

# Central Processing Unit Microarchitecture, Control Logic, and Pipelining Techniques

## 1. Understanding CPU Microarchitecture

The Central Processing Unit (CPU) is the cornerstone of modern computing, executing the instructions that form computer programs. While users interact with software applications, the underlying hardware's design dictates the efficiency and capability of these executions. A fundamental concept in CPU design is its microarchitecture, which bridges the gap between the abstract definition of what a processor should do and the concrete hardware that does it.

### 1.1. Defining Microarchitecture and its Relation to Instruction Set Architecture (ISA)

At the heart of any processor lies its **Instruction Set Architecture (ISA)**. The ISA is an abstract model that defines the processor's programming interface. It specifies the set of instructions the CPU can execute, the data types it supports, the registers available to programmers, the memory addressing modes, and the memory and exception models.<sup>1</sup> In essence, the ISA is the "contract" or specification that software (including operating systems and compilers) relies upon to interact with the hardware.<sup>1</sup> Examples of ISAs include Arm, x86, RISC-V, MIPS, and PowerPC.

**CPU Microarchitecture**, in contrast, refers to the specific hardware implementation that fulfills the ISA contract.<sup>2</sup> It is the detailed description of how the processor's functional units are organized and interconnected to execute the instructions defined by the ISA. The microarchitecture dictates how instructions are fetched, decoded, executed, and how data is managed and moved within the processor.<sup>2</sup> Key aspects determined by the microarchitecture include the pipeline structure, the design of arithmetic logic units (ALUs), the configuration of the cache memory hierarchy, and the control mechanisms.<sup>1</sup>

The distinction is crucial: a single ISA can be implemented by many different microarchitectures. For instance, various Arm processors, such as the high-performance Cortex-A series or the energy-efficient Cortex-M series, all adhere to versions of the Arm ISA but possess distinct microarchitectures tailored for different power, performance, and area (PPA) targets.<sup>1</sup> Similarly, Intel and AMD produce processors that implement the x86 ISA but feature vastly different internal microarchitectures. This separation allows for innovation in hardware design without

breaking software compatibility, as long as the microarchitecture correctly implements the ISA.

The relationship between ISA and microarchitecture is not merely unidirectional. While the microarchitecture serves to implement the ISA, advancements and limitations at the microarchitectural level can, in turn, influence the evolution of the ISA itself. For instance, the introduction of new ISA features, such as Arm's Scalable Vector Extension 2 (SVE2) for enhanced Digital Signal Processing (DSP) and Machine Learning (ML) capabilities <sup>1</sup>, necessitates corresponding microarchitectural innovations. These might include wider datapaths or specialized execution units to efficiently handle the new instructions. Conversely, a breakthrough in microarchitectural design, perhaps a novel cache hierarchy that drastically reduces memory latency, could make certain existing instruction types far more effective or even enable the introduction of new ISA features that were previously impractical due to performance constraints. This dynamic interplay underscores a co-evolutionary path where both the architectural contract and its physical realization drive progress in processor design.

**Table 1: Key Differences: CPU Architecture (ISA) vs. Microarchitecture**

Feature	CPU Architecture (ISA)	CPU Microarchitecture
Definition	Abstract model defining the processor's programming interface; the "what."	Concrete hardware implementation realizing the ISA; the "how."
Focus	Instruction set, registers, memory model, data types, exception handling.	Pipeline design, functional units (ALU, FPU), cache organization, control logic, interconnections, PPA trade-offs.
Level of Abstraction	Higher level, visible to programmers and compilers.	Lower level, internal hardware details, generally hidden from programmers.
Example Aspects	ADD, LOAD, STORE	5-stage pipeline vs. 14-stage

	instructions; general-purpose register count; virtual memory system.	superpipeline; size and associativity of L1 cache; branch prediction algorithm.
<b>Variability</b>	Standardized for compatibility (e.g., Armv8-A, x86-64).	Can vary significantly between different processors implementing the same ISA (e.g., different Intel Core generations).

## 1.2. Key Components and Design Considerations in Microarchitecture

Designing a CPU microarchitecture involves making critical decisions about various hardware components and their organization to achieve specific design goals, primarily balancing power, performance, and area (PPA).<sup>1</sup> Key components and considerations include:

- Functional Units:** These are the hardware blocks that perform specific operations. Major functional units include:
  - Arithmetic Logic Units (ALUs):** Perform arithmetic (addition, subtraction) and logical (AND, OR, NOT) operations.
  - Floating-Point Units (FPUs):** Handle operations on floating-point numbers.
  - Register Files:** Store temporary data and operands close to the execution units for fast access.
  - Load/Store Units:** Manage data transfers between the CPU registers and the memory hierarchy.
- Pipeline Structure:** Modern processors use instruction pipelining to improve performance by overlapping the execution of multiple instructions. Microarchitectural decisions include the number of pipeline stages (pipeline depth), and the tasks performed in each stage.<sup>1</sup> Deeper pipelines can allow for higher clock frequencies but may increase penalties from hazards like branch mispredictions.
- Memory Hierarchy:** This includes multiple levels of cache memory (L1, L2, L3) that store frequently accessed data and instructions closer to the CPU core to reduce memory access latency.<sup>1</sup> Microarchitects decide the size, associativity, replacement policies, and write policies for each cache level.
- Control Mechanisms:** The microarchitecture defines how instructions are fetched from memory, decoded into control signals, and how these signals

manage the datapath operations.<sup>2</sup> This includes the design of the control unit itself.

- **Branch Prediction:** To mitigate performance losses from control hazards (discussed later), sophisticated branch prediction units are incorporated to guess the outcome of branch instructions.
- **Parallelism Techniques:** Beyond pipelining, microarchitectures may implement other forms of parallelism, such as superscalar execution (issuing multiple instructions per clock cycle), out-of-order execution (executing instructions based on data availability rather than program order), and Single Instruction, Multiple Data (SIMD) units for data parallelism.
- **Customization and Extensibility:** Some ISAs, like RISC-V, are designed to be extensible, allowing for custom instructions and tailored microarchitectures. This is particularly relevant for specialized applications where designers can modify the microarchitecture to optimize for specific tasks, often using hardware description languages like Verilog.<sup>2</sup> This contrasts with licensing pre-designed, fixed RTL (Register Transfer Level) cores from IP providers, where microarchitectural modifications are typically impractical or prohibited.<sup>2</sup>

The choices made in these areas collectively define the processor's microarchitecture and significantly impact its overall efficiency and suitability for different applications, from low-power embedded systems to high-performance servers.<sup>1</sup>

## 2. The Core of Execution: Datapath and Control Path

Within any CPU microarchitecture, two fundamental interacting components enable instruction execution: the datapath and the control path (or control unit). These components work in concert to perform the fetch-decode-execute cycle that is the basis of computer operation.<sup>3</sup>

### 2.1. The Datapath: Hardware for Operations

The **datapath** consists of the functional hardware elements that perform the actual data processing operations within the CPU.<sup>3</sup> It is often described as the "brawn" of the processor, as it contains the circuitry for storing data, moving data, and performing calculations.<sup>3</sup> The datapath components are interconnected by buses and controlled by signals generated by the control unit.

Typical components found in a CPU datapath include <sup>3</sup>:

- **Program Counter (PC):** A register that holds the memory address of the

instruction to be fetched next. It is updated after each instruction fetch to point to the subsequent instruction, or loaded with a target address in the case of branches or jumps.

