

Exercise 2

Introduction to Computational Astrophysics, SoSe 2024

Seha Lee

Task 1. Integer representation

Solution. (a) $2^{10} - 1 = 1023$

(b) The largest positive integer has 0 for its most significant digit to indicate it's positive, and then the rest is 1 to be the largest, i.e., $n_{max} = 01111111111111111111111111111111_2$. If you add 1 to this value, the result is 33 digit number, the most significant digit being 1 and the rest being 0, i.e., $n_{max} + 1 = 10000000000000000000000000000000_2$. Utilizing two's complement we can avoid the redundancy of zeros and inefficiency of computations caused by the use of sign bit.

□

Task 2. Floating point number representation

Solution. (a)

The completed code:

```
1 #include <stdio.h>
2
3 int main() {
4     float y, x = 1.043E-13 ;
5     x = x / 10. ;
6     y = x*x - (1.043E-14)*(1.043E-14) ;
7
8     printf("x= %.15f\n", x);
9     printf("y= %.15f\n", y);
10
11     return 0;
12 }
```

It was expected that x and y are 1.043E-14 and 0 respectively. However, computed results are x= 0.000000000000010 and y= -0.000000000000000. There is a clear discrepancy between expected and actual computed result for y, which is a very small negative number close to zero but not exactly zero. This is due to the limitation of computer memory using only finite number of bits to store real numbers, leading to inaccuracy in floating point arithmetic as in this case,

resulting in catastrophic cancellation.

(b)

The bit patterns for the numbers expressed by standard IEEE 754 is as following, the sign bit, exponent, and mantissa respectively:

- 7 is 0 10000001 110000000000000000000000
- $1.E - 7$ is 0 01110111 00001010001111010111000.
- The sum of 7 and $1.E - 7$ retains the bit pattern of 7. This is due to absorption.

□

Task 3. GRS bits

Solution. (a)

$$32 - 2.125 = 29.875$$

(b)

$$\begin{array}{r}
 1.000\,000 \times 2^5 \\
 - 0.000\,101 \times 2^5 \\
 \hline
 0.111\,011 \times 2^5 \\
 1.110\,11 \times 2^4 \\
 1.111 \times 2^4
 \end{array}
 \begin{array}{l}
 \text{GRS 101} \\
 \\
 \text{shifted left} \\
 \text{rounded}
 \end{array}$$

□

Task 4. Radius calculation

Solution. By using the relative value to the Sun, we do not have to introduce constants with extreme values such as σ_{SB} , avoiding potential problems of floating point arithmetic that can lead to inaccurate results.

The code is:

```

1 #include <iostream>
2 #include <cmath>
3
4 int main() {
5     const double tempSun = 5778;
6     const double radiusSun = 695700; // km
7
8     double logLsun;
9     double Teff;
10
11     std::cout << "Enter log(L*/Lsun): ";
12     std::cin >> logLsun;

```

```
13     std::cout << "Enter Teff (in Kelvin): ";
14     std::cin >> Teff;
15
16     double radius = sqrt(pow(10, logLsun) * pow(tempSun / Teff, 4));
17
18     std::cout << "The radius is: " << radius << " R_sun" << std::endl;
19
20     return 0;
21 }
```

□