

**Batch: E-2      Roll No.: 16010123325**

**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)



## Code-

```
#include <bits/stdc++.h>
#include <bitset>
using namespace std;
#define uint unsigned int
#define ull unsigned long long

typedef union {
    float f;
    struct {
        uint mantissa: 23;
        uint exponent: 8;
        uint sign: 1;
    } raw;
} myfloat32;

typedef union {
    double f;
    struct {
        ull mantissa: 52;
        uint exponent: 11;
        uint sign: 1;
    } raw;
} myfloat64;

void printBinary(int n) {
    bitset<32> b(n);
    cout << b.to_string() << '\n';
}

void printBinary64(ull n) {
    bitset<64> b(n);
    cout << b.to_string() << '\n';
}

void printIEEE32(myfloat32 var) {
    cout << "Sign: " << var.raw.sign << '\n';
    cout << "Exponent: ";
    printBinary(var.raw.exponent);
    cout << "Mantissa: ";
    printBinary(var.raw.mantissa);
    cout << '\n';
}
```

```
void printIEEE64(myfloat64 var) {
    cout << "Sign: " << var.raw.sign << '\n';
    cout << "Exponent: ";
    printBinary(var.raw.exponent);
    cout << "Mantissa: ";
    printBinary64(var.raw.mantissa);
    cout << '\n';
}

int main() {
    int choice;
    cout << "Enter 1 for 32-bit IEEE 754 or 2 for 64-bit IEEE 754: ";
    cin >> choice;
    if (choice == 1) {
        myfloat32 var;
        cout << "Enter a decimal number: ";
        cin >> var.f;
        printIEEE32(var);
    } else if (choice == 2) {
        myfloat64 var;
        cout << "Enter a decimal number: ";
        cin >> var.f;
        printIEEE64(var);
    } else {
        cout << "Invalid choice.\n";
    }

    return 0;
}
```

```
PS C:\Users\Shrey\OneDrive\Desktop\KJSCE\SEM-3\COA\Programs> cd "c:\Users\Shrey\OneDrive\Desktop\KJSCE\SEM-3\COA\Programs\" ; if ($?) { g++ IEEE-floating-point.cpp -o IEEE-floating-point } ; if ($?) { .\IEEE-floating-point }
Enter 1 for 32-bit IEEE 754 or 2 for 64-bit IEEE 754: 2
Enter a decimal number: 0.0635
Sign: 0
Exponent: 0000000000000000000000001001101101
Mantissa: 0000000000000000001000011000100110110100111100011010101000
PS C:\Users\Shrey\OneDrive\Desktop\KJSCE\SEM-3\COA\Programs>
```



**Example (Single Precision- 32 bit representation )/ (Double Precision- 64 bit representation )**

DATE:         