- **Instruction Register (IR):** A register that stores the instruction currently being decoded and executed. In multicycle datapaths, the IR holds the instruction fetched from memory for the duration of its execution across multiple clock cycles.<sup>3</sup>
- **Register File (RF):** A collection of general-purpose registers used for storing operands and results of operations. A typical register file has multiple read ports (e.g., two for reading two source operands simultaneously) and one or more write ports (for writing back a result).<sup>3</sup>
- **Arithmetic Logic Unit (ALU):** The core computational engine of the datapath. It performs arithmetic operations (like addition, subtraction, increment, decrement) and logical operations (like AND, OR, XOR, NOT) on its input operands.<sup>3</sup> The ALU also typically generates status flags (e.g., Zero, Carry, Overflow, Negative) that can be used by control logic for conditional operations.
- **Memory Interface Units:** These include registers and logic to interact with the memory system.
  - **Memory Address Register (MAR):** Holds the address of the memory location to be accessed (for read or write).
  - **Memory Data Register (MDR) or Memory Buffer Register (MBR):** Temporarily stores data being read from memory or data to be written to memory.<sup>3</sup> specifically mentions the MDR for multicycle datapaths.
- **Multiplexers (Muxes):** These are data selectors that route data from one of several input sources to a single output line. Muxes are extensively used in datapaths to select operands for the ALU, choose data to be written into the register file, or select the next value for the PC.<sup>3</sup>
- **Sign/Zero Extenders:** Hardware units that extend shorter data values (e.g., immediate operands encoded in an instruction) to the full width of the datapath (e.g., 16-bit immediate to 32-bit), preserving the sign for signed numbers or padding with zeros for unsigned numbers.<sup>3</sup>
- **Internal Buses and Interconnections:** These are the electrical pathways that allow data to be transferred between the various components of the datapath (e.g., from register file to ALU, from ALU to MDR, from MDR to register file).

The design of the datapath is directly influenced by the ISA, as it must support all the operations and addressing modes defined by the architecture.

## 2.2. The Control Path (Control Unit): Orchestrating Operations

The **control path**, more commonly known as the **control unit**, is the component responsible for directing and coordinating the activities of the datapath.<sup>3</sup> It is often referred to as the "brain" of the processor. The control unit interprets the instructions fetched from memory and generates a sequence of control signals that manage the flow of data through the datapath and command its components to perform the necessary operations in the correct order.<sup>3</sup>

The primary functions of the control unit include:

- **Instruction Decoding:** Interpreting the opcode of the current instruction to determine the operation to be performed and the resources required.
- **Control Signal Generation:** Producing the specific electrical signals that control the datapath elements. These signals include <sup>3</sup>:
  - Register read/write enable signals for the register file.
  - ALU operation selection codes (e.g., ALUOp) to specify the function the ALU should perform.
  - Select signals for multiplexers to route data correctly.
  - Memory read/write commands for the memory interface units.
  - Signals to control the updating of the Program Counter.
- **Sequencing Operations:** Ensuring that operations occur in the correct sequence, especially for instructions that require multiple clock cycles or for managing the stages in a pipelined processor.
- **Handling Exceptions and Interrupts:** Managing unexpected events or external requests by altering the normal flow of instruction execution.

The control signals generated by the control unit dictate, for example, which registers supply operands to the ALU, what operation the ALU performs, and where the ALU result is stored.<sup>3</sup>

## 2.3. Interaction and Data Flow between Datapath and Control Path

The datapath and control unit are tightly coupled and interact continuously during CPU operation. The general flow is as follows:

1. **Instruction Fetch:** The control unit signals the datapath to fetch an instruction from the memory location pointed to by the PC. The fetched instruction is placed into the IR (often via the datapath).
2. **Instruction Decode:** The instruction in the IR is fed to the control unit. The control unit decodes the instruction's opcode and other relevant fields (e.g.,

function code for R-type MIPS instructions <sup>3</sup>).

3. **Operand Fetch (if needed):** Based on the decoded instruction, the control unit may signal the datapath to read operands from the register file or to prepare immediate values.
4. **Execution:** The control unit issues signals to the datapath to perform the required operation. For example, it sends an operation code to the ALU and selects the appropriate inputs for the ALU via multiplexers. The datapath executes the operation.
5. **Result Storage (if needed):** The control unit signals the datapath to write the result of the operation back to a register in the register file or to memory.
6. **Status Feedback:** The datapath may send status signals (e.g., ALU flags like Zero, Overflow) back to the control unit. These flags can influence subsequent control decisions, particularly for conditional branch instructions where the next PC value depends on the outcome of a comparison.<sup>3</sup>

For example, executing an R-format arithmetic instruction (e.g., add \$t1, \$t0, \$s0) involves the control unit orchestrating the datapath to read registers \$t0 and \$s0, commanding the ALU to perform an addition, and then directing the result to be written into register \$t1.<sup>3</sup> Similarly, for a load instruction (e.g., lw \$t1, offset(\$t2)), the control unit directs the ALU to calculate the memory address (\$t2 + offset), signals a memory read operation, and then routes the data from memory to register \$t1.<sup>3</sup>

The precise timing and synchronization of these control signals are critical. Incorrect timing can lead to data corruption, erroneous computations, or race conditions within the datapath. Each datapath component (registers, ALU, memory) has specific timing requirements (e.g., setup time, hold time, propagation delay) that the control signals must respect. For instance, a register write-enable signal must be asserted at the correct clock edge and held for an appropriate duration relative to when the data to be written is stable at the register's input. This implies that control unit design is not merely about *what* signals to generate, but crucially *when* and for *how long* they should be active, especially in multi-cycle or pipelined implementations where operations are carefully choreographed across clock cycles.<sup>3</sup> The complexity of ensuring correct timing and synchronization significantly influences the design of the control unit itself.

## Table 2: Datapath vs. Control Path: Roles and Characteristics



Aspect	Datapath	Control Path (Control Unit)
Primary Role	Performs data processing operations; the "brawn."	Generates control signals to direct datapath operations; the "brain."
Key Components/Elements	PC, IR, Register File, ALU, Memory Interface (MAR, MDR), Muxes, Extenders, Internal Buses.	Instruction Decoder, Sequencing Logic, State Registers (in FSM-based control), Control Signal Generation Logic.
Nature of Operations	Data storage, data movement, arithmetic and logical computations.	Instruction interpretation, signal generation, timing and sequencing of operations.
Analogy	The muscles, bones, and organs of the body that perform actions.	The nervous system that coordinates the actions of the body.
Output	Processed data (results of computations, data fetched from memory). Status flags (e.g., ALU Zero flag).	Control signals (e.g., RegWrite, ALUOp, MemRead, MemWrite, Mux selects).

### 3. Implementing Control: The Role of Finite State Machines (FSMs)

The control unit, responsible for orchestrating the datapath, is fundamentally a sequential logic circuit. Its outputs (control signals) at any given time depend not only on its current inputs (like the instruction opcode) but also on its internal state, which represents the current step in processing an instruction. Finite State Machines (FSMs) provide a formal and structured model for designing such sequential control logic.<sup>6</sup>

#### 3.1. Fundamentals of FSMs in Digital Design

A Finite State Machine is an abstract computational model characterized by a finite



number of defined **states**. At any point in time, the FSM can be in exactly one of these states.<sup>7</sup> The machine transitions from one state to another based on its current state and a set of **inputs**. These changes are called **transitions**.<sup>6</sup> An FSM also produces **outputs**, which can depend on the current state and, in some types of FSMs, also on the current inputs. Key components of an FSM include <sup>6</sup>:

- **States:** A finite set representing different conditions or stages of operation.
- **Initial State:** The state the FSM starts in.
- **Inputs:** External signals or conditions that influence state transitions.
- **Outputs:** Signals generated by the FSM.
- **Transition Function:** Logic that determines the next state based on the current state and inputs.
- **Output Function:** Logic that determines the outputs based on the current state (and possibly inputs).

