# Module 4: Arithmetic Logic Unit (ALU) Design

**Module Objective:** This module offers an exhaustive exploration into the design and fundamental operation of the Arithmetic Logic Unit (ALU), the central computational engine of a Central Processing Unit (CPU). It thoroughly covers the hardware implementation of basic arithmetic operations, with a particular focus on the intricate design of integer multiplication and division units. The module then transitions to a comprehensive discussion of floating-point number representation and the complexities of floating-point arithmetic, culminating in an in-depth analysis of the universally adopted IEEE 754 standard and its profound implications for numerical accuracy.

## 4.1 General ALU Design Principles

The **Arithmetic Logic Unit (ALU)** stands as the computational heart of any digital computer. It is a highly specialized combinational digital circuit, meaning its outputs are solely determined by its current inputs, without any memory of past inputs. The ALU is solely responsible for performing all the basic arithmetic and logical operations requested by the CPU's control unit, serving as the essential workhorse that executes the core computations of a program.

**ALU Function: Performing Arithmetic and Logical Operations**

The ALU's functional repertoire can be broadly categorized into two major groups:

1. **Arithmetic Operations:** These operations perform standard mathematical computations on numerical data. The fundamental operations almost universally supported by an ALU include:
   ○ **Addition:** This is the most basic arithmetic operation, summing two binary numbers. All other arithmetic operations often rely on or are derived from addition (e.g., subtraction via two's complement addition).
   ○ **Subtraction:** Finding the difference between two binary numbers. In modern ALUs, subtraction is typically implemented by taking the two's complement of the subtrahend (the number being subtracted) and then performing addition. For instance, A - B is computed as A + (-B), where -B is the two's complement of B.
   ○ **Increment:** Adding the value 1 to a number. This is a very common operation, frequently used for loop counters or memory address manipulation. It can be implemented as a specialized adder or by simply adding a constant 1.
   ○ **Decrement:** Subtracting the value 1 from a number. Similar to increment, this is frequently used for counters.
   ○ While more complex operations like multiplication and division can be achieved through repeated additions and subtractions, for performance reasons, modern CPUs often employ dedicated, specialized hardware units (like integer multipliers and dividers) that are separate from or highly

integrated with the main ALU for these operations, as they are too complex to be handled bit-serially by a general-purpose ALU without substantial performance penalties. However, the basic arithmetic capabilities (add/subtract) are always core to the ALU.

2. **Logical Operations:** These operations perform bitwise manipulations on binary data. Unlike arithmetic operations, they treat their inputs as collections of individual bits rather than numerical values. They are indispensable for tasks such as setting or clearing specific bits, extracting parts of a binary word, or performing comparisons. Common logical operations include:

   ○ **AND:** The bitwise logical AND operation. For each corresponding bit position in the two input operands, the output bit is 1 if and only if *both* input bits are 1; otherwise, it is 0. It is often used for "masking" or clearing specific bits (e.g., ANDing with 00001111 would clear the upper four bits of an 8-bit number).

   ○ **OR:** The bitwise logical OR operation. For each corresponding bit position, the output bit is 1 if *at least one* of the corresponding input bits is 1; otherwise, it is 0. It is often used for "setting" specific bits (e.g., ORing with 10000000 would set the most significant bit).

   ○ **NOT (Invert/Complement):** The bitwise logical NOT operation. This is a unary operation (takes only one input). It inverts every bit of its operand: a 0 becomes a 1, and a 1 becomes a 0. It is crucial for generating two's complement numbers (which involves inverting all bits and then adding 1).

   ○ **XOR (Exclusive OR):** The bitwise logical XOR operation. For each corresponding bit position, the output bit is 1 if the input bits are *different* (one is 0 and the other is 1); otherwise, it is 0. XOR can be used to toggle specific bits, or to compare two values for equality (if two numbers are identical, their XOR result is zero).

   ○ **Shift Operations:** These operations move the bits within an operand to the left or right by a specified number of positions. They are very efficient for multiplication and division by powers of 2.

      ■ **Logical Shift Left (LSL):** Shifts all bits to the left. The leftmost bit(s) are discarded, and the newly vacated positions on the right are filled with 0s. This is equivalent to multiplying an unsigned integer by 2N for N shifts.

      ■ **Logical Shift Right (LSR):** Shifts all bits to the right. The rightmost bit(s) are discarded, and the newly vacated positions on the left are filled with 0s. This is equivalent to dividing an unsigned integer by 2N for N shifts.

      ■ **Arithmetic Shift Right (ASR):** Shifts all bits to the right. The rightmost bit(s) are discarded. Crucially, the newly vacated positions on the left are filled with a copy of the *original sign bit* (the most significant bit). This preserves the sign of a signed (two's complement) integer while performing division by powers of 2.

   ○ **Rotate Operations:** These are similar to shifts, but bits shifted off one end "wrap around" and reappear at the other end of the number, meaning no bits are lost.

      ■ **Rotate Left (ROL):** Bits shift left, the leftmost bit moves to the rightmost position.

- **Rotate Right (ROR):** Bits shift right, the rightmost bit moves to the leftmost position.
- **Rotate with Carry (RC/RCR/RCL):** These rotations involve the CPU's Carry Flag (a status bit), effectively making the rotation operate on one more bit than the operand's width, which is useful in some multi-precision arithmetic operations.

