

1. (+10) Write a C++ program that finds a practical measure of your machine's SP (32 bit) and DP (64 bit) floating point precision by taking the difference of 2 numbers and comparing these to zero in the same precision. What value do you obtain for $\epsilon_{\text{machine}}$ in both precisions? Hint: A loop (j) will be helpful here $1 - (1 + \frac{1}{2^j})$ should do it.

```
int question1() {
    // Float (32 bits)
    std::cout << "##### Single (32 bits) Precision: #####" << std::endl;
    float epsilon_float = 0.0f;
    for (int j = 0; j <= 100; j++) {
        float a = 1.0f;
        float b = 1.0f + 1.0f / pow(2.0f, float(j));
        float diff = a - b;
        if (diff == epsilon_float) {
            epsilon_float = 1.0f / pow(2.0f, float(j-1));
            std::cout << "Estimated single-precision (32 bit) machine epsilon: " << epsilon_float << std::endl;
            break;
        }
    }
    // Double (64 bits)
    std::cout << std::endl;
    std::cout << "##### Double (64 bits) Precision: #####" << std::endl;
    double epsilon_double = 0.0;
    for (int j = 0; j <= 100; j++) {
        double a = 1.0;
        double b = 1.0 + 1.0 / pow(2.0, double(j));
        double diff = a - b;
        if (diff == epsilon_double) {
            epsilon_float = (double)1.0 / pow(2.0, double(j-1));
            std::cout << "Estimated double-precision (64 bit) machine epsilon: " << epsilon_float << std::endl;
            break;
        }
    }
    return 0;
};
```

```
(base) chris@Chris-MacBook-Pro-4 HW % cd "/Users/chris/Library/Clonor/AMATH 583/HW/"HW2
##### Single (32 bits) Precision: #####
Estimated single-precision (32 bit) machine epsilon: 1.19209e-07

##### Double (64 bits) Precision: #####
Estimated double-precision (64 bit) machine epsilon: 2.22045e-16
```

2. (+12) What are the largest and smallest SP (32 bit) and DP (64 bit) numbers that can be represented in IEEE floating point representation? Show work in terms of sign, mantissa, and exponent.

SP

S Sign bit : 1 bit

E Exponent : 8 bits

f fraction : 23 bits

$$(-1)^s \times (1.f) \times 2^{(E-127)}$$

$$1.f = 1 + \frac{\text{fraction bits}}{2^{23}} \quad 1 \leq E \leq 254$$

largest SP: $s = 0$, $2^8 = 256$ values, 0 to 255
 $\Rightarrow (-1)^0 = 1$ \hookrightarrow 255 being represented as ∞ , so $E \leq 254$
 0 being represented as $-\infty$, so $E \geq 1$
 $E = 254 \Rightarrow 2^{(254-127)} = 2^{127}$

Largest fraction: 1.111... - - -

$$0.111... = (1 - 2^{-23})$$

$$1.f = 1 + 1 - 2^{-23} = 2 - 2^{-23}$$

$$(-1)^0 \times (2 - 2^{-23}) \times 2^{127} \approx \underline{3.4028 \times 10^{38}}$$

Smallest SP

$$s = 1 \Rightarrow (-1)^1 = -1$$

$$E = 1$$

smallest fraction 1.000... - - - = 1.0

$$(-1)^1 \times 1.0 \times 2^{-126} \approx \underline{-1.1755 \times 10^{-38}}$$

DP

Sign Bit: 1 bit

Exponent: 11 bits

Fraction: 52 bits

$$1 \leq E \leq 2046$$

$$(-1)^S \times (1.f) \times 2^{(E-1023)}$$

Largest DP:

$$S=0 \Rightarrow (-1)^0 = 1$$

$$E=2046 \Rightarrow 2^{2046-1023} = 2^{1023}$$

$$f = 1.111 \dots = (2 - 2^{-52})$$

$$(-1)^0 \times (2 - 2^{-52}) \times 2^{1023} \approx 1.798 \times 10^{308}$$

Smallest DP:

$$S=0 \Rightarrow (-1)^0 = 1$$

$$E=1 \Rightarrow 2^{1-1023} = 2^{-1022}$$

$$f = 1.0$$

$$(-1)^0 \times 1.0 \times 2^{-1022} \approx 2.2251 \times 10^{-308}$$

3. (+5) Write a C++ program to multiply the integers $200 \times 300 \times 400 \times 500$ on your computer? What is the result? Name the effect you observe.

```
int question3() {
    // Using int type
    int result_int = (int) 200 * (int) 300 * (int) 400 * (int) 500;
    std::cout << "Question 3 (int type): " << result_int << std::endl;
    // Using long type
    long result_long = (long) 200 * (long) 300 * (long) 400 * (long) 500;
    std::cout << "Question 3 (long type): " << (long) result_long << std::endl;
    // Using long long type
    long long result_long_long = (long long) 200 * (long long) 300 * (long long) 400 * (long long) 500;
    std::cout << "Question 3 (long long type): " << (long long) result_long_long << std::endl;
    return 0;
}
```