FSMs are widely used in digital system design to model and implement sequential behavior, such as in communication protocols, vending machines, traffic light controllers, and, critically, CPU control units.<sup>6</sup> A hardware implementation of an FSM typically requires a register to store the current state variables, a block of combinational logic to compute the next state, and another block of combinational logic to generate the outputs.<sup>7</sup>

### 3.2. Mealy vs. Moore Machines in Control Logic

There are two primary types of FSMs, distinguished by how their outputs are generated <sup>6</sup>:

- **Mealy Machine:** In a Mealy machine, the outputs depend on both the **current state** and the **current inputs**.<sup>6</sup> This means that output signals can change as soon as the inputs change, even within the same clock cycle if the inputs are asynchronous to the state transitions.
  - **Characteristics:** Potentially faster response to input changes. Can sometimes be implemented with fewer states than a Moore machine for the same functionality because different outputs can be generated from the same state by varying the inputs. However, the outputs can be less stable and may exhibit glitches if inputs change while the machine is in a particular state, especially if inputs are not synchronized with the clock.<sup>6</sup>
- **Moore Machine:** In a Moore machine, the outputs depend **only on the current state**.<sup>6</sup> The inputs affect the state transitions, but not the outputs directly. Outputs change only when the state changes, typically synchronized with a clock

edge.

- **Characteristics:** Outputs are more stable and predictable, remaining constant for the entire duration of a state (i.e., for one clock cycle in a clocked system). This makes them well-suited for controlling synchronous datapaths where control signals need to be stable throughout a clock cycle. Moore machines might require more states than Mealy machines to implement the same logic if outputs need to differ based on inputs within what would otherwise be a single state.<sup>6</sup>

For CPU control unit design, **Moore machines are often preferred** due to their output stability.<sup>6</sup> Control signals generated by a Moore FSM are valid for the entire clock cycle corresponding to the current state. This synchronicity simplifies the interface with datapath components, which typically latch data or perform operations on clock edges based on stable control signals. While Mealy machines can offer faster reaction times, the potential for asynchronous output changes (glitches) can complicate timing and lead to unreliable operation if not carefully managed.<sup>6</sup>

**Table 3: Comparison of Mealy and Moore FSMs for CPU Control**

Characteristic	Mealy Machine	Moore Machine
<b>Output Dependency</b>	Current State AND Current Inputs	Current State ONLY
<b>Output Timing</b>	Can change asynchronously with input changes	Changes synchronously with state transitions (clocked)
<b>State Count (General)</b>	Potentially fewer states	Potentially more states
<b>Stability of Output</b>	Less stable; can have glitches	More stable; constant within a state
<b>Design Complexity</b>	Can be more complex due to input/output interaction	Generally simpler output logic
<b>Suitability for CPU Control</b>	Less common due to glitch	Often preferred for

	potential	synchronous datapath control
--	-----------	------------------------------

### 3.3. FSMs for Generating CPU Control Signals

In the context of a CPU control unit, each state of the FSM typically corresponds to a specific step or phase in the execution of an instruction (e.g., instruction fetch, decode, operand fetch, execute, memory access, write-back).<sup>3</sup> The outputs generated by the FSM in a particular state are the precise control signals required by the datapath to perform the actions associated with that step.<sup>8</sup>

For instance, in a multicycle CPU design (where an instruction takes multiple clock cycles to execute), the control unit is an FSM where each state corresponds to one clock cycle of an instruction's execution path.<sup>3</sup>

- **State Register:** Holds the current state (e.g., `FETCH_STATE`, `DECODE_STATE`, `EXECUTE_R_TYPE_STATE`, `MEM_ACCESS_LOAD_STATE`).
- **Next-State Logic:** Takes the current state and inputs like the instruction's opcode (from the IR) and status flags (from the ALU) to determine the next state. For example, after decoding an instruction, the FSM transitions to different execution states based on whether it's an R-type, I-type, load, store, or branch instruction.
- **Output Logic:** Based on the current state (for a Moore machine), this logic generates all necessary control signals. For example, in `MEM_ACCESS_LOAD_STATE`, it would assert `MemRead`, set the memory address source multiplexer, and configure the destination for the loaded data.

A concrete example can be seen in some MIPS control unit designs, where the opcode and current state bits are used as inputs to ROMs or Programmable Logic Arrays (PLAs). These combinational logic blocks then output the next state bits and the full set of datapath control signals, such as `ALUSrcA`, `ALUSrcB` (ALU input mux controls), `ALUOp` (ALU operation), `RegWrite` (register file write enable), `MemRead`, `MemWrite`, and `PCWrite` (PC update control).<sup>8</sup>

### 3.4. Hardwired vs. Microprogrammed Control Units

There are two main approaches to implementing the FSM logic for a CPU control unit<sup>5</sup>:

- **Hardwired Control Unit:**

The FSM is implemented directly using combinational logic circuits (e.g., AND, OR, NOT gates, decoders, multiplexers) and state registers (flip-flops).<sup>5</sup> The control signals are generated by this fixed logic.

- **Advantages:** Generally faster because control signals are generated through direct logic paths with minimal delay.
- **Disadvantages:** Less flexible. Modifying the control logic (e.g., to add new instructions or fix bugs) requires redesigning and refabricating the hardware. The design can become very complex and error-prone for ISAs with many instructions or complex instruction formats.
- **Typical Use:** Often found in RISC (Reduced Instruction Set Computer) processors, where the instruction set is simpler and more regular, leading to less complex control logic that can be efficiently hardwired.<sup>5</sup>
- **Microprogrammed Control Unit:**

The control signals are not generated directly by complex combinational logic. Instead, they are stored as binary words, called microinstructions, in a special high-speed memory called the Control Memory (CM) or Control Store (CS), which is often a Read-Only Memory (ROM).<sup>9</sup> Each microinstruction typically specifies the control signals to be asserted in one clock cycle and also contains information to determine the address of the next microinstruction to be fetched (the next "state").

  - **Components:**
    - **Control Memory (CM/CS):** Stores the microprogram (collection of microinstructions).
    - **Control Address Register (CAR) or Microprogram Counter ( $\mu$ PC):** Holds the address of the current microinstruction in the CM.
    - **Microinstruction Register ( $\mu$ IR):** Holds the microinstruction fetched from the CM. The bits of the  $\mu$ IR directly or indirectly provide the control signals.
    - **Microsequencer (Next Address Generator):** Logic that determines the address of the next microinstruction to be loaded into the CAR. This can be based on the current microinstruction, instruction opcode, ALU flags, etc..<sup>9</sup>
  - **Operation:** Executing a machine instruction involves fetching a sequence of microinstructions from the CM. Each fetched microinstruction provides the control signals for one step (micro-operation) in the machine instruction's execution.
  - **Advantages:** More flexible than hardwired control. Changes to the control

logic or instruction set can often be made by modifying the microprogram in the CM, which is easier than redesigning hardware.<sup>9</sup> This approach simplifies the design of control units for complex instruction sets.

- **Disadvantages:** Generally slower than hardwired control due to the time required to fetch microinstructions from the CM.<sup>5</sup>
- **Typical Use:** Commonly found in CISC (Complex Instruction Set Computer) processors, which have large and complex instruction sets that would be very difficult to implement with hardwired logic.<sup>5</sup>
- **Variations:**
  - **Horizontal Microprogramming:** Each microinstruction is wide, with many bits, where each bit (or a small field) directly controls one or a few signals. This allows for more parallelism in micro-operations but requires wider CM.
  - **Vertical Microprogramming:** Microinstructions are narrower, with fields being encoded. These encoded fields must be decoded to generate the actual control signals. This saves CM space but may be slower due to the decoding step.<sup>9</sup>
  - **Nanoprogramming:** An additional level of indirection where microinstructions in a main control store point to even more basic nano-instructions in a nanostore. This can further reduce the total control store size but adds another level of memory access, potentially increasing latency.<sup>5</sup>

The choice between hardwired and microprogrammed control is a fundamental microarchitectural decision. It reflects a core trade-off driven by the complexity of the ISA. Simpler, regular ISAs characteristic of RISC architectures lend themselves to the speed and simplicity of hardwired control. The control logic for these ISAs is manageable enough to be implemented efficiently with gates and PLAs. In contrast, the intricate and varied instructions of CISC ISAs necessitate a more structured and flexible approach; microprogramming provides this by breaking down complex instructions into sequences of simpler micro-operations, making the control unit design more tractable and adaptable. This decision profoundly impacts performance, design effort, and the ability to evolve the processor's functionality.

