



Batch: A3 Roll No.: 16010121045

Experiment / assignment / tutorial No. 5

Grade: AA / AB / BB / BC / CC / CD /DD

Signature of the Staff In-charge with date

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) Example (Double Precision- 64 bit representation)

		139						
12.	23	5						
					I I I			
ste	P 1	; (0	nuert dec i	to bin				
2	2	-	12 = 🕸 1100		0.25			
2	6	0 1			× 2_			
2	3	0	and rates	Constraints and	0.50			
2	1	1	CH	E pringel	×2			
•	1	11		tid	+ 1000 stop			
			1202	5				
			1100 00	the stand of the	they but here			
		1		-				
ste	P 2	.: No	maliration	1.000000	Mar Carl			
		1.	10001 × 23	(scienti	bic notation)			
		Mart.		a glass of	La I			
st	42	3: B1	asing					
single Breasion		pouble	pouble Prefision					
1	E-127			Ē-	E-1029			
3= E-127			27	3	3= E -1029			
	E = 130				E=1026			





/							
130 2 65 0 2 32 1 2 32 1 2 32 1 2 32 1 2 0 0 2 32 1 0 0 0 2 2 2 2 2 0 1 1 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0	ing (32 bit)	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$					
single her		ri lugintiant					
sign bit	Brased exponent	Martisa / sugrigicant					
0	10000010	100 01					
1 bit	8 bit	23 bits					
bigg tot Double precision (64 bit)							
0	1000000000	10001					
1 bit	11 ubits	52 bits.					
		and the second second					





Implementation:

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
int bi[11], f[23], sign[1], expo[8], frac[23];
int expo1[11], fract[52];
int m = 0, fl = 0, i;
void binary(int n)
{
    while (n > 0)
    {
        bi[m] = n \% 2;
        n = n / 2;
        m++;
    }
void floating(float x)
{
    for (i = 0; i < 23; i++)
    {
        x = x * 2;
        f[i] = (int)x;
        x = x - f[i];
    }
void precision(int num)
{
    int e, ee, ee1, k = 0, j = 0, l, r = 0;
    while (m != 0)
    {
        if (bi[m] == 1)
```





```
{
            e = m;
            ee = m + 127;
            ee1 = m + 1023;
            printf("\nSingle precision:\nBiased
exponent:%d\n", ee);
            printf("\nDouble precision:\nBiased
exponent:%d\n", ee1);
            while (ee1 > 0)
            {
                expo1[r] = ee1 \% 2;
                ee1 = ee1 / 2;
                 r++;
            }
            printf("\n");
            printf("%d.", bi[m]);
            m——;
            for (i = m; i \ge 0; i--)
            {
                 frac[k] = bi[i];
                 fract[k] = bi[i];
                printf("%d", frac[k]);
                k++;
            }
            for (i = 0; i < 10; i++)
            {
                 frac[k] = f[i];
                fract[k] = f[i];
                printf("%d", frac[k]);
                k++;
            }
            printf(" x 2^%d", e);
            printf("\n");
            if (num > 0)
```





```
sign[0] = 0;
            else
                sign[0] = 1;
            while (ee > 0)
            {
                expo[i] = ee \% 2;
                ee = ee / 2;
                j++;
            }
            printf("\nSingle bit precision:\n");
            printf("\nSign bit Exponent\t \t \t
Mantissa\n");
            printf("%d", sign[0]);
            printf("\t\t\t");
            for (i = j; i \ge 0; i--)
                printf("%d", expo[i]);
            printf("\t\t\t");
            for (i = 0; i < 23; i++)
                printf("%d", frac[i]);
            printf("\n");
            printf("\nDouble bit precision:\n");
            printf("\nSign bit
                                    Exponent\t \t \t
Mantissa\n");
            printf("%d", sign[0]);
            printf("\t\t\t");
            for (i = r; i \ge 0; i--)
                printf("%d", expo1[i]);
            printf("\t\t\t");
            for (i = 0; i < 52; i++)
                printf("%d", fract[i]);
            break;
        }
        else
            m--:
```





	}
}	
int	main(void)
{	
	float num, x;
	int n;
	<pre>printf("Enter the no.: ");</pre>
	scanf("%f", #);
	n = (int)fabs(num);
	x = fabs(num) - n;
	binary(n);
	<pre>floating(x);</pre>
	<pre>printf("\nIEEE Representation:\n");</pre>
	<pre>precision(num);</pre>
	return 0;
}	

Output:

pargat@Router Programs Enter the no.: 12.25	% cd "/Users/pargat/Docu	ments/COL	LEGE/COA/Programs/"	&& gcc	ieee.c -o	ieee &&	"/Users/pa
IEEE Representation:							
Single precision: Biased exponent:130							
Double precision: Biased exponent:1026							
1.100010000000 × 2^3							
Single bit precision:							
Sign bit Exponent 0	110000010	Mantissa	10001000000000000000	0000			
Double bit precision:							
Sign bit Exponent Ø pargat@Router Programs	110000000010 % <mark>-</mark>	Mantissa	1000100000000000000	0000000	0000000000	00000000	00000008 <mark>8</mark>





Post Lab Descriptive Questions

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

- The IEEE Standard for Floating-Point Arithmetic (IEEE 754) is a technical standard for floating-point computation which was 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 reduced their portability. IEEE Standard 754 floating point is the most common representation today for real numbers on computers, including Intel-based PC's, Macs, and most Unix platforms.
- There are several ways to represent floating point number but IEEE 754 is the most efficient in most cases.

Conclusion : The code for single and double precision formats to represent floating point numbers was executed successfully.

Date:

Signature of faculty in-charge