**Inputs and Outputs of an ALU**

A typical ALU communicates with the rest of the CPU via several input and output lines:

- **Operands (A and B):** These are the primary data inputs to the ALU. An ALU is generally designed to perform operations on two operands simultaneously. The number of bits in these operands defines the **data path width** of the ALU (e.g., a 32-bit ALU has 32 input lines for A and 32 for B). These operands typically come from CPU registers or directly from memory (via data buses and temporary buffers).
- **Function Select Code (Opcode/Control Signals):** This is a set of control input lines that dictate to the ALU *which specific operation* it should perform at any given moment. These signals are generated by the CPU's control unit after it decodes a machine instruction. For example, if an ALU supports 16 different operations (e.g., ADD, SUB, AND, OR, LSL, etc.), it would require at least $\log_2(16)=4$ select lines. These lines drive internal multiplexers and enable/disable various functional units within the ALU.
- **Result (F):** This is the main output of the ALU, representing the outcome of the performed operation. Its width matches the data path width of the ALU (e.g., 32-bit result for a 32-bit ALU).
- **Status Flags (Condition Codes):** These are single-bit output signals (typically 1 bit each) that provide supplementary information about the result of the operation. These flags are crucial for implementing conditional branches and other control flow mechanisms in a program, allowing the CPU to make decisions based on the outcome of arithmetic or logical computations. The most common status flags include:
  - **Zero Flag (Z):** This flag is set to 1 if the result of the operation is exactly zero. It is cleared to 0 otherwise. This is critical for checking for equality or determining if a counter has reached zero.
  - **Carry Flag (C):** This flag is set to 1 if an arithmetic operation produced a carry-out from the most significant bit (for addition) or a borrow into the most significant bit (for subtraction). It is crucial for multi-precision arithmetic (adding numbers larger than the ALU's bit width) and for certain shift/rotate operations.
  - **Sign Flag (N):** This flag is set to 1 if the result of the operation is negative. In two's complement representation, this means the most significant bit (MSB) of the result is 1. It is cleared to 0 if the result is positive. This allows for checking the sign of a number after computation.
  - **Overflow Flag (V):** This flag is set to 1 if a signed arithmetic operation resulted in an **overflow**. Overflow occurs when the true result of an operation

exceeds the largest positive or smallest negative number that can be represented within the given number of bits using signed (two's complement) representation. This is distinct from a carry-out; an overflow indicates that the result is incorrect for signed interpretation, even if a carry did or did not occur. For example, adding two large positive numbers could result in a negative number if an overflow occurs.

**Basic Logic Gates as Building Blocks: AND, OR, NOT, XOR**

The ALU, at its most fundamental level, is composed of interconnected **digital logic gates**. These gates are the atomic components that perform elementary Boolean logic functions on binary inputs.

- **AND Gate:** Produces a 1 output only if *all* its inputs are 1. Otherwise, the output is 0.
- **OR Gate:** Produces a 1 output if *any* of its inputs are 1. Only if all inputs are 0 is the output 0.
- **NOT Gate (Inverter):** Takes a single input and produces its opposite. If the input is 0, the output is 1; if the input is 1, the output is 0.
- **XOR Gate (Exclusive OR):** Produces a 1 output if its inputs are *different*. If the inputs are the same (both 0 or both 1), the output is 0.

These basic gates are then combined in various configurations to construct more complex functional units within the ALU, such as full adders, multiplexers, decoders, and specialized bit-slice logic for each operation. For example, the logical AND unit of a 32-bit ALU would simply be 32 independent 2-input AND gates operating in parallel, one for each bit position.

**Full Adder and Ripple-Carry Adder: Basic Arithmetic Circuits**

The ability to add binary numbers is the cornerstone of all arithmetic operations within the ALU.

- **Half Adder:** This is the most elementary adding circuit. It takes two single binary inputs (A and B) and produces two outputs: a **Sum (S)** and a **Carry-out (Cout)**. It cannot accept a carry-in from a previous stage of addition. Its logic is:
  - S=AoplusB (XOR gate)
  - Cout=AcdotB (AND gate)
- **Full Adder:** The **Full Adder** is the fundamental building block for constructing multi-bit binary adders. It takes three single binary inputs: two data bits (A and B) and a **Carry-in (Cin)** from the less significant bit position. It produces two outputs: a **Sum (S)** and a **Carry-out (Cout)** to the next more significant bit position. A full adder can be built using two half adders and an OR gate. Its logic is:
  - S=AoplusBoplusCin
  - Cout=(AcdotB)+(Cincdot(AoplusB)) (This expression means: a carry-out is generated if both A and B are 1, OR if Cin is 1 AND one of A or B is 1.)
- **Ripple-Carry Adder (RCA):** The simplest and most straightforward way to construct an N-bit binary adder is by cascading N single-bit full adders. The Carry-out (Cout) of each full adder is directly connected as the Carry-in (Cin) to the immediately next more significant full adder.

- ○ **Operation:** The addition effectively proceeds bit by bit, starting from the least significant bit (LSB) position and propagating the carry signal sequentially towards the most significant bit (MSB). The sum bit ($S_i$) and carry-out bit ($Cout_i$) for bit position i cannot be fully determined until the carry-in ($Cin_i$) from the previous position (i−1) is available.
- ○ **Advantage:** This design is remarkably simple and requires a minimal amount of hardware (gates).
- ○ **Disadvantage: Speed is its major drawback.** The critical path (the longest delay path from input to output) is determined by the "ripple" effect of the carry signal. For an N-bit adder, the carry might have to propagate through all N stages. This means the sum and carry outputs for the most significant bits are not stable until the carries have cascaded through all the preceding stages. This cumulative delay, known as **carry propagation delay**, can significantly limit the clock speed of the CPU, especially for wide (e.g., 32-bit or 64-bit) adders.

**Look-Ahead Carry Adder: Improving Adder Speed**

To overcome the inherent speed limitation of the ripple-carry adder, more sophisticated adder designs were developed. The **Look-Ahead Carry Adder (LCA)** is a widely used technique to significantly accelerate the carry propagation process.