## 4. Instruction Pipelining: Enhancing CPU Throughput

To improve the performance of CPUs, microarchitects employ various techniques to execute more instructions in a given amount of time. One of the most fundamental

and widely used techniques is **instruction pipelining**.

#### 4.1. The Concept of Pipelining

Instruction pipelining is a technique that implements instruction-level parallelism within a single processor.<sup>11</sup> It works by dividing the processing of an instruction into a sequence of smaller, independent steps or **stages**. Each stage is handled by a dedicated segment of hardware. Multiple instructions can be in different stages of execution simultaneously, much like an assembly line in a factory where different products are at various stations of completion at the same time.<sup>11</sup>

The primary goal of pipelining is to increase **instruction throughput** – the number of instructions completed per unit of time.<sup>13</sup> While the total time an individual instruction takes to pass through all pipeline stages (its **latency**) might be slightly longer than in a non-pipelined multicycle processor (due to pipeline register overhead), the overall execution time for a sequence of instructions is significantly reduced because instructions are completed at a much faster rate (ideally, one instruction per clock cycle in a balanced pipeline).<sup>11</sup>

To enable this parallel processing, **pipeline registers** (also called latches or buffers) are placed between adjacent stages.<sup>11</sup> At each clock cycle, each stage performs its task on the instruction it currently holds and passes its result (along with any necessary control information) to the pipeline register for the next stage. The next stage then picks up this data from the pipeline register in the subsequent clock cycle. These registers are crucial for isolating the stages and synchronizing their operation to a common clock.<sup>12</sup>

#### 4.2. The Classic Five-Stage RISC Pipeline

A common and illustrative example of pipelining is the classic five-stage pipeline found in many RISC (Reduced Instruction Set Computer) processors, such as early MIPS architectures.<sup>11</sup> These five stages are:

1. **Instruction Fetch (IF)**
2. **Instruction Decode & Register Fetch (ID)**
3. **Execute & Address Calculation (EX)**
4. **Memory Access (MEM)**
5. **Write Back (WB)**

The operations performed in each stage are detailed below.<sup>13</sup>

**Table 4: The Classic Five-Stage Pipeline: Operations per Stage**

Stage Abbreviation	Stage Name	Key Operations & Purpose	Datapath Components Involved (Examples)
<b>IF</b>	Instruction Fetch	Fetch the instruction from memory (instruction cache) using the Program Counter (PC). Increment PC ( $PC = PC + 4$ ). Store fetched instruction in IF/ID pipeline register.	PC, Instruction Memory/Cache, Adder (for PC increment), IF/ID Register.
<b>ID</b>	Instruction Decode & Register Fetch	Decode the instruction to determine operation and operands. Read source register values from the Register File. Sign-extend immediate values. Calculate branch target address (if branch). Generate control signals. Store info in ID/EX register.	ID/EX Register, Register File, Instruction Decoder, Sign Extender, Adder.
<b>EX</b>	Execute & Address Calculation	Perform ALU operation (for R-type/I-type). Calculate effective memory address (for Load/Store). Evaluate branch condition and select target PC (for branches). Store result/address in	EX/MEM Register, ALU, Multiplexers.



		EX/MEM register.	
<b>MEM</b>	Memory Access	Load data from data memory (for Load). Store data to data memory (for Store). Pass through data (for ALU/Branch). Store loaded data or pass-through in MEM/WB register.	MEM/WB Register, Data Memory/Cache.
<b>WB</b>	Write Back	Write result back to the Register File (for ALU instructions or Load instructions).	Register File, Multiplexer (to select data for write back).

#### 4.2.1. Instruction Fetch (IF)

In the IF stage, the CPU retrieves the next instruction to be executed from memory (typically an instruction cache for speed).<sup>15</sup> The **Program Counter (PC)** holds the address of this instruction. The instruction is fetched and stored in the **IF/ID pipeline register**. Concurrently, the PC is updated to point to the next instruction in sequence. For RISC architectures with fixed-length instructions (e.g., 4 bytes), this usually involves incrementing the PC by 4 ( $PC = PC + 4$ ).<sup>14</sup>

#### 4.2.2. Instruction Decode & Register Fetch (ID)

The instruction fetched in the IF stage moves to the ID stage via the IF/ID pipeline register. Here, the instruction is **decoded** by the control unit to determine the operation to be performed (e.g., add, load, branch) and to identify the operands.<sup>15</sup> If the instruction uses source registers, their values are read from the **Register File**.<sup>13</sup> If the instruction contains an immediate value, this value is typically **sign-extended** to the full datapath width. For branch instructions, the potential **branch target address** is often calculated in this stage by adding the sign-extended offset to the (already incremented) PC value.<sup>13</sup> All necessary control signals for subsequent stages are also generated. The decoded instruction, operand values, immediate value, and control signals are then passed to the **ID/EX pipeline register**.

#### 4.2.3. Execute & Address Calculation (EX)

The EX stage is where the actual computation or address calculation takes place.<sup>15</sup>

The **Arithmetic Logic Unit (ALU)** is the primary component in this stage.

- For **arithmetic or logical instructions** (R-type or I-type), the ALU performs the specified operation (e.g., addition, subtraction, AND, OR) on the operands received from the ID/EX register.<sup>13</sup>
- For **load or store instructions**, the ALU calculates the effective memory address by, for example, adding a base register value and an offset.<sup>13</sup>
- For **branch instructions**, the ALU may be used to evaluate the branch condition (e.g., comparing two registers for equality). If the condition is met, the branch target address (calculated in ID) is selected as the next PC value. The result of the ALU operation (or the calculated memory address) is passed to the **EX/MEM pipeline register**.

#### 4.2.4. Memory Access (MEM)

This stage handles interactions with the data memory (typically a data cache).<sup>15</sup>

- For **load instructions**, the effective address calculated in the EX stage is used to read data from memory. The fetched data is then placed in the **MEM/WB pipeline register**.<sup>16</sup>
- For **store instructions**, the data to be stored (which was read from the register file in ID and passed through EX) and the effective address are used to write to memory.
- For **ALU or branch instructions** that do not access data memory, this stage might be idle, or simply pass the ALU result (if any) through to the MEM/WB register.

#### 4.2.5. Write Back (WB)

In the final WB stage, the result of the instruction is written back into the **Register File**.<sup>15</sup>

- For **ALU instructions**, the result computed in the EX stage (and passed through MEM) is written to the destination register specified in the instruction.
- For **load instructions**, the data fetched from memory in the MEM stage is written to the destination register. Store and branch instructions typically do not write to the register file in this stage (stores write to memory in MEM; branches modify the PC).

The choice of pipeline depth is a significant microarchitectural decision. While the classic 5-stage pipeline is a common baseline, modern high-performance processors often employ much deeper pipelines, sometimes called "superpipelines".<sup>11</sup> For instance, Intel's Skylake microarchitecture features a 14-19 stage pipeline, and ARM's Cortex-A77 has a 13-stage pipeline.<sup>17</sup> Deeper pipelines divide the work of each instruction into smaller, simpler pieces per stage. Simpler logic in each stage generally has a shorter propagation delay, which allows the processor to be clocked at a higher frequency.<sup>11</sup> A higher clock frequency can, in turn, lead to higher overall performance, provided the pipeline can be kept full of useful instructions. However, this potential for increased clock speed comes with trade-offs. Deeper pipelines tend to increase the penalty associated with pipeline stalls, such as those caused by branch mispredictions, because more stages (and thus more instructions) need to be flushed or cleared. The latency of individual instructions (the time taken for one instruction to complete all stages) may also increase. Therefore, determining the optimal pipeline depth involves a complex balance between the benefits of higher clock rates and the increased costs and complexities of managing hazards in a deeper pipeline.

