

# NUMERISCHE METHODEN DER PHYSIK

WiSE 2023-2024 – PROF. MARC WAGNER

MICHAEL EICHBERG: [eichberg@itp.uni-frankfurt.de](mailto:eichberg@itp.uni-frankfurt.de)

LASSE MÜLLER: [lmuller@itp.uni-frankfurt.de](mailto:lmuller@itp.uni-frankfurt.de)

## Exercise sheet 1

*Not to be handed in. To be discussed in the tutorials on 20.10.23 and 23.10.23*

### Exercise 1 [Floating-point numbers]

Let us assume we have a (rather primitive) computer that uses 8-bit floating-point arithmetic. The first bit represents the sign, the next 4 bits the exponent with bias  $b = 7$  and the last 3 bits for the mantissa (normalized representation with leading 1 before the comma). With this we have the representation of a real number  $x$  as

$$x = (-1)^s \left[ 1 + \sum_{n=1}^3 m_n \cdot 2^{-n} \right] \cdot 2^{(\sum_{i=0}^3 e_{3-i} \cdot 2^{3-i}) - b}$$

which can be stored in a bit string  $\mathbf{s e_3 e_2 e_1 e_0 m_1 m_2 m_3}$ . Assume that non-representable numbers are rounded to the nearest representable number (as usually it happens).

- (i) Which number is represented by the bit-string 10111000?
- (ii) Which is the bit-string for the number  $-26$ ? And for the number 0?
- (iii) How many different numbers can be represented exactly in this way? Which is the smallest and which is the largest positive number?
- (iv) What are the numerical results of the differences  $(\frac{35}{32} - \frac{33}{32})$  and  $(\frac{37}{32} - \frac{35}{32})$ ?
- (v) Which number(s) have the largest absolute error? Which have the largest relative error in the interval between the smallest and the largest representable positive number?
- (vi) Repeat (iii) setting  $b = 3$ . Which role does the bias play? What happens, when you vary the bias?
- (vii) How could you determine the smallest positive representable number on your computer? Try to write a simple program which prints the result to the screen using `single` and the `double` precision.
- (viii) Do you think it is a good idea, in a program, to check for equality between two floating-point numbers using the equality operator? When is it safe and when not? What could be an alternative?