- ● **Motivation:** The core problem in the RCA is that the carry-in for a given stage ($C_i$) depends on the carry-out from the previous stage ($C_{i-1}$), which depends on $C_{i-2}$, and so on. This sequential dependency creates delay. The LCA aims to compute the carries for multiple stages *in parallel* or at least much faster, by generating them directly from the input bits.
- ● **Principle:** The LCA introduces two new signals for each bit position i:
  - ○ **Generate ($G_i$):** A carry is *generated* at position i if $A_i$ and $B_i$ are both 1. This carry originates at this stage, irrespective of any carry-in from previous stages. $G_i = A_i \cdot B_i$.
  - ○ **Propagate ($P_i$):** A carry is *propagated* through position i if a carry-in ($C_i$) would cause a carry-out ($C_{i+1}$). This happens if either $A_i$ or $B_i$ (or both) is 1. A common definition for $P_i$ is $P_i = A_i \oplus B_i$ (if only one input is 1, it will propagate a carry). Another common definition is $P_i = A_i + B_i$ (if at least one input is 1, it will propagate a carry). The choice often depends on the specific logic circuit.
  - ○ Using these G and P signals, the carry-out for any stage ($C_{i+1}$) can be expressed directly in terms of the initial carry-in ($C_0$) and the G and P signals of preceding stages.
    - ■ $C_1 = G_0 + (P_0 \cdot C_0)$
    - ■ $C_2 = G_1 + (P_1 \cdot C_1) = G_1 + (P_1 \cdot (G_0 + (P_0 \cdot C_0))) = G_1 + P_1 \cdot G_0 + P_1 \cdot P_0 \cdot C_0$
    - ■ $C_3 = G_2 + (P_2 \cdot C_2) = G_2 + P_2 \cdot G_1 + P_2 \cdot P_1 \cdot G_0 + P_2 \cdot P_1 \cdot P_0 \cdot C_0$
  - ○ Notice the key aspect: $C_2$, $C_3$, and higher carries are computed directly using a two-level AND-OR logic from the input bits ($A_i, B_i$) and the initial

C_0. There is no ripple delay for these carry computations. This parallel computation drastically reduces the overall delay.

- **Advantage:** Significantly faster addition performance, especially for wider data paths, because the carry logic is not sequential. This is critical for high-performance CPUs.
- **Disadvantage:** Requires considerably more complex hardware (a larger number of gates and more intricate interconnections) compared to a simple ripple-carry adder. For very wide adders (e.g., 64-bit), a single, monolithic LCA becomes too large and complex. Thus, practical LCAs are often implemented in a hierarchical manner, where smaller LCAs generate carries within blocks, and another level of look-ahead logic generates carries *between* blocks.

**Multi-bit ALUs: Combining Basic Units to Handle Wider Data Paths**

A complete ALU designed for a modern CPU is a highly integrated circuit capable of performing various operations on data words that are typically 8, 16, 32, or 64 bits wide.

- A multi-bit ALU is constructed by arranging multiple identical **single-bit ALU slices** in parallel. Each slice is responsible for processing one specific bit position of the input operands.
- Each single-bit ALU slice is a sophisticated combinational circuit itself. It incorporates the necessary logic gates (AND, OR, XOR, NOT) for logical operations, a full adder for arithmetic operations, and potentially components for shift/rotate operations for its specific bit position.
- All these single-bit slices receive the same **function select control signals** from the CPU's control unit. These signals simultaneously activate the appropriate functional units (e.g., enable the adders, disable the logical gates, or vice-versa) across all slices.
- At the output stage, a large **multiplexer** (or a network of multiplexers) receives the results from the various functional units (adder output, logical AND output, shifter output, etc.) within the ALU. The function select signals then control this multiplexer to route the correct computed result to the main ALU output.
- The **status flags** (Zero, Carry, Sign, Overflow) are derived from logic that monitors the bits of the final result and/or the carry signals from the most significant adder stage.
  - The Zero flag is set if all bits of the result are 0.
  - The Sign flag is directly the MSB of the result (for two's complement).
  - The Carry flag is the carry-out from the MSB of the adder/subtractor.
  - The Overflow flag is determined by examining the carries into and out of the MSB of the adder/subtractor. For signed addition, overflow occurs if (Carry-in to MSB) ne (Carry-out from MSB).

## 4.2 Integer Multiplication Design

Integer multiplication, at its core, is a process of repeated addition and shifting. While simple in concept, its efficient hardware implementation for speed and area is a complex design challenge.

**Basic Principles of Multiplication: Repeated Addition**

Let's consider multiplying two unsigned binary numbers: a Multiplicand (M) and a Multiplier (Q). If the multiplicand is N bits long and the multiplier is N bits long, the product can be up to 2N bits long.

The manual "long multiplication" method illustrates the principle:

   M (Multiplicand)

x   Q (Multiplier)

------------------

   P0 = M * Q0 (Q0 is the LSB of Q)

  P1 = M * Q1  (shifted left by 1)

  P2 = M * Q2   (shifted left by 2)

  ...

------------------

  Product

Each "partial product" (P0, P1, P2, etc.) is either zero (if the corresponding multiplier bit $Q_i$ is 0) or a shifted copy of the Multiplicand (if $Q_i$ is 1). The final product is the sum of all these partial products.

**Hardware Implementation of Unsigned Multiplication**

- **Array Multiplier (Combinational/Parallel Implementation):**
    - **Concept:** An array multiplier is a purely combinational circuit designed to compute the product in a single clock cycle (after a propagation delay through its gates). It achieves this by generating *all* partial products simultaneously and then summing them up in parallel using a dedicated network of adders.
    - **Structure:** For an N-bit multiplicand and an N-bit multiplier, an array multiplier consists of:
        - **N x N AND gates:** Each AND gate computes one bit of a partial product. Specifically, AND(M_j,Q_i) computes the jth bit of the ith partial product.
        - **A Matrix of Adders:** The outputs of these AND gates (the partial product bits) are then fed into a grid-like arrangement of full adders. These adders accumulate the partial products diagonally. The simplest form uses ripple-carry adders in each row, summing partial products sequentially. More advanced array multipliers use **Carry-Save Adders (CSAs)** to speed up the partial product accumulation. A CSA takes three inputs (two numbers and a carry-in) and produces two outputs (a