## 5. Navigating Challenges: Pipeline Hazards

While pipelining significantly enhances CPU throughput, it also introduces a set of challenges known as **pipeline hazards**. Hazards are conditions that prevent the next instruction in the instruction stream from executing during its designated clock cycle in the pipeline.<sup>13</sup> If not handled correctly, hazards can lead to incorrect program execution or, at best, require the pipeline to stall, thereby reducing its efficiency and negating some of the performance gains. There are three main types of pipeline hazards: structural, data, and control hazards.<sup>18</sup>

### 5.1. Defining Pipeline Hazards

A pipeline hazard occurs when an instruction cannot proceed to its next stage in the pipeline as planned due to a conflict or dependency related to a preceding instruction still in the pipeline.<sup>20</sup> Such situations disrupt the smooth flow of instructions and can force the pipeline to pause (stall) or take corrective actions to ensure correct execution.<sup>21</sup>

### 5.2. Structural Hazards: Resource Conflicts

**Structural hazards** arise when two or more instructions in the pipeline require simultaneous access to the same hardware resource, but the hardware is not designed to support such concurrent access.<sup>19</sup> This is essentially a resource conflict.

- **Examples:**

- A unified memory unit for both instruction fetches (IF stage) and data accesses for load/store instructions (MEM stage). If one instruction is in IF and another is in MEM, and both need the memory unit in the same clock cycle, a structural hazard occurs.<sup>19</sup>
- A single ALU that is needed by two different instructions in their respective EX stages at the same time (e.g., in a superscalar processor issuing multiple instructions, or if an ALU operation takes multiple cycles).<sup>18</sup>
- A register file with an insufficient number of read or write ports to service all concurrent requests from different pipeline stages (e.g., ID stage reading operands, WB stage writing a result).<sup>19</sup>

- **Resolution:**

- **Stalling:** The simplest solution is to stall one of the conflicting instructions, allowing the other to use the resource. The stalled instruction (and subsequent instructions) waits until the resource becomes available.<sup>18</sup> This, however, reduces pipeline throughput.
- **Resource Duplication:** A more effective, albeit more costly, solution is to duplicate the contended resource. For example, using separate instruction caches and data caches (a Harvard cache architecture or modified Harvard architecture) can eliminate memory conflicts between IF and MEM stages.<sup>19</sup> Providing multiple ALUs or designing register files with more read/write ports can also alleviate structural hazards. These are microarchitectural design choices that trade hardware cost for performance.

### 5.3. Data Hazards: Dependencies on Data

**Data hazards** occur when an instruction's execution depends on the result of a previous instruction that is still being processed in the pipeline and has not yet completed its operation or written its result to the register file or memory.<sup>13</sup> If not handled, the dependent instruction might use stale or incorrect data.

#### 5.3.1. Read After Write (RAW) or True Data Dependency

This is the most common type of data hazard.<sup>22</sup> A RAW hazard occurs when an instruction attempts to **read** an operand (e.g., from a register) before a preceding instruction has **written** its result to that same operand.<sup>18</sup>

- **Example:**

Code snippet

```
I1: ADD R0, R1, R2 ; Instruction I1 writes to R0
```

I2: SUB R4, R0, R3 ; Instruction I2 reads from R0

In a 5-stage pipeline, I2 will be in its ID stage (ready to read R0) when I1 is in its EX stage (R0 has not yet been updated by ADD). If I2 proceeds, it will read the old value of R0. The result of the ADD operation is typically not available until the end of I1's EX stage or even its WB stage.

- **Resolution:**

- **Stalling:** I2 can be stalled in the ID stage until I1 completes its WB stage.
- **Forwarding (Bypassing):** The result from I1's ALU (available at the end of its EX stage) can be directly "forwarded" to the input of I2's ALU in its EX stage, bypassing the register file. This is a common and effective solution discussed in detail later.

### 5.3.2. Write After Read (WAR) or Anti-Dependency

A WAR hazard occurs when an instruction attempts to **write** to a destination (register or memory location) before a preceding instruction has finished **reading** its original value from that same location.<sup>18</sup>

- **Example:**

Code snippet

I1: ADD R2, R1, R0 ; Instruction I1 reads from R0

I2: SUB R0, R3, R4 ; Instruction I2 writes to R0

If I2 were allowed to execute and write to R0 before I1 reads the original value of R0, I1 would use the incorrect, prematurely updated value.

- **Occurrence and Resolution:** WAR hazards are less common in simple, in-order pipelines like the classic 5-stage model because reads typically occur in early stages (ID) and writes in a late stage (WB). However, they can become an issue in processors with out-of-order execution (where instructions are not necessarily executed in program sequence) or when instructions have variable latencies and can complete out of order.<sup>19</sup> A common technique to resolve WAR hazards in such advanced processors is **register renaming**, where physical registers are dynamically allocated to logical registers, effectively eliminating the name conflict.<sup>17</sup>

### 5.3.3. Write After Write (WAW) or Output Dependency

A WAW hazard occurs when two instructions attempt to **write** to the same destination (register or memory location), and the second instruction might complete its write before the first one, leading to the wrong final value being stored.<sup>18</sup>

- **Example:**

Code snippet

I1: MUL R0, R1, R2 ; Instruction I1 (e.g., a slow multiply) writes to R0

I2: ADD R0, R4, R5 ; Instruction I2 (a fast add) also writes to R0

If I2 (ADD) completes and writes to R0 before I1 (MUL) does (e.g., in an out-of-order execution scenario or if MUL takes many more cycles than ADD), then the result of the ADD will be overwritten by the MUL, which is incorrect if I1 appears before I2 in the program order.

- **Occurrence and Resolution:** Like WAR hazards, WAW hazards are primarily a concern in pipelines with out-of-order execution or instructions with significantly different execution times.<sup>19</sup> **Register renaming** is also effective in resolving WAW hazards for registers.<sup>17</sup> For memory, ensuring writes are committed in program order is crucial.

#### 5.4. Control Hazards (Branch Hazards)

**Control hazards**, also known as branch hazards, arise from instructions that change the flow of program execution, such as conditional branches, unconditional jumps, and procedure calls/returns.<sup>18</sup> The problem is that the pipeline may have already fetched and started processing instructions sequentially following the control flow instruction before the outcome of that instruction (e.g., whether a branch is taken or not, or the target address of a jump) is known.

- **Example:**

Code snippet

I1: BEQ R1, R2, TARGET\_LABEL ; Branch if R1 == R2

I2: NEXT\_INSTR\_IN\_SEQ ; Fetched assuming branch not taken

...

TARGET\_LABEL:

I3: TARGET\_INSTR ; Should be fetched if branch is taken

...

