This can include unused variable assignments, entire functions that are never called, or blocks of code guarded by conditions that are always false.
- **Instruction Selection:** As part of the backend, the compiler chooses the most efficient sequence of machine instructions from the target ISA to implement a given operation or sequence of IR operations.<sup>33</sup> This requires knowledge of instruction latencies, throughput, and encoding sizes for the target CPU.
- **Peephole Optimizations:** Examining a small "window" (the peephole) of adjacent generated machine instructions and replacing them with a shorter or faster sequence that achieves the same result.<sup>33</sup> For example, a store instruction followed immediately by a load from the same location might be eliminated.
- **Tail Call Removal/Optimization:** A tail call occurs when a function's last action is to call another function (or itself). If a function ends with a call to itself (tail recursion), this can often be transformed into an iterative loop by the compiler.<sup>46</sup> This avoids creating a new stack frame for each recursive call, preventing stack overflow for deep recursions and reducing function call overhead.

Many of these optimization problems, when considered formally, are NP-hard (meaning finding a perfect, globally optimal solution is computationally intractable for all but trivial cases).<sup>33</sup> Consequently, compilers rely on sophisticated heuristics—rules of thumb and approximation algorithms—and detailed models of CPU behavior to make good optimization decisions within a reasonable compilation time. This reliance on heuristics and CPU models means that the effectiveness of optimizations can vary between different compilers (e.g., GCC, Clang, MSVC) and even between different versions of the same compiler, as these heuristics and models are continually refined.<sup>49</sup> It also implies that code optimally compiled for one specific CPU microarchitecture might not be equally optimal for another, even if they share the same ISA.

#### Leveraging CPU Microarchitecture

Beyond general ISA-level optimizations, compilers can achieve significant performance gains by tailoring code to the specific microarchitectural features of the target CPU. This requires the compiler to have a model of the CPU's internal workings.

- **Instruction Scheduling for Pipelined Execution:** Modern CPUs use instruction



pipelines to execute multiple instructions concurrently, with each instruction passing through several stages (fetch, decode, execute, memory access, write-back). To keep this pipeline full and operating efficiently, compilers perform instruction scheduling.<sup>33</sup> This involves reordering machine instructions in a way that preserves the original program semantics but minimizes pipeline stalls. Stalls can occur due to data dependencies (e.g., an instruction needing the result of a previous, not-yet-completed instruction – Read-After-Write or RAW dependency), resource conflicts (e.g., too many instructions needing the same execution unit simultaneously), or control flow changes (branches). The scheduler considers instruction latencies (how long an instruction takes to complete) and the availability of CPU functional units to find an optimal or near-optimal ordering. This aims to maximize Instruction-Level Parallelism (ILP).

- **Branch Prediction Optimization:** Conditional branches (if-statements, loops) can disrupt the smooth flow of instructions through a pipeline because the CPU doesn't know which path to take until the condition is evaluated. To mitigate this, CPUs employ branch predictors that try to guess the outcome of a branch before it's resolved, allowing speculative execution down the predicted path.<sup>50</sup> If the prediction is correct, performance is maintained. If incorrect (a misprediction), the speculatively executed instructions must be discarded, and the pipeline flushed and refilled from the correct path, incurring a significant performance penalty (10-30+ cycles).<sup>50</sup> Compilers can help improve branch prediction accuracy by:
  - Arranging code so that the more frequently executed path of a conditional branch is the "fall-through" path (the one taken if the branch is not taken), which is often easier for predictors to handle.
  - Using compiler hints provided by the programmer (e.g., `__builtin_expect` in GCC/Clang, or C++20 attributes `[[likely]]` and `[[unlikely]]`) to inform the compiler about the probable outcome of a branch.<sup>50</sup> The compiler can then optimize code layout accordingly.
  - Employing Profile-Guided Optimization (PGO), where the program is compiled with instrumentation, run with typical workloads to gather data on branch behavior, and then recompiled using this profile data to make more informed optimization decisions, including branch ordering.<sup>45</sup>
- **Loop Optimizations:** Since programs often spend a significant portion of their execution time in loops, optimizing loops is crucial for performance.<sup>46</sup> Compilers employ a wide array of loop transformations:
  - **Loop Unrolling:** Duplicates the loop body multiple times within a single

iteration of the modified loop, reducing the overhead of loop control instructions (incrementing and testing the loop counter, branching) and potentially exposing more instructions for parallel execution or other optimizations.<sup>46</sup>