(base) chris@Chriss-MacBook-Pro-4 HW % cd "/Users/chris/Library/CloudStorage/OneDrive-UW/4. Seignor/AMATH 583 nor/AMATH 583/HW/"HW2

HW2.cpp:39:56: warning: overflow in expression; result is -884901888 with type 'int' [-Winteger-overflow]
 39 | int result_int = (int) 200 * (int) 300 * (int) 400 * (int) 500;

1 warning generated.

Question 3 (int type): -884901888

Question 3 (long type): 12000000000

Question 3 (long long type): 120000000000

Integer multiplication effect = overflow

4. (+5) Given C++ code segment below, what is the final value of *counter*?

```
unsigned int counter = 0;  
for (int i = 0; i < 3; ++i) --counter;
```

```
int question4() {  
    unsigned int counter = 0;  
    for (int i = 0; i < 3; i++) --counter;  
    std::cout << "Final Counter Value: " << counter << std::endl;  
    return 0;  
}
```

clang++. Error: linker command failed with

- (base) chris@Chriss-MacBook-Pro-4 HW % cd nor/AMATH 583/HW/"HW2
Final Counter Value: 4294967293
- (base) chris@Chriss-MacBook-Pro-4 HW %

5. (+10) Count and report how many IEEE SP (32 bit) normalized and denormalized floating point numbers there are. Please count and label infinities and NaNs as well. Show work.

Normalized numbers

sign bit = 1 bit (0 or 1)

Exponent: 1 to 254 $\Rightarrow 254$ values

Fraction: 23 bits $\Rightarrow 2^{23}$ values

$$\text{Total \# of normalized numbers} = 2 \times 254 \times 2^{23} = \underline{4,261,412,864}$$

Denormalized numbers

sign bit = 1 bit

Exponent: 0

Fraction: 23 bits $\Rightarrow 2^{23} - 1$ (minus 0 value)

$$\text{Total \# of denormalized numbers} = 2 \times (2^{23} - 1) = \underline{16,777,214 \text{ (denormalized)}}$$

$$\text{Zeros} = +0 \text{ or } -0 \Rightarrow \underline{2 \text{ (zeros)}} \quad \text{Infinity} = +\infty \text{ or } -\infty \Rightarrow \underline{2 \text{ (Infinities)}}$$

$$\text{NaN} = \left. \begin{array}{l} \text{sign bit} : 1 \text{ bit} \\ \text{Exponent} : 255 \\ \text{fraction} : 23 \text{ bits} \Rightarrow 2^{23} - 1 \end{array} \right\} 2 \times (2^{23} - 1) = \underline{16,777,214 \text{ (NaN)}}$$

6. (+15) Consider a 6 bit floating point system with $s = 1$ (1 sign bit), $k = 3$ (3 bit exponent field), and $n = 2$ (2 bit mantissa).

- Calculate by hand all the representable normalized numbers. Show work.
- Calculate by hand all the representable denormalized numbers. Show work.
- Plot both sets of numbers (ignoring NaNs and infinities) as a number line to see the gaps of (un)representable numbers.

$$\text{Value} = (-1)^s \times \left(1 + \frac{\text{fraction bits}}{2^n}\right) \times 2^{(e - \text{bias})}$$

$$s = 0/1, n = 2, \text{bias} = 2^{k-1} - 1 = 2^2 - 1 = 3$$

$$E = e - \text{bias} = \{1, 2, 3, 4, 5, 6\} - 3 = \{-2, -1, 0, 1, 2, 3\}$$

$$M = 1 + \frac{fb}{4}$$

$$\text{Value} = (-1)^s \times M \times 2^E \quad fb = \begin{matrix} 0 & 1 & 2 & 3 \\ 00 & 01 & 10 & 11 \end{matrix}$$