sum and a carry-out) without propagating the carry, thus delaying carry propagation until the final addition stage.
- ○ **Advantages:**
  - ■ **Extremely Fast:** The product is available after a fixed combinational logic delay. There are no iterative clock cycles involved once the inputs are stable. This makes them ideal for applications requiring very high throughput multiplication, such as in high-performance CPUs or Digital Signal Processors (DSPs).
- ○ **Disadvantages:**
  - ■ **High Hardware Cost:** The number of gates and adders grows quadratically with the input bit width (roughly N2). For a 32-bit multiplier, this means hundreds of AND gates and hundreds of full adders, consuming significant silicon area and power. This makes large array multipliers costly to implement.
- ○ **Analogy:** Imagine drawing a multiplication table on paper. An array multiplier essentially has dedicated hardware (an AND gate and adder) for *every single entry* in that table, calculating them all simultaneously.
- ● **Sequential Multiplier (Iterative/Sequential Implementation):**
  - ○ **Concept:** A sequential multiplier computes the product iteratively, usually over N clock cycles (for N-bit operands). It uses a single adder, some registers, and shifters, much like how you would perform long multiplication by hand. It reuses the same hardware components in each step.
  - ○ **Structure:** A typical sequential multiplier involves:
    - ■ **Multiplicand Register:** Stores the Multiplicand (M).
    - ■ **Multiplier Register (Q):** Stores the Multiplier (Q). This register is typically shifted right during the process.
    - ■ **Accumulator / Product Register (A):** A register, often twice the width of the operands, used to store the partial product being accumulated. It is often paired with the Multiplier register.
    - ■ **Adder:** A standard N-bit adder (e.g., a Look-Ahead Carry Adder) to perform the addition of the Multiplicand to the partial product.
    - ■ **Shifters:** To shift the Multiplier (right) and the partial product/accumulator (right).
    - ■ **Control Unit:** A finite state machine or sequential logic that controls the sequence of operations (add/shift) over N clock cycles.
  - ○ **Algorithm (Simplified for unsigned N-bit Multiplier):**
    - ■ Initialize Product/Accumulator register (A) to 0.
    - ■ Initialize Multiplier register (Q) with the Multiplier.
    - ■ For N clock cycles (or steps):
      a. Check LSB of Multiplier (Q0): Look at the least significant bit of the Multiplier register.
      b. Add/No Add:
      * If $Q\_0=1$, add the Multiplicand (M) to the Accumulator (A).
      * If $Q\_0=0$, do nothing (effectively add 0).
      c. Shift:
      * Shift the Multiplier register (Q) one bit to the right (the LSB is discarded, the next bit becomes the new LSB).
      * Shift the combined Accumulator-Multiplier register pair (A and Q

together) one bit to the right. The bit shifted out of A's LSB goes into Q's MSB. This efficiently merges the partial product with the shifting of the multiplier bits.
- After N cycles, the final product is contained in the combined Accumulator and Multiplier registers.
  - **Advantages:**
    - **Low Hardware Cost:** Reuses the same adder and registers multiple times, leading to a much smaller and more area-efficient hardware implementation compared to an array multiplier.
  - **Disadvantages:**
    - **Slower:** Takes N clock cycles to complete the multiplication (for N-bit operands). This means an N-bit sequential multiplier is N times slower than an N-bit array multiplier.

**Booth's Algorithm: Efficient Multiplication for Signed (Two's Complement) Numbers**

Booth's algorithm is a powerful technique for multiplying signed binary numbers, specifically those represented in **two's complement** format. It's often more efficient than standard sequential multiplication, particularly when the multiplier contains long strings of 0s or 1s.

- **Motivation:** Standard unsigned multiplication algorithms need special handling for signed numbers (e.g., converting to unsigned, multiplying, then adjusting the sign of the product). Booth's algorithm inherently handles two's complement and also aims to reduce the number of addition/subtraction operations. A sequence of '1's in a binary number (e.g., ...011110...) can be thought of as ...100000... - ...000010.... This observation is key.
- **Principle:** Booth's algorithm examines pairs of bits in the multiplier, including an implied bit 0 to the right of the least significant bit. This allows it to identify strings of 1s or 0s and perform fewer operations.
  - It uses the following rules, considering the current multiplier bit ($Q_i$) and the bit to its right ($Q_{i-1}$):
    - If $Q_i Q_{i-1}$ is **0 0**: No operation (shift only).
    - If $Q_i Q_{i-1}$ is **0 1**: Add Multiplicand (M) to the partial product. Then shift. This signifies the *start* of a string of 1s (e.g., ...01...).
    - If $Q_i Q_{i-1}$ is **1 0**: Subtract Multiplicand (M) from the partial product. Then shift. This signifies the *end* of a string of 1s (e.g., ...10...).
    - If $Q_i Q_{i-1}$ is **1 1**: No operation (shift only). This signifies being *within* a string of 1s.
  - The shifts involved are arithmetic right shifts to preserve the sign of the partial product.
- **Advantages:**
  - **Directly Handles Signed Numbers:** Multiplies two's complement numbers correctly without requiring separate sign handling or conversion to unsigned magnitudes.
  - **Potentially Faster:** Can significantly reduce the number of additions/subtractions compared to simple sequential multiplication,

especially when the multiplier contains long sequences of identical bits (e.g., `0000`, `1111`). For example, multiplying by `00111100` requires only two operations (subtract M, add M shifted) instead of four additions.
- **Disadvantages:**
    - **More Complex Control Logic:** The logic to implement Booth's algorithm (decoding the bit pairs and controlling add/subtract/shift operations) is more complex than a simple sequential multiplier.
    - **Performance Variability:** While it can be faster in some cases, for multipliers with alternating 0s and 1s (e.g., `10101010`), it might not offer much advantage and could even be slightly slower due to the overhead of the more complex control.


## 4.3 Integer Division Design

Integer division is fundamentally the inverse of multiplication, involving repeated subtraction. Like multiplication, efficient hardware implementation of division is complex and iterative.

**Basic Principles of Division: Repeated Subtraction**

Similar to manual long division, binary division involves iteratively determining quotient bits by seeing if the divisor can be subtracted from a portion of the dividend.

Let's divide a Dividend (N) by a Divisor (D) to get a Quotient (Q) and Remainder (R).

$N = Q \times D + R$, where $0 \le R < D$.

The process in binary long division involves:

1. Aligning the divisor with the most significant part of the dividend.
2. Attempting to subtract the divisor.
3. If successful (result is non-negative), the corresponding quotient bit is 1. The result becomes the new partial remainder.
4. If unsuccessful (result is negative), the corresponding quotient bit is 0. The partial remainder does not change (or is restored).
5. Shift the partial remainder left and repeat the process.

**Hardware Implementation of Unsigned Division**

Sequential division algorithms are generally used in hardware, mimicking the manual process over several clock cycles.