- **Loop Fusion (Jamming/Combining):** Combines two or more adjacent loops that iterate over the same range into a single loop, reducing loop overhead and potentially improving data locality.<sup>38</sup>
- **Loop Fission (Distribution):** Splits a single loop with a complex body into multiple simpler loops, each handling a part of the original body. This can improve data locality for different parts of the computation or enable vectorization for specific parts.<sup>48</sup>
- **Loop-Invariant Code Motion (LICM):** Identifies computations within a loop whose results do not change from one iteration to the next (loop-invariants) and moves them outside the loop to be computed only once before the loop starts.<sup>46</sup>
- **Loop Interchange:** Swaps nested loops (e.g., making an inner loop an outer loop and vice-versa). This can significantly improve cache performance by changing the memory access patterns, especially for multi-dimensional array processing.<sup>48</sup>
- Other techniques include loop parallelization (restructuring loops for multi-core execution), loop scheduling, loop skewing (for nested loops with dependencies), loop splitting/peeling (handling initial/final iterations differently), and loop reversal.<sup>48</sup>
- **SIMD Instructions and Auto-Vectorization:**
  - Modern CPUs include Single Instruction, Multiple Data (SIMD) execution units and registers (e.g., MMX, SSE, AVX, AVX-512 on x86; NEON on ARM).<sup>49</sup> A single SIMD instruction can perform the same operation (e.g., addition, multiplication) on multiple data elements (packed into wide SIMD registers) simultaneously. This is a form of data-level parallelism (DLP).
  - Compilers can perform **auto-vectorization**, which involves analyzing loops (typically inner loops) and automatically transforming them to use SIMD instructions if it's safe and deemed beneficial.<sup>53</sup> This often requires the loop iterations to be independent and the data to be laid out appropriately. Auto-vectorization is usually enabled by default at higher optimization levels (e.g., -O2 or -O3 with flags like -vec in some compilers).<sup>53</sup>
  - The compiler needs to know which SIMD instruction sets are available on the target CPU to generate appropriate code. Flags like -xHost (Intel compilers) or

-march=native (GCC/Clang) instruct the compiler to generate code optimized for the instruction set of the machine performing the compilation.<sup>53</sup> Other flags like -xfeature (Intel) or specific -march=<architecture> flags allow targeting particular CPU features or generating multiple code paths for different CPU capabilities (dynamic dispatch).<sup>53</sup>

#### Importance of Specifying Target CPU Architecture to the Compiler

Providing the compiler with precise information about the target CPU architecture (e.g., via command-line flags like -march=skylake or -mtune=zen2) is crucial for achieving optimal performance.<sup>44</sup> This allows the compiler to:

1. **Utilize specific ISA extensions:** Generate code that uses advanced instructions (e.g., AVX2, AVX-512, FMA - Fused Multiply-Add) available on that CPU but perhaps not on older models.
2. **Tune for microarchitectural parameters:** Optimize instruction scheduling based on the known pipeline depth, number and type of execution units, instruction latencies, and cache sizes/associativity of that specific microarchitecture. Without this specific information, the compiler might generate generic code for a baseline ISA, or tune for a generic microarchitecture, potentially leaving significant performance untapped on newer or more specialized CPUs.<sup>47</sup>

The relationship between clear, well-structured source code and the compiler's ability to optimize is symbiotic. While compilers are powerful, overly complex or "tricky" manual optimizations by the programmer can sometimes confuse the compiler and prevent it from applying more effective transformations.<sup>45</sup> The "golden rule" is often to write readable and maintainable code that clearly expresses the program's logic, and then to guide the compiler with appropriate flags and, where necessary, profile data (PGO) to achieve the best performance.<sup>47</sup> This acknowledges the tension between code clarity for human understanding and the aggressive transformations compilers might perform, which can sometimes make the generated machine code appear very different from the original source, complicating debugging.

