

Question 2 (20 points): When building models for deep learning, lower-precision floating-point representation can be used. Thus, there is a growing interest in floating-point representation that use fewer bits. Last year a minifloat representation was proposed. Minifloat borrows the same ideas used in floating-point formats to represent integer values in a large range using only 8 bits. An 8-bit minifloat has 1 sign bit, 4 exponent bits and 3 mantissa bits and an exponent bias equal +2 for normalized numbers.

- Minifloat has representations for +infinity (0 1111 000) and -infinity (1 1111 000).
- Minifloat has a representation for Not-a-Number (NaN): x 1111 yyy where yyy \neq 000.
- When the exponent is 0000 the number represented is subnormal, and the value is $0.\text{mmm} \times 2^3$ where mmm are the three bits of the mantissa.
- When the exponent is not zero, then the number represented is normalized. Given a representation s eeee mmm, the value represented is $1.\text{mmm} \times 2^{e+2}$ where e is the value of the exponent eeee in the representation.
- The standard rules of rounding to the nearest even apply to minifloat.

1. (**5 points**) Which is the smallest, non-zero, positive value that can be represented in minifloat? Provide the binary representation and the decimal value.

2. (**5 points**) Which is the largest positive value that can be represented in minifloat? Provide the binary representation and the decimal value.

3. (**5 points**) There are many numbers within the range between the smallest and the largest value that cannot be represented precisely in minifloat. In such a case the value is rounded to the nearest value. What is the binary representation of 46_{10} in minifloat? Which is the actual value of the number represented after rounding?

4. (**5 points**) A hardware that supports minifloat has an adder with a guard bit, a round bit, and a sticky bit. Show how the operation $5+64$ is performed in this hardware, and what is the value of the result both in binary and in decimal. What is the minifloat representation of the result?