- **Restoring Division Algorithm:**
    - **Concept:** This algorithm directly follows the manual long division procedure. In each step, it subtracts the divisor from the current partial remainder. If the subtraction results in a negative value (meaning the divisor was too large for that partial remainder), the original partial remainder is "restored" (by adding the divisor back), and the quotient bit for that position is set to 0. If the subtraction results in a non-negative value, the result becomes the new partial remainder, and the quotient bit is 1.

- ○ **Structure:** Requires:
    - ■ **Registers:** A register for the Dividend, a register for the Divisor, a register for the Quotient, and a register to hold the Remainder (often initialized with the Dividend).
    - ■ **Adder/Subtractor:** To perform the subtraction (and restoration if needed).
    - ■ **Shifters:** To shift the Remainder left and the Quotient bits into position.
    - ■ **Control Unit:** To manage the iterative steps.
- ○ **Algorithm (Simplified for unsigned N-bit numbers):**
    - ■ Initialize the Remainder register (R) with the Dividend.
    - ■ Initialize the Quotient register (Q) to 0.
    - ■ Initialize the Divisor register (D) with the Divisor.
    - ■ For N iterations (where N is the number of bits):
      a. Left shift the Remainder register (R) by one bit. (This brings down the next bit of the original dividend).
      b. Perform a "trial subtraction": $R\_temp = R - D$.
      c. Check Result:
      * If $R\_temp \ge 0$ (non-negative): This means the divisor fits. Set the least significant bit of the Quotient register ($Q\_0$) to 1. Update the Remainder: $R = R\_temp$.
      * If $R\_{temp} \< 0$ (negative): This means the divisor doesn't fit. Set the least significant bit of the Quotient register ($Q\_0$) to 0. Restore the Remainder: $R = R + D$ (effectively undoing the trial subtraction).
      d. Left shift the Quotient register (Q) by one bit to prepare for the next quotient bit.
    - ■ After N iterations, the Quotient register holds the final quotient, and the Remainder register holds the final remainder.
- ○ **Advantages:** Conceptually simple and directly mimics manual long division, making it easier to understand.
- ○ **Disadvantages:**
    - ■ **Slower:** The "restoring" step (adding the divisor back if the subtraction was negative) adds an extra operation cycle whenever a trial subtraction fails. This can add significant time, especially if many trial subtractions result in negative numbers.
    - ■ Can take up to 2N operations (N shifts + N subtractions, potentially N restorations) in the worst case.
- ● **Non-Restoring Division Algorithm (More Efficient):**
    - ○ **Concept:** This algorithm improves upon restoring division by eliminating the "restoring" step, making it faster. Instead of undoing a negative result, it uses the negative remainder directly in the next iteration, but changes the operation to addition instead of subtraction.
    - ○ **Algorithm (Simplified for unsigned N-bit numbers):**
        - ■ Initialize Remainder register (R) with Dividend, Quotient register (Q) to 0, Divisor register (D) with Divisor.
        - ■ **Initial Step (outside loop):** Left shift R by one bit.
        - ■ For N iterations:
          a. Decide Operation:

\* If the Remainder (R) is currently negative (or was negative from the previous iteration, meaning the previous subtraction failed), ADD the Divisor (D) to the Remainder: R=R+D.

\* If the Remainder (R) is currently non-negative (or was non-negative from the previous iteration, meaning the previous subtraction succeeded), SUBTRACT the Divisor (D) from the Remainder: R=R−D.

b. Set Quotient Bit:

\* If the result of the operation in step (a) is non-negative: Set the LSB of the Quotient register (Q_0) to 1.

\* If the result of the operation in step (a) is negative: Set the LSB of the Quotient register (Q_0) to 0.

c. Left shift the Quotient register (Q) by one bit.

d. Prepare for next iteration: If this is not the last iteration, left shift the Remainder (R) by one bit.

- **Final Adjustment (after N iterations):** If the final Remainder (R) is negative, add the Divisor (D) back to it one last time to make it positive. This one final restoration is sometimes needed.
  - **Advantages:**
    - **Faster:** Eliminates the extra "restore" cycle. In each iteration, only one arithmetic operation (add or subtract) is performed, reducing the total execution time compared to restoring division.
  - **Disadvantages:**
    - **More Complex Logic:** The control logic to decide whether to add or subtract, and to handle the final remainder adjustment, is more intricate than restoring division.

**Signed Division Considerations: Handling Signs of Dividend, Divisor, Quotient, and Remainder**

Most hardware division units simplify the design by performing unsigned division internally. The signs of the final quotient and remainder are then determined based on the signs of the original dividend and divisor using predefined rules. This avoids the complexity of two's complement arithmetic within the core division algorithm.

- **Rules for Signs:**
  1. **Quotient Sign:**
     - If the original Dividend and Divisor have the same sign (both positive or both negative), the Quotient will be **positive**.
     - If the original Dividend and Divisor have different signs (one positive and one negative), the Quotient will be **negative**.
  2. **Remainder Sign:**
     - The sign of the Remainder is typically defined to be the **same as the sign of the original Dividend**. This is a convention that simplifies number theory consistency.
- **Typical Implementation Strategy for Signed Division:**

1. **Convert to Absolute Values:** Take the absolute (unsigned) value of both the Dividend and the Divisor. Store their original signs separately.
2. **Perform Unsigned Division:** Execute the unsigned division algorithm (either restoring or non-restoring) using these absolute values. This will yield an unsigned quotient and an unsigned remainder.
3. **Adjust Quotient Sign:** Apply the sign rule for the quotient. If the final quotient should be negative, convert the unsigned quotient result to its two's complement representation.
4. **Adjust Remainder Sign:** Apply the sign rule for the remainder. If the original Dividend was negative, convert the unsigned remainder result to its two's complement representation.

This approach ensures that the division logic itself only deals with positive numbers, and the sign correction is applied as a final step, simplifying the hardware design.

## 4.4 Floating Point Arithmetic

While integers are excellent for exact counting, they are inadequate for representing a vast range of numbers encountered in scientific, engineering, and graphical applications: numbers that are very large, very small, or contain fractional components. **Floating-point numbers** address this limitation by adopting a system analogous to scientific notation.