Q -13.54

2	13		
2	6	1	↑
2	3	0	
2	1	1	
2	0	1	

$0.54 \times 2 = 1.08$     1

$0.08 \times 2 = 0.16$     0

$0.16 \times 2 = 0.32$     0

$0.32 \times 2 = 0.64$     0

$0.64 \times 2 = 1.28$     1

$0.28 \times 2 = 0.56$     0

$0.56 \times 2 = 1.12$     1

$0.12 \times 2 = 0.24$     0

$0.24 \times 2 = 0.48$     0

$0.48 \times 2 = 0.96$     0

$0.96 \times 2 = 1.92$     1

$0.92 \times 2 = 1.84$     1

$= (-13.54)_{10} = (1101.100010100011)_2$

$= (1.N) = (1.101100010100011) \times 2^3$

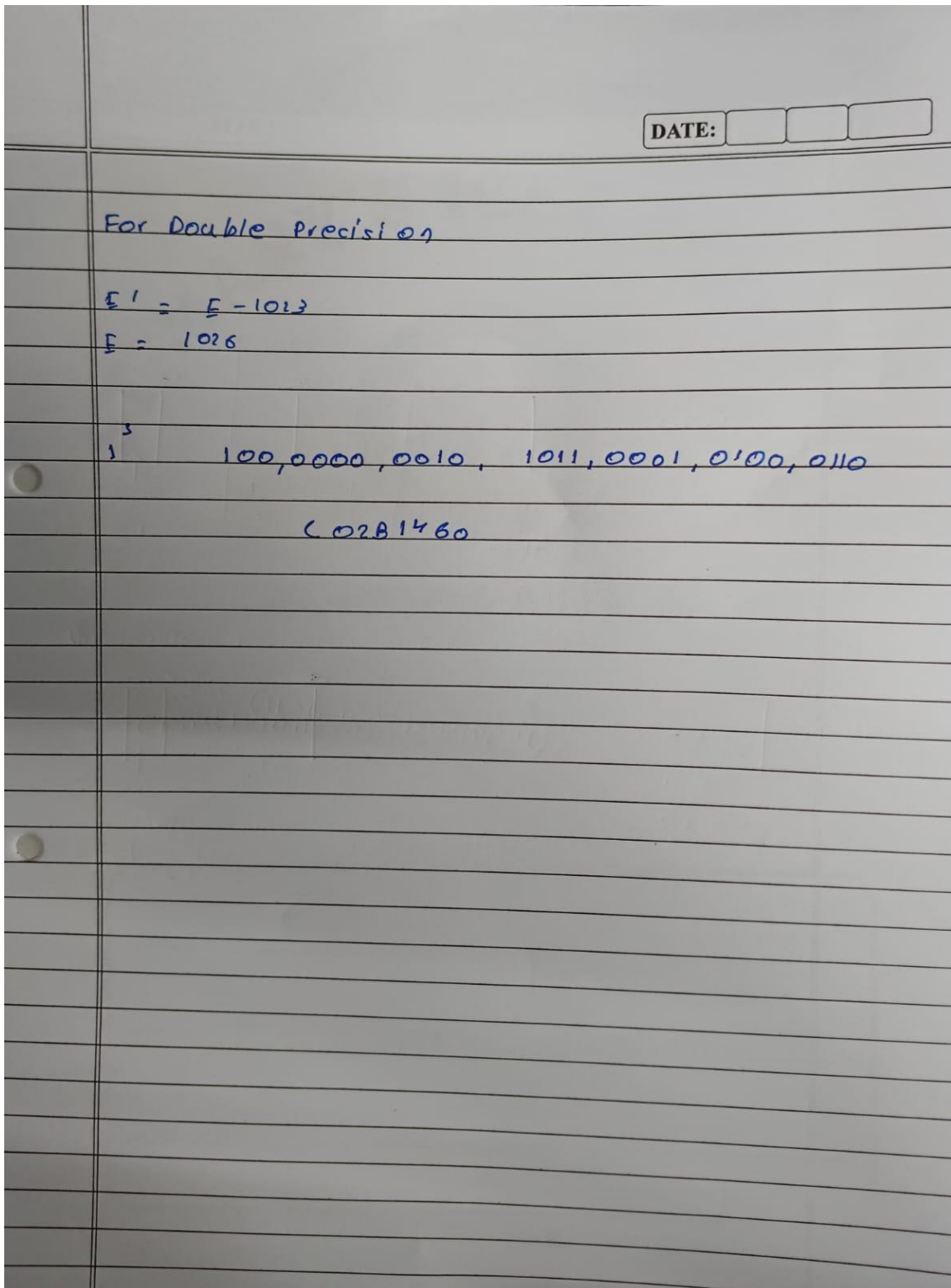
→ For single precision

$E' = E - 127$

$E = 130$

1<sup>1</sup> 100,000 1,0 101,1000,1010,0011

C 158A380



### **Post Lab Descriptive Questions**

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

The **IEEE-754** standard is crucial for representing floating-point numbers consistently across various computing systems. It defines how real numbers are stored using a combination of sign, exponent, and mantissa, allowing for both very large and very small numbers to be represented with a degree of precision. This standard also defines special cases like **NaN** (Not a Number), **Infinity**, and **subnormal numbers**, which help handle arithmetic errors, overflow, and underflow gracefully. The use of IEEE-754 ensures portability, meaning floating-point calculations give the same results on different hardware and software platforms.

### **Conclusion**

Hence, the above experiment highlights the importance of IEEE-754 representation of floating point numbers and their implementation in code using Java.

**Date:** \_\_\_\_\_