$$M = \{1.0, 1.25, 1.50, 1.75\}$$

a)

$$s = 1 \text{ (Positive)}$$

$$E = -2, \text{Value} = \begin{aligned} (1.0) \times 2^{-2} &= 0.25 \\ (1.25) \times 2^{-2} &= 0.3125 \\ (1.50) \times 2^{-2} &= 0.375 \\ (1.75) \times 2^{-2} &= 0.4375 \end{aligned}$$

$$E = 2, \text{Value} = \begin{aligned} (1.0) \times 2^2 &= 4.0 \\ (1.25) \times 2^2 &= 5.0 \\ (1.50) \times 2^2 &= 6.0 \\ (1.75) \times 2^2 &= 7.0 \end{aligned}$$

$$E = -1, \text{Value} = \begin{aligned} (1.0) \times 2^{-1} &= 0.5 \\ (1.25) \times 2^{-1} &= 0.625 \\ (1.50) \times 2^{-1} &= 0.75 \\ (1.75) \times 2^{-1} &= 0.875 \end{aligned}$$

$$E = 3, \text{Value} = \begin{aligned} (1.0) \times 2^3 &= 8.0 \\ (1.25) \times 2^3 &= 10.0 \\ (1.50) \times 2^3 &= 12.0 \\ (1.75) \times 2^3 &= 14.0 \end{aligned}$$

$$E = 0, \text{Value} = \begin{aligned} (1.0) \times 2^0 &= 1.0 \\ (1.25) \times 2^0 &= 1.25 \\ (1.50) \times 2^0 &= 1.50 \\ (1.75) \times 2^0 &= 1.75 \end{aligned}$$

Same value with negative sign for $s = 0$

$$E = 1, \text{Value} = \begin{aligned} (1.0) \times 2^1 &= 2.0 \\ (1.25) \times 2^1 &= 2.5 \\ (1.50) \times 2^1 &= 3.0 \\ (1.75) \times 2^1 &= 3.5 \end{aligned}$$

h) Denormalized value

$$M = \frac{fb}{4} = \{0.0, 0.25, 0.50, 0.75\}$$

$$E = 1 - \text{bias} = 1 - 3 = -2$$

$$\text{Value} = (-1)^S \times M \times 2^E$$

$$S = 1 \text{ (Positive)}$$

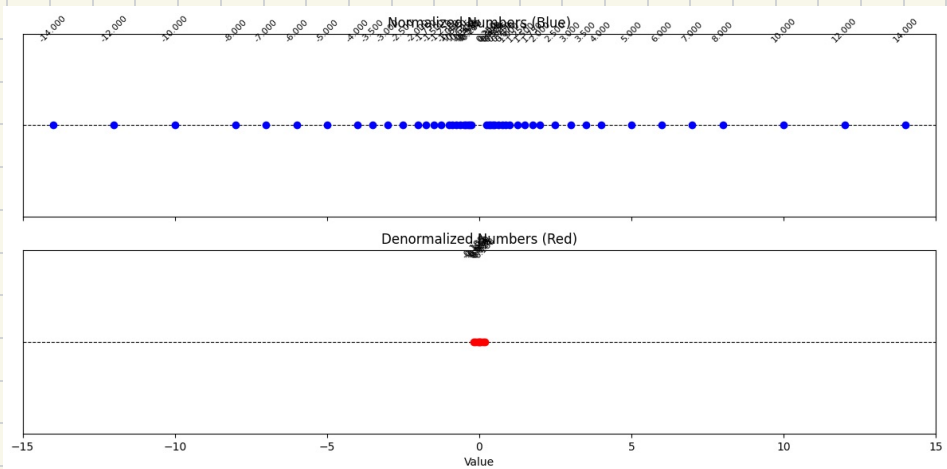
$$S = -1 \text{ (Negative)}$$

$$\begin{aligned}\text{Value} &= 0.0 \times 2^{-2} = 0 \\ 0.25 \times 2^{-2} &= 0.0625 \\ 0.50 \times 2^{-2} &= 0.125 \\ 0.75 \times 2^{-2} &= 0.1875\end{aligned}$$

$$\begin{aligned}\text{Value} &= 0 \\ &= -0.0625 \\ &= -0.125 \\ &= -0.1875\end{aligned}$$

Zero $\Rightarrow \pm 0$,

c)



7. (+10) Conversions. Show work.

(a) Write $(D3B701)_{16}$ as an integer in base-10.