**Motivation for Floating Point Numbers: Representing Very Large, Very Small, and Fractional Numbers**

- **Representation of Fractional Values:** Unlike integers, floating-point numbers can accurately represent values with decimal (or binary) fractions, such as 3.14159, 0.001, or 2.718. This is indispensable for calculations that involve measurements, percentages, or non-whole quantities. Fixed-point numbers can represent fractions but have a limited range and fixed decimal point.
- **Representation of Very Large Numbers:** Floating-point numbers use an exponent to scale a base number, much like scientific notation ($M \times 10^E$). This allows them to represent extremely large magnitudes, such as the number of atoms in a mole ($6.022 \times 10^{23}$) or astronomical distances, which would overflow even a 64-bit integer.
- **Representation of Very Small Numbers:** Conversely, they can represent numbers incredibly close to zero, such as the mass of an electron ($9.109 \times 10^{-31} \text{ kg}$) or a tiny electrical current. These small values would underflow to zero in fixed-point or integer systems.
- **Dynamic Range:** The exponential scaling inherent in floating-point representation provides an enormous "dynamic range" – the ratio between the largest and smallest non-zero numbers that can be represented. This allows calculations to span many orders of magnitude while maintaining a relatively consistent level of *relative* precision across that range.

**Structure of a Floating Point Number: Sign, Exponent, Mantissa (Significand)**

A binary floating-point number in a computer is typically composed of three distinct parts, inspired by the scientific notation $S \times M \times B^E$ (Sign times Mantissa times Base $^{\text{Exponent}}$):

1. **Sign (S):** This is a single bit that indicates the polarity of the number.
    - 0 typically represents a positive number.
    - 1 typically represents a negative number.
2. **Exponent (E):** This field represents the power to which the **base** (almost always 2 for binary floating-point numbers) is raised. This exponent determines the "magnitude" of the number by "floating" the binary (or decimal) point. A large exponent shifts the binary point far to the right, making a large number; a large negative exponent shifts it far to the left, making a very small number.
3. **Mantissa (M) or Significand:** This field represents the significant digits or the "precision" of the number. It's the fractional part of the number, typically normalized to have a leading 1 (or 0, for special cases) before the binary point. The more bits allocated to the mantissa, the higher the precision of the floating-point number.

The numerical value of a floating-point number is generally calculated using the formula:

$$\text{Value} = (-1)^S \times \text{Mantissa} \times 2^{\text{TrueExponent}}$$

**Normalization: Standardizing the Mantissa**

**Normalization** is a crucial step in floating-point representation that ensures a unique binary representation for most numbers and maximizes the precision within the available bits.

- **Principle:** For a non-zero binary floating-point number, it is always possible (and desirable) to shift the mantissa bits and adjust the exponent such that the binary point is immediately to the right of the first non-zero bit. In binary, this means the mantissa will always have a leading '1' before the binary point (e.g., $1.xxxx_2$).
- **The "Implied Leading 1" (Hidden Bit):** Since the first bit of a normalized binary mantissa is *always* 1, there's no need to store it explicitly in memory. This "implied leading 1" (or "hidden bit") effectively gives an extra bit of precision for the mantissa without consuming any storage space. For example, if a mantissa field is 23 bits, with the implied 1, it provides 24 bits of precision.
- *Example:*
    - The binary number $101.11_2$ is equivalent to $1.0111_2 \times 2^2$.
    - The binary number $0.00101_2$ is equivalent to $1.01_2 \times 2^{-3}$.
      In both cases, the mantissa is shifted until it is in the form $1.xxxx..._2$. The '1' before the binary point is implied, and only the fractional part (xxxx...) is stored.

**Bias in Exponent: Representing Both Positive and Negative Exponents**

The exponent field in floating-point numbers typically uses a **biased representation** (also called "excess-K" or "excess-N" representation) rather than two's complement for handling both positive and negative exponents.

- **Motivation:** Standard binary number systems (like unsigned or two's complement) have specific ranges and complexities for signed comparisons. By adding a fixed "bias" value to the true exponent, the entire range of exponents (positive and negative) is mapped to a range of positive unsigned integers. This simplifies hardware design, particularly for comparing floating-point numbers.
- Principle: A constant "bias" value is chosen. The actual numerical exponent (the "true exponent") has this bias added to it before being stored in the exponent field.
  $$\text{StoredExponent} = \text{TrueExponent} + \text{Bias}$$
  To retrieve the true exponent:
  $$\text{TrueExponent} = \text{StoredExponent} - \text{Bias}$$
- *Example:* If the bias is 127 (as in IEEE 754 single-precision):
  - A true exponent of 0 would be stored as $0+127=127$.
  - A true exponent of +1 would be stored as $1+127=128$.
  - A true exponent of -1 would be stored as $-1+127=126$.
  - A true exponent of -126 (the minimum) would be stored as $-126+127=1$.
  - A true exponent of +127 (the maximum) would be stored as $127+127=254$.
- **Benefits of Biased Exponent:**
  - **Simplified Comparison:** When comparing two floating-point numbers, if their signs are the same, a simple unsigned integer comparison of their biased exponent fields (followed by mantissa comparison) will correctly determine which number is larger. This is because a larger biased exponent directly corresponds to a larger true exponent and thus a larger magnitude.
  - **No Special Handling for Negative Exponents:** The hardware logic for handling the exponent becomes simpler as it operates only on unsigned numbers.
  - **Consistent Sorting:** Numbers can be sorted lexicographically (like text strings) based on their sign, then biased exponent, then mantissa, which generally holds true for their numerical values (with caveats for negative numbers).

## 4.5 IEEE 754 Floating Point Formats

The **IEEE 754 standard** (formally ANSI/IEEE Std 754-1985, later revised as IEEE 754-2008 and IEEE 754-2019) is a cornerstone of modern computing. It is the universally accepted technical standard for floating-point computation, defining consistent representations and arithmetic operations across diverse computer systems and programming languages. Its adoption ensures that floating-point calculations produce predictable and reproducible results, which is critical for portability and reliability in numerical software.

**Single-Precision (32-bit) Format**

The IEEE 754 single-precision format uses a total of 32 bits to represent a floating-point number.