Furthermore, there's an ongoing "arms race" between CPU hardware advancements and compiler optimization technology. As CPU manufacturers introduce new features (wider SIMD units, more sophisticated branch predictors, novel instructions, heterogeneous cores), compiler developers must continually adapt and create new optimization strategies to effectively exploit these features. This means that the performance benefits of new hardware are not always realized automatically; they

often depend on using up-to-date compilers that are aware of and can target these new capabilities. Using an older compiler with newer hardware might result in suboptimal performance because the compiler is unaware of the hardware's full potential.

**Table: Key Compiler Optimization Techniques and Their Impact**

Optimization Technique	Primary Goal(s)	How it Works (Briefly)	Snippet Refs
Inlining	Reduce call overhead, enable further opts	Replaces function call with function body	38
Loop Unrolling	Reduce loop overhead, improve ILP	Duplicates loop body, reducing iterations/tests	46
Loop-Invariant Code Motion	Reduce redundant computations	Moves calculations constant within a loop outside the loop	46
Strength Reduction	Replace expensive ops with cheaper ones	E.g., <code>i*4</code> to <code>i&lt;&lt;2</code> ; <code>x/const</code> to <code>x*reciprocal_const</code>	47
Dead Code Elimination	Reduce code size, remove useless work	Removes code that doesn't affect output or is unreachable	38
Common Subexpression Elim.	Avoid recomputing same values	Calculates an expression once, reuses result	38
Instruction Scheduling	Improve pipeline usage, avoid stalls,	Reorders instructions respecting data	33

	increase ILP	dependencies for target microarchitecture	
Register Allocation	Minimize memory access, speed up data access	Assigns frequently used variables/temporaries to CPU registers	33
Auto-Vectorization (SIMD)	Exploit Data-Level Parallelism (DLP)	Transforms loops to use SIMD instructions operating on multiple data elements at once	49
Branch Prediction Optimization	Reduce branch misprediction penalties	Arranges code (e.g., via hints, PGO) to make branch outcomes more predictable	50

*Sources: As listed in the table.*

## Conclusion: The Interplay and Its Impact on Modern Computing

The intricate dance between the Operating System, the Compiler, and the CPU forms the very bedrock of modern computation. This report has elucidated the critical mechanisms that govern their interactions, revealing a sophisticated synergy that enables complex software to harness the power of underlying hardware efficiently and securely.

The OS acts as the CPU's primary steward, establishing distinct privilege levels (user and kernel modes) to protect system integrity and manage resources. It handles asynchronous events through hardware interrupts, allowing for responsive I/O and preemptive multitasking. Central to this multitasking capability is context switching, a process where the OS saves and restores CPU states to create the illusion of concurrent execution, albeit with inherent performance overhead. The runtime structure for program execution is meticulously managed through stack frames, which encapsulate the state of individual function calls, and standardized calling conventions, which ensure seamless interoperability between different code modules.

The compiler, in turn, serves as the crucial bridge between human-readable source code and the CPU's native instruction set. It translates high-level languages into ISA-specific machine code, a process heavily reliant on its backend components for instruction selection, register allocation, and instruction scheduling. More profoundly, the compiler is an engine of optimization, applying a vast array of techniques—from strength reduction and function inlining to sophisticated loop transformations and auto-vectorization for SIMD units—to tailor the generated code to the specific microarchitectural nuances of the target CPU. This optimization process is vital for unlocking the hardware's full performance potential.

This triad is not static; it is in a state of continuous evolution. Advances in CPU architecture, such as the proliferation of multi-core designs, the introduction of new ISA extensions (like advanced SIMD capabilities or hardware support for virtualization and security), and the development of deeper, more complex pipelines, invariably drive corresponding innovations in OS design and compiler technology. Operating systems must adapt their scheduling algorithms to manage heterogeneous cores effectively and provide new abstractions for emerging hardware features. Compilers must develop new analysis techniques and optimization passes to exploit these new CPU capabilities, leading to more intelligent code generation for instruction scheduling, register allocation, and automatic parallelization.

Understanding these fundamental concepts offers tangible benefits to both software developers and system designers.