The decision for BEQ (whether R1 equals R2) is typically made in the EX stage (or sometimes ID). By the time this decision is made, I2 (and possibly more sequential instructions) might already be in earlier pipeline stages (IF, ID). If the branch is actually taken, these fetched instructions (I2, etc.) are incorrect and must be discarded (flushed) from the pipeline, and the correct instruction (I3 from TARGET\_LABEL) must be fetched. This flushing and redirection create bubbles or stalls in the pipeline.<sup>23</sup>

- **Resolution:**
  - **Stalling (Pipeline Freeze):** The pipeline can be stalled after fetching a branch

instruction until its outcome and target address are determined. This is simple but introduces significant delays, as branches are frequent in programs.<sup>19</sup>

- **Branch Prediction:** The processor predicts the outcome of the branch (e.g., always predict not taken, or use more sophisticated dynamic prediction based on past behavior stored in a Branch Target Buffer - BTB) and speculatively fetches instructions from the predicted path.<sup>17</sup> If the prediction is correct, no time is lost. If incorrect, the speculatively fetched instructions are flushed, and the pipeline restarts from the correct path. Modern CPUs use highly accurate dynamic branch predictors.
- **Delayed Branch:** An ISA-level technique where one or more instruction slots immediately following a branch instruction (the "delay slots") are always executed, regardless of the branch outcome. The compiler tries to fill these slots with useful instructions that are independent of the branch or that would be executed anyway.<sup>19</sup> This approach was common in earlier RISC architectures but is less prevalent in modern deep pipelines due to the difficulty of finding enough useful instructions to fill potentially many delay slots.
- **Flushing:** The process of discarding incorrectly fetched instructions from the pipeline stages by converting them into NOPs or clearing the pipeline registers.<sup>23</sup>

The types and likelihood of hazards encountered are not static; they are influenced by microarchitectural choices. A simple, in-order pipeline primarily contends with RAW and control hazards. However, introducing advanced features like out-of-order execution to boost performance inherently creates more scenarios for WAR and WAW hazards, as instructions no longer necessarily complete or write results in their original program sequence.<sup>17</sup> Similarly, superscalar designs, which issue multiple instructions per cycle, increase the demand on hardware resources, making structural hazards more probable if not managed through sufficient resource duplication. This implies that as microarchitects strive for higher instruction-level parallelism, they must concurrently develop more sophisticated mechanisms for detecting and resolving the evolving landscape of pipeline hazards.

**Table 5: Overview of Pipeline Hazards: Types, Causes, and Primary Mitigation Approaches**

Hazard Type	Description/Cause	Example Scenario	Common Mitigation
-------------	-------------------	------------------	-------------------



		(brief)	Techniques
<b>Structural</b>	Hardware resource conflict; multiple instructions need the same resource in the same cycle.	Two instructions need memory access (IF vs. MEM) simultaneously with a single memory port.	Stalling, Resource Duplication (e.g., separate I/D caches, multiple ALUs).
<b>Data (RAW)</b>	Read After Write: Instruction tries to read an operand before a preceding instruction writes to it.	ADD R0,R1,R2 followed by SUB R3,R0,R4. SUB needs R0 from ADD.	Stalling, Data Forwarding (Bypassing).
<b>Data (WAR)</b>	Write After Read: Instruction tries to write to a location before a preceding instruction reads its old value.	ADD R2,R1,R0 followed by SUB R0,R3,R4. SUB writes R0 before ADD reads old R0.	Register Renaming (in OoO), Ensuring reads complete before writes by later instructions.
<b>Data (WAW)</b>	Write After Write: Two instructions write to the same location; writes must occur in program order.	MUL R0,R1,R2 (slow) followed by ADD R0,R3,R4 (fast). ADD might write R0 first.	Register Renaming (in OoO), Ensuring writes commit in program order.
<b>Control</b>	Branch/Jump instruction; next instruction address not known immediately, leading to fetching wrong path.	BEQ R1,R2,TARGET. Sequential instructions fetched before BEQ outcome is known.	Stalling, Branch Prediction, Delayed Branch, Flushing.

## 6. Ensuring Smooth Flow: Hazard Mitigation Techniques

To maintain the efficiency of a pipelined processor, it is crucial to detect and resolve

hazards effectively. The two primary hardware techniques for mitigating data hazards are data forwarding (or bypassing) and pipeline stalling (or inserting bubbles). Control hazards are typically managed through branch prediction and flushing.

### 6.1. Data Forwarding (Bypassing)

Data forwarding is a hardware technique that significantly reduces stalls caused by RAW data hazards.<sup>21</sup> The fundamental idea is to provide the result of an instruction directly to a subsequent dependent instruction as soon as it is computed, rather than waiting for it to be written back to the register file in the WB stage and then read again in the ID stage of the dependent instruction.<sup>24</sup>

Mechanism:

Forwarding involves adding extra data paths (wires) within the datapath that route a result from the output of a functional unit (e.g., ALU output at the end of the EX stage, or data from memory at the end of the MEM stage) directly to the input of a functional unit (typically the ALU inputs) for a subsequent instruction that is in an earlier pipeline stage (usually EX).<sup>22</sup>

**Hardware:**

- **Forwarding Paths:** These are physical connections from the output of pipeline registers (like EX/MEM and MEM/WB) or directly from functional unit outputs, back to the inputs of earlier stages, most commonly the ALU inputs in the EX stage.<sup>22</sup>
- **Multiplexers (Muxes):** These are placed at the inputs of the ALU (and potentially other units that need forwarded data). For each ALU input, a multiplexer selects its data source from several possibilities:
  1. The value read from the register file (passed via the ID/EX pipeline register – the normal path).
  2. The result from an instruction currently in the MEM stage (forwarded from the EX/MEM pipeline register's ALU output field).
  3. The result from an instruction currently in the WB stage (forwarded from the MEM/WB pipeline register's ALU output field or loaded data field).<sup>22</sup>

Control Logic: The Forwarding Unit

A dedicated piece of control logic, often called the Forwarding Unit, is responsible for detecting RAW hazards that can be resolved by forwarding and generating the select signals for these ALU input multiplexers.<sup>24</sup>

- **Inputs to the Forwarding Unit:** The forwarding unit continuously monitors the source and destination registers of instructions in different pipeline stages. For a typical 5-stage pipeline, its inputs would include <sup>26</sup>:

- ID/EX.RegisterRs, ID/EX.RegisterRt: The source register numbers for the instruction currently in the EX stage.
- EX/MEM.RegisterRd, EX/MEM.RegWrite: The destination register number and write-enable signal for the instruction currently in the MEM stage.
- MEM/WB.RegisterRd, MEM/WB.RegWrite: The destination register number and write-enable signal for the instruction currently in the WB stage.
- **Outputs from the Forwarding Unit:** The outputs are control signals (e.g., ForwardA, ForwardB) that drive the select lines of the ALU input multiplexers.<sup>26</sup> These signals tell the multiplexers which source to choose for each ALU operand.
- **Logic:** The forwarding unit's logic implements rules like the following (conceptual, specific MIPS conditions shown in Table 6):
  1. **Forward from EX/MEM to EX:** If the instruction in the MEM stage (EX/MEM.RegWrite is true and EX/MEM.RegisterRd is not zero) is writing to a register that the instruction in the EX stage needs as a source (EX/MEM.RegisterRd == ID/EX.RegisterRs or EX/MEM.RegisterRd == ID/EX.RegisterRt), then forward the ALU result from the EX/MEM pipeline register.
  2. **Forward from MEM/WB to EX:** If the instruction in the WB stage (MEM/WB.RegWrite is true and MEM/WB.RegisterRd is not zero) is writing to a register that the instruction in the EX stage needs, AND this dependency is not satisfied by forwarding from the EX/MEM stage (as the EX/MEM result is "newer"), then forward the result from the MEM/WB pipeline register.
  3. **No Forwarding:** If neither condition above is met, the ALU inputs come from the register file (via the ID/EX register).<sup>21</sup>

Table 6: Forwarding Conditions and Control Signals (Illustrative for MIPS-like EX stage sources)

(Assuming ForwardA/B values: 00=No Forward (from RegFile via ID/EX), 01=Forward from EX/MEM.ALUOut, 10=Forward from MEM/WB.Result)