- **Bit Allocation:**
  - **Sign Bit (1 bit):** This is the most significant bit (bit 31).
    - 0 indicates a positive number.
    - 1 indicates a negative number.

- ○ **Exponent Field (8 bits):** These bits (from bit 30 down to bit 23) store the biased exponent.
    - ■ The **bias** for single-precision is **127**.
    - ■ The actual value of the true exponent is calculated as: $\text{True\_Exponent} = \text{Stored\_Exponent} - 127$.
    - ■ The range of stored exponents is $00000000_2$ (0) to $11111111_2$ (255). However, the values 0 (all zeros) and 255 (all ones) are reserved for special cases (explained below).
    - ■ Therefore, for normal numbers, the $\text{Stored\_Exponent}$ ranges from 1 to 254, meaning the $\text{True\_Exponent}$ ranges from $1-127=-126$ to $254-127=+127$.
  - ○ **Mantissa (Significand) Field (23 bits):** These bits (from bit 22 down to bit 0) store the fractional part of the mantissa.
    - ■ **Implied Leading 1:** For **normalized numbers** (the vast majority of representable numbers), there is an **implied leading 1** before the binary point. So, the actual mantissa value is $1.f\_22f\_21...f\_0$, where f_i are the bits stored in the mantissa field. This effectively gives a 24-bit precision ($1 \text{ implied bit} + 23 \text{ stored bits}$).
- ● **Range and Precision (for normalized numbers):**
  - ○ **Smallest Positive Normalized Number:** When the true exponent is -126 (stored as 1) and the mantissa is $1.00...0_2$. This results in approximately $1.18 \times 10^{-38}$.
  - ○ **Largest Positive Normalized Number:** When the true exponent is +127 (stored as 254) and the mantissa is $1.11...1_2$. This results in approximately $3.40 \times 10^{38}$.
  - ○ **Precision:** With an effective 24-bit mantissa, single-precision numbers can represent about **6 to 7 decimal digits** of precision reliably. This means if you write a decimal number with 7 significant digits, it can usually be represented exactly (or very close to exactly).
- ● **Special Values (defined by reserved exponent values):**
  - ○ **Zero ($\pm 0.0$):** Represented by an exponent field of **all zeros** (00000000) and a mantissa field of **all zeros**. The sign bit distinguishes between +0.0 and -0.0, though they typically compare as equal.
  - ○ **Infinity ($\pm \infty$):** Represented by an exponent field of **all ones** (11111111) and a mantissa field of **all zeros**. The sign bit indicates positive or negative infinity. Infinity results from operations like division by zero (e.g., 1.0/0.0).
  - ○ **NaN (Not a Number):** Represented by an exponent field of **all ones** (11111111) and a **non-zero mantissa field**. NaNs are used to represent the results of invalid or indeterminate operations, such as 0.0/0.0, $\infty - \infty$, or $\sqrt{-1}$. NaNs are "sticky" – once a NaN is produced, most operations involving it will also result in a NaN. There are two types: Quiet NaN (QNaN) and Signaling NaN (SNaN). QNaNs propagate without signalling, while SNaNs typically raise an exception when accessed.
  - ○ **Denormalized (or Subnormal) Numbers:** Represented by an exponent field of **all zeros** (00000000) and a **non-zero mantissa field**. Unlike normalized numbers, these numbers have an **implied leading 0** (i.e., $0.f\_22f\_21...f\_0 \times 2^{\text{True\_Exponent\_Min}}$). They are used to represent numbers very close to zero that would otherwise "underflow" directly to zero.

Denormalized numbers allow for "gradual underflow," meaning the precision gracefully degrades as numbers approach zero, which helps in preventing unexpected errors in certain algorithms. The smallest denormalized number is smaller than the smallest normalized number.

## Double-Precision (64-bit) Format

The IEEE 754 double-precision format uses 64 bits, offering a significantly wider range and much higher precision compared to single-precision.

- **Bit Allocation:**
  - **Sign Bit (1 bit):** Bit 63.
  - **Exponent Field (11 bits):** Bits 62-52.
    - The **bias** for double-precision is **1023**.
    - The Stored_Exponent ranges from 0 to 2047. Reserved values are 0 (for zeros and denormals) and 2047 (for infinities and NaNs).
    - For normal numbers, True_Exponent ranges from $1 - 1023 = -1022$ to $2046 - 1023 = +1023$.
  - **Mantissa (Significand) Field (52 bits):** Bits 51-0.
    - **Implied Leading 1:** Similar to single-precision, there is an implied leading 1 for normalized numbers, resulting in an effective 53-bit mantissa ($1 \text{ implied bit} + 52 \text{ stored bits}$).
- **Extended Range and Precision (for normalized numbers):**
  - **Smallest Positive Normalized Number:** Approximately $2.22 \times 10^{-308}$.
  - **Largest Positive Normalized Number:** Approximately $1.80 \times 10^{308}$.
  - **Precision:** With an effective 53-bit mantissa, double-precision numbers can represent about **15 to 17 decimal digits** of precision reliably. This makes them suitable for demanding scientific and engineering calculations where accuracy is paramount.

## Floating Point Arithmetic Operations

Floating-point arithmetic is considerably more involved and computationally intensive than integer arithmetic. This is due to the separate exponent and mantissa components, the need for alignment, normalization, and precise rounding. These operations are typically handled by a dedicated hardware unit called the **Floating-Point Unit (FPU)**, which may be integrated into the main CPU or exist as a separate co-processor.

- Addition and Subtraction:
  These are the most complex floating-point operations.
  1. **Extract Components:** The sign, exponent, and mantissa are extracted from both operands.
  2. **Handle Special Cases:** Check for operands being zero, infinity, or NaN. If any are present, special rules apply (e.g., $X + \infty = \infty$).
  3. **Align Exponents:** For addition/subtraction, the exponents *must be the same*. The mantissa of the number with the *smaller* exponent is shifted right until its exponent matches the larger exponent. Each right shift of the mantissa

effectively divides the number by 2, and incrementing the exponent multiplies it by 2, maintaining the number's value. This process ensures the binary points are aligned before addition/subtraction.

