

High Performance Computing
Homework 3 - Question 1

I ran my program on the COE Vector Linux platform (Intel® Xeon® CPU E5-2698 v4 @ 2.20GHz, 20 cores per socket, 2 sockets, with 791957684 KB of memory, and Linux 4.18.0).

a. Result shown below:

```
Single Precision Floating Point Calculations:
sin(0.05) = 0.049979172646999359130859375
sin(0.5) = 0.479425519704818725585937500
sin(1) = 0.841470956802368164062500000
sin(1.5) = 0.997325062751770019531250000

Double Precision Floating Point Calculations:
sin(0.05) = 0.049979169270678337755331455127816298045217990875244140625
sin(0.5) = 0.479425538604187073676854424775228835642337799072265625000
sin(1) = 0.841470969086639208889266683399910107254981994628906250000
sin(1.5) = 0.997324949172103236705311246623750776052474975585937500000

|Single Precision Error|:
sin(0.05): 3.376e-09
sin(0.5): 1.890e-08
sin(1): 2.801e-08
sin(1.5): 1.699e-04

|Double Precision Error|:
sin(0.05): 6.939e-18
sin(0.5): 1.593e-14
sin(1): 1.572e-08
sin(1.5): 1.700e-04
```

Figure 1. Taylor Series Sin(x) Approximations with Floats and Doubles

The implementation with doubles provides answers with more decimal places specified. This is because double precision floating point values have a longer mantissa. Additionally, the absolute error of results was consistently greater when single precision floats were used. Since not all decimal numbers can be exactly specified using IEEE 754 floating point, there are often rounding errors in their binary representations. These rounding errors are more severe for single precision because of the 23 bit mantissa compared to the 52 bit mantissa of a double precision floating point value. When many computations are made, these errors compound. Therefore, when accuracy is crucial and/or there is a multi step calculation, double precision floating points will produce a more desirable result in most cases.

b. IEEE 754 Floating Point Representations

• 2.1

a. Single Precision

Sign Bit: 0 (+)

Exponent: $1 + 127_bias = 128 = 10000000_2$

Mantissa: $1.05 - 1_implied = 0.05 = 0.000011001100_2$
 $\approx 0.00001100110011001100$

Binary: 01000000000001100110011001100110₂

Hexadecimal: 40066666

b. Double Precision

Sign Bit: 0 (+)

Exponent: $1 + 1023_bias = 1024 = 10000000000_2$

Mantissa: $1.05 - 1_implied = 0.05 = 0.000011001100_2$
 $\approx 0.00001100110011001100110011001100110011001101_2$

Binary:

01000000000000011001100110011001100110011001100110011001100110011001101₂

Hexadecimal: 4000CCCCCCCCCCCD

• 6300

a. Single Precision

Sign Bit: 0 (+)

Exponent: $12 + 127_bias = 139 = 10001011_2$

Mantissa: $1.5380859375 - 1_implied = 0.5380859375$
 $= 0.1000100111000000000000_2$

Binary: 0100010111000100111000000000000₂

Hexadecimal: 45C4E000

b. Double Precision

Sign Bit: 0 (+)

Exponent: $12 + 1023_bias = 1035 = 10000001011_2$

Mantissa: $1.5380859375 - 1_implied = 0.5380859375$
 $= 0.100010011100_2$

Binary:

010000001011100010011100₂

Hexadecimal: 40B89C0000000000

- -1.044

- a. Single Precision

Sign Bit: 1 (-)

Exponent: $0 + 127_{\text{bias}} = 127 = 01111111_2$

Mantissa: $1.044 - 1_{\text{implied}} = 0.044 = 0.000010110100001110010101\dots_2$

$\approx 0.00001011010000111001011$

Binary: 10111111100001011010000111001011

Hexadecimal: BF85A1CB

- b. Double Precision

Sign Bit: 1 (-)

Exponent: $0 + 1023_{\text{bias}} = 1023 = 01111111111_2$

Mantissa: $1.044 - 1_{\text{implied}} = 0.044$

$= 0.0000101101000011100101011000000100000110001001001110\dots_2$

Binary:

101111111110000101101000011100101011000000100000110001001001110

Hexadecimal: BFF0B4395810624E