Hazard Condition (Data Source for EX Stage's Rs or Rt operand)	Condition for ForwardA (ALU Input A)	Condition for ForwardB (ALU Input B)	ForwardA Value	ForwardB Value
1a. Result from	EX/MEM.RegWri	(Depends on Rt	01	(Varies)

instruction in MEM stage needed for Rs	te AND EX/MEM.RegisterRd!= 0 AND (EX/MEM.RegisterRd == ID/EX.RegisterRs )	conditions)		
<b>1b.</b> Result from instruction in MEM stage needed for Rt	(Depends on Rs conditions)	EX/MEM.RegWrite AND EX/MEM.RegisterRd!= 0 AND (EX/MEM.RegisterRd == ID/EX.RegisterRt )	(Varies)	01
<b>2a.</b> Result from instruction in WB stage needed for Rs (and not covered by 1a)	MEM/WB.RegWrite AND MEM/WB.RegisterRd!= 0 AND (MEM/WB.RegisterRd == ID/EX.RegisterRs ) AND NOT (1a condition met for Rs)	(Depends on Rt conditions)	10	(Varies)
<b>2b.</b> Result from instruction in WB stage needed for Rt (and not covered by 1b)	(Depends on Rs conditions)	MEM/WB.RegWrite AND MEM/WB.RegisterRd!= 0 AND (MEM/WB.RegisterRd == ID/EX.RegisterRt ) AND NOT (1b condition met for Rt)	(Varies)	10
<b>3.</b> No hazard for Rs or Rt	None of the above	None of the above	00	00

requiring forwarding (or data comes from RegFile)	conditions for Rs are met.	conditions for Rt are met.		
---	----------------------------	----------------------------	--	--

*Note: This table is illustrative. Actual forwarding logic also handles cases where an instruction in MEM might be a load, so the data comes from MEM/WB.LoadData instead of MEM/WB.ALUOut. Priority is given to the most recent result (EX/MEM over MEM/WB if both could supply the same register).*

Limitations (e.g., Load-Use Hazards):

Forwarding cannot eliminate all RAW data hazard stalls. A classic example is the load-use hazard. Consider a LW (load word) instruction followed immediately by an instruction that uses the loaded data:

Code snippet

I1: LW R1, 0(R2) ; Load data into R1 from memory

I2: ADD R3, R1, R4 ; Use R1 as a source

The data loaded by LW is available only *after* the MEM stage of I1. If I2 needs this data for its EX stage, even forwarding from the MEM stage output (i.e., from the MEM/WB pipeline register at the beginning of the WB stage of I1) to the EX stage of I2 is too late. I2 would be in its EX stage when I1 is in its MEM stage. Thus, I2 must be stalled for one clock cycle to wait for the loaded data to become available for forwarding.<sup>13</sup>

## 6.2. Pipeline Stalling (Inserting Bubbles)

When a hazard cannot be resolved by forwarding (like the load-use hazard) or for structural or control hazards, the pipeline must be **stalled**. Stalling involves temporarily halting the progression of instructions in some pipeline stages, effectively inserting one or more "do-nothing" cycles, often called **bubbles** or **NOPs (No-Operations)**, into the pipeline.<sup>20</sup>

Mechanism and Control:

A Hazard Detection Unit is responsible for identifying conditions that necessitate a stall.<sup>23</sup> This unit often works in conjunction with or as part of the main control unit.

- **Inputs to Hazard Detection Unit (for load-use hazard):**
  - IF/ID.RegisterRs, IF/ID.RegisterRt: Source registers of the instruction currently in the ID stage (which is a candidate for stalling if it depends on a load).
  - ID/EX.MemRead: A control signal indicating if the instruction currently in the EX stage is a memory read (load) operation.
  - ID/EX.RegisterRt (or ID/EX.RegisterRd depending on ISA format): The destination register of the instruction in the EX stage (the register being loaded into).<sup>23</sup>
- **Stall Implementation (Control Actions):** When a stall is required (e.g., for a load-use hazard where ID/EX.MemRead is true, and ID/EX.RegisterRt matches IF/ID.RegisterRs or IF/ID.RegisterRt):
  1. **Freeze Early Stages:** The PC is prevented from being updated (PCWrite = 0), and the IF/ID pipeline register is prevented from being written (IF/ID.Write = 0). This effectively holds the instruction in the IF stage and the instruction in the ID stage (the one causing the stall detection) in their current stages for an extra cycle.<sup>23</sup>
  2. **Insert NOP (Bubble):** The control signals for the instruction in the ID stage (which is now stalled) are de-asserted or "squashed" as it (or rather, a NOP representing it for this cycle) enters the ID/EX pipeline register for the next clock cycle. This means the instruction in the EX stage during the stall cycle becomes a NOP, performing no operation and not writing any results.<sup>23</sup> In MIPS, a common NOP is sll \$0, \$0, 0 (shift left logical of register zero by zero, writing to register zero), which has an all-zero binary encoding and no side effects.<sup>23</sup>
- **Effect:** The NOP propagates through the EX, MEM, and WB stages, creating a one-cycle delay (a bubble). This delay allows the preceding instruction (e.g., the load) to complete its critical stage (MEM for a load) so its result can then be correctly forwarded to the previously stalled instruction, which now proceeds one cycle later.

Stalling for Control Hazards:

For control hazards, if a branch is mispredicted, or if the branch outcome is not known early enough, instructions fetched along the incorrect path must be flushed. Flushing is a form of stalling where the incorrectly fetched instructions in the early pipeline stages (IF, ID) are converted into NOPs.<sup>23</sup> For example, an IF.Flush control signal can be used to force the IF/ID pipeline register to load a NOP instruction.<sup>23</sup> This prevents them from having any effect on the processor state.

Impact on Performance:

Pipeline stalls directly reduce performance because, during a stall cycle, no useful instruction completes its execution.<sup>12</sup> Each bubble inserted increases the overall Cycles Per Instruction (CPI) of the program. While stalls are necessary for correctness, minimizing their frequency and duration is a key goal of pipeline design and hazard mitigation strategies like forwarding and branch prediction.

The interplay between forwarding and stalling is crucial. The control logic generally attempts to resolve data hazards using forwarding first, as this avoids bubbles. Stalling is then used as a fallback when forwarding is insufficient (e.g., load-use) or for other types of hazards. This implies a sophisticated decision-making process within the control unit, weighing different hazard conditions and applying the most efficient resolution strategy. While the common stall mechanism of freezing early stages and injecting NOPs is relatively simple to implement, it can be a blunt approach, halting the entire front-end of the pipeline. More advanced microarchitectures, particularly those with out-of-order execution, may employ more complex buffering and scheduling mechanisms to allow independent instructions to continue processing even when another instruction stream is stalled, thereby recovering some of the lost parallelism. This highlights a trade-off between the simplicity of the stall mechanism and the potential for higher performance through more fine-grained pipeline control.

## 7. Synergy of Concepts: An Integrated View

The concepts of CPU microarchitecture, datapath, control path, FSM-based control, pipelining, hazards, and mitigation techniques are not isolated topics but are deeply interconnected elements of processor design. Understanding their synergy is key to appreciating how a modern CPU functions.

The **microarchitecture** provides the high-level blueprint for how an **ISA** is implemented. This blueprint details the organization of the **datapath**—the collection of hardware units (ALU, registers, memory interfaces, muxes) that actually manipulate data—and the **control path** (control unit), which directs the datapath's actions.

The control unit itself is often designed as a **Finite State Machine (FSM)**, either hardwired for speed (common in RISC) or microprogrammed for flexibility (common in CISC). This FSM interprets instructions and generates a timed sequence of control signals. These signals orchestrate the datapath components, dictating data flow, ALU operations, and memory accesses.

**Pipelining** is a microarchitectural technique that divides instruction processing into stages (e.g., IF, ID, EX, MEM, WB) to increase instruction throughput. While the main



control unit still provides the basic control signals for each instruction type as it flows through these stages (often passed along via pipeline registers), the pipelined nature introduces **hazards**—structural, data, and control—that can disrupt the flow and lead to incorrect execution if not managed.

This is where specialized pipeline control logic, including **forwarding units** and **hazard detection units**, comes into play. These units, also part of the overall control system, monitor the state of instructions in various pipeline stages.