- **For Software Developers:** A grasp of OS-CPU-compiler interactions empowers developers to write more efficient and robust code. Knowing the cost associated with system calls and context switches can inform decisions about algorithmic design and concurrency patterns. Understanding memory layout (stack vs. heap) and the mechanics of stack frames aids in debugging and avoiding vulnerabilities like buffer overflows. Awareness of how compilers optimize code (and what can hinder optimization) allows developers to write "compiler-friendly" code and effectively use compiler flags, intrinsics, or profile-guided optimization to achieve better performance. This knowledge is invaluable for diagnosing complex performance bottlenecks or elusive crashes.
- **For System Designers:** A deep understanding of these interactions informs critical decisions regarding hardware selection, OS configuration, and the choice of development tools, including compilers. It enables system-level performance tuning by identifying bottlenecks that may lie at the interface of these

components. Architectural decisions, such as how to partition tasks across processes and threads, or how to design inter-process communication, can be made more effectively when the underlying mechanics are well understood.

In essence, the seamless execution of billions of operations per second on our computing devices is a testament to this silent, intricate choreography. The complex interplay between the CPU's raw processing power, the OS's meticulous management and abstraction, and the compiler's intelligent translation and optimization is largely hidden from the end-user and often even from the high-level application programmer. Yet, it is this very sophistication that underpins the high-performance, secure, and versatile computing experiences that define the modern era.

### Works cited

1. What is the difference between user and kernel modes in operating ..., accessed June 4, 2025, <https://stackoverflow.com/questions/1311402/what-is-the-difference-between-user-and-kernel-modes-in-operating-systems>
2. How are different privilege levels (user mode and kernel mode) implemented in a processor, if at all? At the data path level, is anything different? - Quora, accessed June 4, 2025, <https://www.quora.com/How-are-different-privilege-levels-user-mode-and-kernel-mode-implemented-in-a-processor-if-at-all-At-the-data-path-level-is-anything-different>
3. Difference Between User Mode and Kernel Mode | GeeksforGeeks, accessed June 4, 2025, <https://www.geeksforgeeks.org/difference-between-user-mode-and-kernel-mode/>
4. 2. Privilege Levels - Learning | Linux Journey, accessed June 4, 2025, <https://linuxjourney.com/lesson/kernel-privilege-levels>
5. Dual Mode Operation in OS- Know Types & Why Switching is Required - Testbook, accessed June 4, 2025, <https://testbook.com/operating-system/dual-mode-operation-in-os>
6. Dual Mode Operations in OS- Scaler Topics, accessed June 4, 2025, <https://www.scaler.com/topics/dual-mode-operation-in-os/>
7. System Calls, Signals, & Interrupts - CS 3410 - CS@Cornell, accessed June 4, 2025, [https://www.cs.cornell.edu/courses/cs3410/2025sp/notes/syscall\\_signal\\_interrupt.html](https://www.cs.cornell.edu/courses/cs3410/2025sp/notes/syscall_signal_interrupt.html)
8. System Calls Explained: Bridging Programs and OS Tasks - Lenovo, accessed June 4, 2025, <https://www.lenovo.com/us/en/glossary/what-is-system-call/>
9. System call - Wikipedia, accessed June 4, 2025,