- *Example:* Adding (1.011_2times25) and (1.101_2times23). The second number has a smaller exponent. To match the exponent of 5, we shift its mantissa right by 5−3=2 positions: 1.101_2times23=0.01101_2times25.

4. **Add/Subtract Mantissas:** Once exponents are aligned, the mantissas are added or subtracted as if they were integers (using an integer adder/subtractor). The sign of the result is determined.

5. **Normalize Result:** The result of the mantissa operation might not be normalized (e.g., it might be 0.xxxx_2 if it underflowed, or 10.xxxx_2 if it overflowed during addition). The mantissa is then shifted left or right, and the exponent is adjusted accordingly, until the mantissa is in the 1.xxxx_2 normalized form.

6. **Round Result:** After normalization, the result's mantissa may have more bits than the target format (e.g., 23 bits for single-precision). The mantissa must be rounded to fit the available precision according to the chosen rounding mode.

7. **Check for Over/Underflow:** After rounding and final normalization, the exponent is checked to ensure it falls within the representable range. If it's too large, the result becomes pminfty. If it's too small, it might become a denormalized number or pm0.0.

- Multiplication and Division:

These operations are generally simpler than addition/subtraction because exponent alignment is not required in the same way.

1. **Extract Components:** Separate sign, exponent, and mantissa.

2. **Handle Special Cases:** Check for zeros, infinities, NaNs.

3. **Multiply/Divide Signs:** The sign of the result is determined by XORing the sign bits of the two operands. (Same signs rightarrow positive (0); Different signs rightarrow negative (1)).

4. **Add/Subtract Exponents:**

- For **Multiplication:** The true exponents are added. To account for the bias, the formula is usually: Result_Exponent_Biased = (Exp1_Biased + Exp2_Biased) - Bias.

- For **Division:** The true exponents are subtracted. The formula is usually: Result_Exponent_Biased = (Exp1_Biased - Exp2_Biased) + Bias.

5. **Multiply/Divide Mantissas:** The mantissas are multiplied or divided as if they were unsigned integers. This typically produces a mantissa result with double the precision of the input mantissas (e.g., 24-bit * 24-bit multiplication yields a 48-bit product).

6. **Normalize Result:** The resulting mantissa is normalized (shifted and exponent adjusted).

7. **Round Result:** The normalized mantissa is rounded to the target format's precision.

8. **Check for Over/Underflow:** Verify that the final exponent is within the valid range, otherwise set the result to pminfty, pm0.0, or a denormalized number.

**Rounding Modes**

The IEEE 754 standard specifies four primary rounding modes to manage the precision limitation when an exact result cannot be represented:

- **Round to Nearest Even (RoundTiesToEven):** This is the **default and most commonly used rounding mode**. It rounds the result to the nearest representable floating-point number. If the exact result falls precisely halfway between two representable numbers, it rounds to the one whose least significant bit (LSB) of the mantissa is 0 (i.e., the "even" one). This strategy helps to prevent a cumulative bias in a long sequence of operations (e.g., consistently rounding up) by ensuring that roughly half the time, halfway cases round down, and half the time they round up.
- **Round to Zero (Chop/Truncate):** This mode rounds the result towards zero. This means simply discarding (truncating) any bits beyond the specified precision. For positive numbers, it effectively rounds down; for negative numbers, it effectively rounds up towards zero. This is often the fastest rounding mode but introduces a consistent bias towards zero.
- **Round to Plus Infinity (RoundUp):** This mode rounds the result towards positive infinity. For any unrounded result, it rounds to the smallest representable floating-point number that is greater than or equal to the unrounded value.
- **Round to Minus Infinity (RoundDown):** This mode rounds the result towards negative infinity. For any unrounded result, it rounds to the largest representable floating-point number that is less than or equal to the unrounded value.

**Impact of Floating Point Arithmetic on Numerical Accuracy and Precision**

While indispensable, floating-point arithmetic introduces inherent limitations that must be understood to avoid common pitfalls in numerical computation:

- **Finite Precision:** Floating-point numbers represent a continuous range of real numbers using a finite number of bits. This means that only a discrete subset of real numbers can be represented exactly. Most real numbers, especially irrational numbers (like pi or sqrt2) or even simple decimal fractions that do not have a finite binary representation (like 0.1), cannot be stored precisely. They are instead approximated by the closest representable floating-point number.
- **Rounding Errors:** Due to this finite precision, almost every arithmetic operation on floating-point numbers involves some degree of rounding. These small rounding errors, though tiny individually, can accumulate over a long sequence of computations. This accumulation can lead to a significant loss of accuracy in the final result, especially in iterative algorithms or when many operations are performed.
- **Loss of Significance (Catastrophic Cancellation):** A particularly problematic form of rounding error occurs when two floating-point numbers of nearly equal magnitude are subtracted. The most significant bits, which are identical, cancel each other out, leaving a result with far fewer significant digits. The remaining bits (the less significant ones) may then largely consist of accumulated rounding errors from prior operations, leading to a drastically reduced effective precision and a highly inaccurate result. For example, $(1.0000001 \times 10^5) - (1.0000000 \times 10^5)$ should be $(0.0000001 \times 10^5)$, but the subtraction loses most of the precise bits.

- **Non-Associativity of Addition/Multiplication:** Unlike true real number arithmetic, floating-point arithmetic is not always strictly associative. This means that (A+B)+C might not yield precisely the same result as A+(B+C) due to intermediate rounding. The order of operations can influence the final accuracy.
- **Limited Exact Integer Representation:** While floating-point numbers can represent integers, they can only do so exactly up to a certain magnitude (e.g., up to 224 for single-precision, or 253 for double-precision). Beyond this range, integers also become subject to rounding when stored as floating-point numbers, as the gaps between representable floating-point numbers become larger than 1.
- **Special Values and Their Behavior:** The existence of pminfty and NaN means that mathematical operations can produce non-numerical results. This necessitates careful handling in software to prevent these special values from propagating unexpectedly and invalidating further computations.

Understanding these implications is paramount for anyone working with numerical data in computing. Programmers and engineers must be aware of the potential for numerical instability and use techniques like numerical analysis, higher-precision data types (if available), or specifically designed algorithms to mitigate the effects of limited precision and rounding errors in critical applications.