- When a **data hazard** (specifically RAW) is detected, the forwarding unit attempts to route the required data directly from a later pipeline stage (EX/MEM or MEM/WB) to an earlier stage (EX), bypassing the register file and avoiding a stall. This involves controlling multiplexers in the datapath.
- If forwarding is not possible or insufficient (as in a load-use hazard), or if a **structural hazard** or **control hazard** occurs, the hazard detection unit initiates a **pipeline stall**. This involves freezing earlier pipeline stages (e.g., by disabling PC and IF/ID register updates) and often inserting NOPs (bubbles) into the pipeline to create the necessary delay or to flush incorrectly fetched instructions.

The ultimate goal is a harmonious and efficient operation where the datapath performs the computational work, the main FSM-based control unit directs the fundamental sequence of operations for each instruction, and the specialized pipeline control logic (for hazard detection, forwarding, and stalling) dynamically intervenes to ensure correctness while maximizing the pipeline's throughput. The complexity of this interaction and the sophistication of these control mechanisms escalate significantly with more advanced microarchitectural features such as deeper pipelines, superscalar instruction issue, and out-of-order execution, all aimed at extracting ever-greater performance from the processor.

## 8. Conclusions

The design of a modern Central Processing Unit is a complex interplay of architectural specifications and detailed hardware implementation. The **microarchitecture** defines the specific organization of hardware components, such as the **datapath** and **control path**, to realize the functionality prescribed by the Instruction Set Architecture. The datapath provides the operational hardware, while the control path, often implemented using **Finite State Machines**, generates the precisely timed signals necessary to orchestrate these operations.

**Instruction pipelining** stands as a cornerstone technique for enhancing CPU throughput by overlapping the execution of multiple instructions. The classic five-stage pipeline (Instruction Fetch, Instruction Decode, Execute, Memory Access, and Write Back) provides a foundational model for understanding this parallelism. However, pipelining introduces **hazards**—structural, data (Read-After-Write, Write-After-Read, Write-After-Write), and control—which can impede performance and threaten correct execution.

Effective hazard mitigation is therefore critical. **Data forwarding (bypassing)** allows results to be routed directly between pipeline stages, significantly reducing stalls from many data dependencies. This requires additional datapath connections, multiplexers, and sophisticated forwarding unit control logic to detect dependencies and select appropriate data sources. When forwarding is insufficient, or for other hazard types, **pipeline stalling (bubble insertion)** is employed. This involves pausing parts of the pipeline and injecting NOP instructions, managed by a hazard detection unit that controls pipeline registers and program counter updates. While necessary for correctness, stalls inherently reduce pipeline efficiency.

The choice of microarchitectural features, such as pipeline depth and control unit implementation (hardwired vs. microprogrammed), involves intricate trade-offs between performance, power consumption, area, and design complexity. As processors evolve with deeper pipelines and more aggressive parallelism techniques like out-of-order execution, the challenges of hazard detection and resolution become even more pronounced, demanding increasingly sophisticated control mechanisms to maintain both correctness and high performance. A thorough understanding of these fundamental concepts is essential for anyone involved in the design, analysis, or optimization of computer systems.

## Works cited

1. Arm CPU Architecture, accessed June 4, 2025, <https://www.arm.com/architecture/cpu>
2. Microarchitecture - Cudasip, accessed June 4, 2025, <https://codasip.com/glossary/microarchitecture/>
3. Organization of Computer Systems: Processor & Datapath - UF CISE, accessed June 4, 2025, <https://www.cise.ufl.edu/~mssz/CompOrg/CDA-proc.html>
4. 3.2.6. Control Path - Intel, accessed June 4, 2025, <https://www.intel.com/content/www/us/en/docs/programmable/683152/22-1/control-path.html>
5. Control unit – Knowledge and References - Taylor & Francis, accessed June 4,

2025,

[https://taylorandfrancis.com/knowledge/Engineering\\_and\\_technology/Electrical\\_%26\\_electronic\\_engineering/Control\\_unit/](https://taylorandfrancis.com/knowledge/Engineering_and_technology/Electrical_%26_electronic_engineering/Control_unit/)

6. Finite State Machine (FSM): A Comprehensive Guide - SiliconVeda, accessed June 4, 2025,  
<https://www.siliconveda.com/2024/10/finite-state-machine-fsm-comprehensive.html>
7. Finite-state machine - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Finite-state\\_machine](https://en.wikipedia.org/wiki/Finite-state_machine)
8. Finite State Machine: Control Implementation, accessed June 4, 2025,  
<http://class.ece.iastate.edu/arun/Cpre305/lectures/week09.pdf>
9. Microprogrammed Control Unit in Computer Organization - NxtWave, accessed June 4, 2025,  
<https://www.ccbp.in/blog/articles/microprogrammed-control-unit-in-computer-organization>
10. Microprogrammed Control Unit - BYJU'S, accessed June 4, 2025,  
<https://byjus.com/gate/microprogrammed-control-unit-notes/>
11. Instruction pipelining - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Instruction\\_pipelining](https://en.wikipedia.org/wiki/Instruction_pipelining)
12. Computer Organization and Architecture | Pipelining | Set 1 ..., accessed June 4, 2025,  
<https://www.geeksforgeeks.org/computer-organization-and-architecture-pipelining-set-1-execution-stages-and-throughput/>
13. www.cs.fsu.edu, accessed June 4, 2025,  
[https://www.cs.fsu.edu/~zwang/files/cda3101/Fall2017/Lecture7\\_cda3101.pdf](https://www.cs.fsu.edu/~zwang/files/cda3101/Fall2017/Lecture7_cda3101.pdf)
14. Classic RISC pipeline - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Classic\\_RISC\\_pipeline](https://en.wikipedia.org/wiki/Classic_RISC_pipeline)
15. CPU Pipeline Stages to Know for Intro to Computer Architecture, accessed June 4, 2025, <https://library.fiveable.me/lists/cpu-pipeline-stages>
16. Delving Deeper into the MIPS Pipeline - Jyotiprakash's Blog, accessed June 4, 2025, <https://blog.jyotiprakash.org/delving-deeper-into-the-mips-pipeline>
17. Advanced Pipelining: Techniques & Hazards | Advanced Computer ..., accessed June 4, 2025, <https://library.fiveable.me/advanced-computer-architecture/unit-3>
18. Pipelining Hazards and Performance - TAMU-CC Repository, accessed June 4, 2025,  
<https://tamucc-ir.tdl.org/bitstreams/d0f794d3-812b-43f0-8dfc-137161b504c4/download>
19. Pipeline Hazards | Computer Architecture - Witscad, accessed June 4, 2025,  
<https://witscad.com/course/computer-architecture/chapter/pipeline-hazards>
20. Pipeline stall - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Pipeline\\_stall](https://en.wikipedia.org/wiki/Pipeline_stall)
21. LECTURE 9 Pipeline Hazards, accessed June 4, 2025,  
[https://www.cs.fsu.edu/~zwang/files/cda3101/Fall2017/Lecture9\\_cda3101.pdf](https://www.cs.fsu.edu/~zwang/files/cda3101/Fall2017/Lecture9_cda3101.pdf)

22. Data Hazards and Forwarding | Advanced Computer Architecture ..., accessed June 4, 2025,  
<https://library.fiveable.me/advanced-computer-architecture/unit-3/data-hazards-forwarding/study-guide/UEfsh4VI6bvLQ7dI>
23. courses.cs.washington.edu, accessed June 4, 2025,  
<https://courses.cs.washington.edu/courses/cse378/09wi/lectures/lec13.pdf>
24. Operand forwarding - Wikipedia, accessed June 4, 2025,  
[https://en.wikipedia.org/wiki/Operand\\_forwarding](https://en.wikipedia.org/wiki/Operand_forwarding)
25. Video 51: Forwarding or Bypassing, CS/ECE 3810 Computer Organization - YouTube, accessed June 4, 2025,  
<https://www.youtube.com/watch?v=k3osCp8ppiE>
26. courses.cs.vt.edu, accessed June 4, 2025,  
<https://courses.cs.vt.edu/cs2506/Spring2013/Assignments/8/ForwardingUnit.pdf>