- [https://en.wikipedia.org/wiki/System\\_call](https://en.wikipedia.org/wiki/System_call)
10. User Mode vs Kernel Mode - Tutorialspoint, accessed June 4, 2025,  
<https://www.tutorialspoint.com/User-Mode-vs-Kernel-Mode>
  11. docs.oracle.com, accessed June 4, 2025,  
<https://docs.oracle.com/cd/E19253-01/816-4854/interrupt-15678/index.html#:~:text=Interrupt%20Handler%20Overview,-An%20interrupt%20is&text=An%20interrupt%20tells%20the%20CPU,CPU%20suspends%20the%20current%20thread.>
  12. Chapter 8 Interrupt Handlers, accessed June 4, 2025,  
<https://docs.oracle.com/cd/E19253-01/816-4854/interrupt-15678/index.html>
  13. www.geeksforgeeks.org, accessed June 4, 2025,  
<https://www.geeksforgeeks.org/user-mode-and-kernel-mode-switching/#:~:text=This%20switching%20between%20user%20mode,new%20context%20into%20the%20processor.>
  14. How does an operating system manage system resources? - TutorChase, accessed June 4, 2025,  
<https://www.tutorchase.com/answers/ib/computer-science/how-does-an-operating-system-manage-system-resources>
  15. How the Operating System Manages Hardware Resources in a ..., accessed June 4, 2025,  
<https://dev.to/adityabhuyan/how-the-operating-system-manages-hardware-resources-in-a-complete-system-4pn7>
  16. Technical Guide To System Calls: Implementation And Signal ..., accessed June 4, 2025,  
<https://mohitmishra786.github.io/chessman/2025/03/31/Technical-Guide-to-System-Calls-Implementation-and-Signal-Handling-in-Modern-Operating-Systems.html>
  17. What Is Context Switching in Operating System | LambdaTest, accessed June 4, 2025, <https://www.lambdatest.com/blog/context-switching/>
  18. Context Switches in Operating Systems | Baeldung on Computer ..., accessed June 4, 2025, <https://www.baeldung.com/cs/os-cpu-context-switch>
  19. Mastering Context Switching: Key to Peak Operating System Efficiency, accessed June 4, 2025,  
<https://krishparekh.hashnode.dev/the-art-of-context-switching-maximizing-os-efficiency>
  20. Context Switch - YouTube, accessed June 4, 2025,  
<https://www.youtube.com/watch?v=vTgccrbYHYs&pp=0gcJCdgAo7VqN5tD>
  21. Context Switching - IBM, accessed June 4, 2025,  
<https://www.ibm.com/docs/en/zvm/7.3?topic=exits-context-switching>
  22. Context switch: what decides when? - Stack Overflow, accessed June 4, 2025,  
<https://stackoverflow.com/questions/46450058/context-switch-what-decides-when>
  23. The Role Of The Stack Frame In Function Calls - FasterCapital, accessed June 4, 2025,

- <https://fastercapital.com/topics/the-role-of-the-stack-frame-in-function-calls.html>
24. Function Call Stack in C | GeeksforGeeks, accessed June 4, 2025, <https://www.geeksforgeeks.org/function-call-stack-in-c/>
  25. www.math.utah.edu, accessed June 4, 2025, [https://www.math.utah.edu/docs/info/gdb\\_7.html#:~:text=The%20call%20stack%20is%20divided,which%20the%20function%20is%20executing.](https://www.math.utah.edu/docs/info/gdb_7.html#:~:text=The%20call%20stack%20is%20divided,which%20the%20function%20is%20executing.)
  26. Debugging with GDB - Examining the Stack - Math.Utah.Edu, accessed June 4, 2025, [https://www.math.utah.edu/docs/info/gdb\\_7.html](https://www.math.utah.edu/docs/info/gdb_7.html)
  27. Calling Conventions | Cratecode, accessed June 4, 2025, <https://cratecode.com/info/calling-conventions>
  28. Calling Conventions in C/C++ | GeeksforGeeks, accessed June 4, 2025, <https://www.geeksforgeeks.org/calling-conventions-in-c-cpp/>
  29. Calling Conventions Demystified - CodeProject, accessed June 4, 2025, <https://www.codeproject.com/Articles/1388/Calling-Conventions-Demystified>
  30. x64 Calling Convention | Microsoft Learn, accessed June 4, 2025, <https://learn.microsoft.com/en-us/cpp/build/x64-calling-convention?view=msvc-170>
  31. builtin.com, accessed June 4, 2025, <https://builtin.com/software-engineering-perspectives/compiler-vs-interpreter#:~:text=A%20compiler%20translates%20code%20written,time%20before%20the%20program%20runs.>
  32. Difference Between Compiler and Interpreter - GeeksforGeeks, accessed June 4, 2025, <https://www.geeksforgeeks.org/difference-between-compiler-and-interpreter/>
  33. Compiler - Wikipedia, accessed June 4, 2025, <https://en.wikipedia.org/wiki/Compiler>
  34. Instruction set architecture - Wikipedia, accessed June 4, 2025, [https://en.wikipedia.org/wiki/Instruction\\_set\\_architecture](https://en.wikipedia.org/wiki/Instruction_set_architecture)
  35. www.arm.com, accessed June 4, 2025, <https://www.arm.com/glossary/isa#:~:text=The%20ISA%20defines%20the%20supported,model%20of%20multiple%20ISA%20implementations.>
  36. What is Instruction Set Architecture (ISA)? – Arm®, accessed June 4, 2025, <https://www.arm.com/glossary/isa>
  37. Is machine code different for various devices? : r/AskComputerScience - Reddit, accessed June 4, 2025, [https://www.reddit.com/r/AskComputerScience/comments/19cmu5u/is\\_machine\\_code\\_different\\_for\\_various\\_devices/](https://www.reddit.com/r/AskComputerScience/comments/19cmu5u/is_machine_code_different_for_various_devices/)
  38. How Compilers Work | Baeldung on Computer Science, accessed June 4, 2025, <https://www.baeldung.com/cs/how-compilers-work>
  39. Phases of a Compiler - GeeksforGeeks, accessed June 4, 2025, <https://www.geeksforgeeks.org/phases-of-a-compiler/>
  40. Code generation (compiler) - Wikipedia, accessed June 4, 2025,

