

## K. J. Somaiya College of Engineering, Mumbai-77





Batch: C-2 Roll No.: 16010122267

Experiment / assignment / tutorial No. 5

#### TITLE: Implementation of IEEE-754 floating point representation

**AIM:** To demonstrate the single and double precision formats to represent floating point numbers.

**Expected OUTCOME of Experiment: (Mention CO attained here)** 

#### **Books/ Journals/ Websites referred:**

- **1.** Carl Hamacher, Zvonko Vranesic and Safwat Zaky, "Computer Organization", Fifth Edition, TataMcGraw-Hill.
- **2.** William Stallings, "Computer Organization and Architecture: Designing for Performance", Eighth Edition, Pearson.

#### **Pre Lab/ Prior Concepts:**

The IEEE Standard for Floating-Point Arithmetic (IEEE 754) is a technical standard for floating-point computation established in 1985 by the Institute of Electrical and Electronics Engineers (IEEE). The standard addressed many problems found in the diverse floating point implementations that made them difficult to use reliably and portably. Many hardware floating point units now use the IEEE 754 standard.

#### The standard defines:

- arithmetic formats: sets of binary and decimal floating-point data, which consist of finite numbers (including signed zeros and subnormal numbers), infinities, and special "not a number" values (NaNs)
- *interchange formats:* encodings (bit strings) that may be used to exchange floating-point data in an efficient and compact form
- rounding rules: properties to be satisfied when rounding numbers during arithmetic and conversions
- *operations:* arithmetic and other operations (such as trigonometric functions) on arithmetic formats

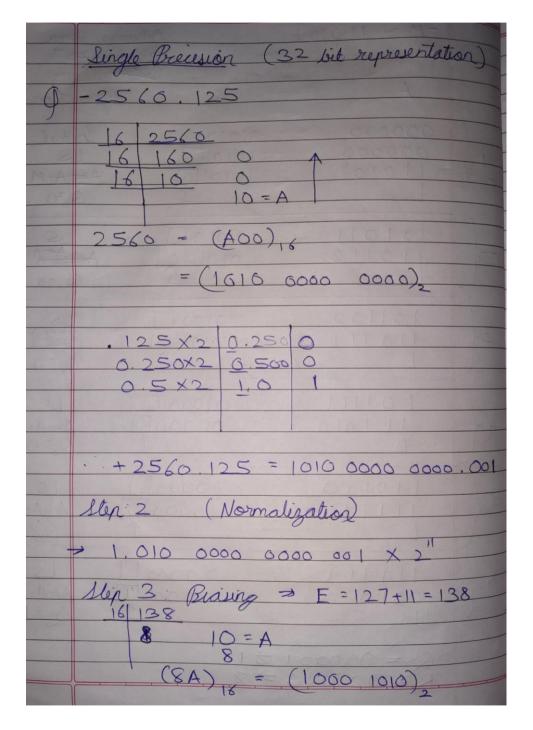






• *exception handling:* indications of exceptional conditions (such as division by zero, overflow, *etc* 

### **Example (Single Precision- 32 bit representation )**









	SP	(32 bits)	
/	Sign Rit	Biased Exponent	Mantissa/Significand
-	1 bit	8 bit	23 bits

## **Example (Double Precision- 64 bit representation )**

Eg.	Double Beusian (64 bit representation)
φ. 100000	102516.0625 16 102516 16 6407 4 16 400 7 16 25 0
	$(9.75)_{16} = 0 + 1 + 6100$ $= (1001 0000 0111 0100)_{2}$ $0.0625 \times 2     0.125     0$ $0.125 \times 2     0.250     0$ $0.25 \times 2     0.5     0$ $0.5 \times 1     1.0     1$
	102516,0625= 1001 0000 0111 0160,0001 step 2 (Abrimalization) 1,001 0000 0111 0100,0001 × 215 Step 3: (Biasing) = E=1023+15=1038





## **Department of Computer Engineering**

	Page
	16 10 38 16 64 14 ≠ E
3-10	
(36)	16 = (40E) = (0000000 1110)2
	Double Breusian (64 bits) Sign Bit Biased Exponent Mantiscal Significant 0 10000001110 001 0000 0111 010000
	1 bit 11 bits 52 bits

#### **Post Lab Descriptive Questions**

Give the importance of IEEE-754 representation for floating point numbers? Ans:

# IEEE-754 representation for floating-point numbers is critically important in the world of computing and numerical analysis for several reasons:

- 1. **Standardization:** IEEE-754 is an industry-standard format for representing floating-point numbers. This standardization ensures that floating-point numbers are consistently represented and interpreted across different hardware platforms and programming languages. It promotes compatibility and portability of numerical computations.
- 2. **Precision and Accuracy:** IEEE-754 defines different formats (e.g., single precision and double precision) that allow for varying levels of precision and accuracy when representing real numbers. This flexibility enables programmers to choose an appropriate format based on their specific needs, balancing computational efficiency and precision.
- 3. **Consistency:** IEEE-754 provides clear rules for performing arithmetic operations (addition, subtraction, multiplication, division) on floating-point numbers. These rules ensure consistent results across different implementations, minimizing



# Somanya TRUST

## **Department of Computer Engineering**

errors due to rounding and other issues.

- 4. **Error Analysis:** The IEEE-754 standard defines mechanisms for error analysis, including the calculation of the relative error and the handling of special values like NaN (Not-a-Number) and infinity. These features are crucial for identifying and managing errors in numerical computations.
- 5. **Compatibility with Hardware:** Most modern CPUs and GPUs support IEEE-754 floating-point arithmetic directly in hardware. This hardware support enables efficient and high-performance floating-point computations, making it essential for scientific and engineering applications.
- 6. **Numerical Stability:** The IEEE-754 standard provides guidelines for ensuring numerical stability in numerical algorithms. It helps programmers avoid common pitfalls that can lead to numerical instability, such as loss of precision in iterative calculations.
- 7. **Interoperability:** IEEE-754 facilitates the exchange of numerical data between different software and hardware components. This interoperability is crucial in a wide range of applications, including scientific computing, graphics, simulations, and data analysis.
- 8. **Portability:** With IEEE-754, numerical algorithms and data can be easily shared and transferred between different platforms and programming languages, reducing the need for custom conversions and adaptations.
- 9. **Education and Research:** IEEE-754 is widely taught in computer science and engineering courses, ensuring that students and researchers have a common foundation for understanding and working with floating-point numbers.
- 10. **Quality Assurance:** The standard helps software developers and quality assurance teams test and verify the correctness of numerical algorithms, making it easier to identify and fix issues related to floating-point precision and rounding. In summary, IEEE-754 representation for floating-point numbers is crucial for ensuring consistency, precision, and interoperability in numerical computations across diverse computing environments. It plays a vital role in scientific, engineering, and computational fields where accurate representation and manipulation of real numbers are essential.

#### Conclusion

Through this experiment, we learnt to convert floating point numbers to their respective IEEE-754 floating point representations (Single Precision and Double Precision).

Date:	26/9/2023	