(b) Write $(1010000100111111)_2$ as an integer in base-16, i.e. as a hexadecimal number.

a)

0 1 2 3 4 5 6 7 8 9 A B C D E F

10 11 12 13 14 15

$$13 \cdot 16^5 + 3 \cdot 16^4 + 11 \cdot 16^3 + 7 \cdot 16^2 + 0 \cdot 16^1 + 1 \cdot 16^0$$

$$= 13874945$$

b)

12

$$8 + 4 + 2 + 1 =$$

$$\begin{array}{cccc} 1010 & 0001 & 0011 & 1111 \\ \downarrow & \downarrow & \downarrow & \downarrow \\ 2^3 + 2 = 10 & 1 & 2 + 1 = 3 & 2^3 + 2^2 + 2^1 + 2^0 = 15 \\ \downarrow & \downarrow & \downarrow & \downarrow \\ A & 1 & 3 & F \end{array}$$

$$(A13F)_{16}$$

8. (+5) Are there $a, b, c \in \mathbb{Z}$ s.t. $6a + 9b + 15c = 107$? Show work.

$$\gcd(6, 9, 15) = 3$$

Any linear combination of 6, 9, 15, $6a + 9b + 15c$ $a, b, c \in \mathbb{Z}$

must be a scalar multiple of 3. Since 107 is not divisible

by 3, there does not exist integers a, b, c s.t. $6a + 9b + 15c = 107$.

```
int question8() {
    int a_max = 107 / 6 + 1;
    int b_max = 107 / 9 + 1;
    int c_max = 107 / 15 + 1;
    int target = 107;
    for (int i = 0; i <= a_max; i++) {
        for (int j = 0; j <= b_max; j++) {
            for (int k = 0; k <= c_max; k++) {
                if ((6 * i + 9 * j + 15 * k) == target) {
                    std::cout << i << ", " << j << ", " << k << std::endl;
                    break;
                }
            }
        }
    }
    std::cout << "DNE" << std::endl;
    return 0;
}
```

• (base) chris@Chriss-MacBook-Pro
nor/AMATH 583/HW/"HW2
Question 8: DNE

9. (+10) Equivalence classes modulo n . $\forall a, b \in \mathbb{Z}$ then $a \equiv b \pmod{n}$ means $n \mid (a - b)$ or $a = b + k \cdot n$ and $k \in \mathbb{Z}$. \mathbb{Z}_n is the set of equivalence classes $\{[0], [1], \dots, [n-1]\}$. Is $(\mathbb{Z}_n, +, \cdot)$ a ring? Hint: If $s \in [i]$, then $n \mid (s - i)$. Show work (use ring properties).

$$\text{If } [a] = [a'] \text{ and } [b] = [b']$$

$$\text{show } [a+b] = [a'+b'] \text{ and } [a \times b] = [a' \times b']$$

$$\begin{array}{lcl} \text{Since } [a] = [a'] & \text{means} & a \equiv a' \pmod{n}, \quad n \mid (a - a') \\ [b] = [b'] & & b \equiv b' \pmod{n}, \quad n \mid (b - b') \end{array}$$

$$(a+b) - (a' - b') = (a - a') + (b - b') \quad \text{rs also divisible by } n.$$

Operations well-defined.

Ring Axiom Verification

R1: Closure under $+$ and \times

$$\text{If } [a], [b] \in \mathbb{Z}_n, \text{ then } \left. \begin{array}{l} [a] + [b] = [a+b] \\ [a] \cdot [b] = [ab] \end{array} \right\} \text{ are still equivalence class}$$

R2: Associativity

$$([a] + [b]) + [c] = [(a+b) + c] = [a + (b+c)] = [a] + ([b] + [c]).$$

R3: Commutativity $+$

$$[a] + [b] = [b] + [a] \quad \text{because } a+b = b+a \text{ in } \mathbb{Z}$$

R4: Existence of additive identity R5: Inverse additive

$$[a] + [0] = [a+0] = [a]$$

$$[a] + [-a] = [a+(-a)] = [0]$$

R6: Commutativity of \times

$$[a] \times [b] = [ab] = [ba] = [b] \times [a]$$

R7: Existence of Multiplicative Identity

$$[a] \times [1] = [a \times 1] = [a]$$

R8: Distributivity

$$([a] + [b]) \times [c] = [(a+b) \times c] = [a \times c + b \times c] = [a] \times [c] + [b] \times [c]$$