- [https://en.wikipedia.org/wiki/Code\\_generation\\_\(compiler\)](https://en.wikipedia.org/wiki/Code_generation_(compiler))
41. en.wikipedia.org, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Register\\_allocation#:~:text=In%20compiler%20optimization%2C%20register%20allocation.limited%20number%20of%20processor%20registers.](https://en.wikipedia.org/wiki/Register_allocation#:~:text=In%20compiler%20optimization%2C%20register%20allocation.limited%20number%20of%20processor%20registers.)
  42. Register allocation - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Register\\_allocation](https://en.wikipedia.org/wiki/Register_allocation)
  43. en.wikipedia.org, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Instruction\\_scheduling#:~:text=In%20computer%20science%2C%20instruction%20scheduling,on%20machines%20with%20instruction%20pipelines.](https://en.wikipedia.org/wiki/Instruction_scheduling#:~:text=In%20computer%20science%2C%20instruction%20scheduling,on%20machines%20with%20instruction%20pipelines.)
  44. Instruction scheduling - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Instruction\\_scheduling](https://en.wikipedia.org/wiki/Instruction_scheduling)
  45. Compilers - What Every Programmer Should Know About Compiler Optimizations, accessed June 4, 2025,  
<https://learn.microsoft.com/en-us/archive/msdn-magazine/2015/february/compilers-what-every-programmer-should-know-about-compiler-optimizations>
  46. Optimizing compiler - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Optimizing\\_compiler](https://en.wikipedia.org/wiki/Optimizing_compiler)
  47. Optimizations in C++ Compilers - ACM Queue, accessed June 4, 2025,  
<https://queue.acm.org/detail.cfm?id=3372264>
  48. Loop Optimization in Compiler Design | GeeksforGeeks, accessed June 4, 2025,  
<https://www.geeksforgeeks.org/loop-optimization-in-compiler-design/>
  49. different compilers ability to optimize SIMD instructions : r/cpp\_questions - Reddit, accessed June 4, 2025,  
[https://www.reddit.com/r/cpp\\_questions/comments/194t3s4/different\\_compilers\\_ability\\_to\\_optimize\\_simd/](https://www.reddit.com/r/cpp_questions/comments/194t3s4/different_compilers_ability_to_optimize_simd/)
  50. Branch Prediction: The Definitive Guide for High-Performance C++ ..., accessed June 4, 2025,  
<https://johnfarrier.com/branch-prediction-the-definitive-guide-for-high-performance-c/>
  51. Branch Prediction in C (using GCC) - Jyotiprakash's Blog, accessed June 4, 2025,  
<https://blog.jyotiprakash.org/branch-prediction-in-c-using-gcc>
  52. Loop optimization - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Loop\\_optimization](https://en.wikipedia.org/wiki/Loop_optimization)
  53. www.danysoft.com, accessed June 4, 2025,  
<https://www.danysoft.com/free/16%20-%20Best%20Practices%20for%20vectorization.pdf>
  54. Understanding CPU Microarchitecture to Increase Performance - YouTube, accessed June 4, 2025, <https://www.youtube.com/watch?v=rglmJ6Xyj1c>