

There were a number of questions about the exaggerated floating-point arithmetic system presented in class, with  $-1 \leq e \leq 3$ ,  $t = 3$  significant digits, and base  $\beta = 2$ . There is an equally brief discussion of this system in the text on p. 194, but the paper by Goldberg (see the link on the course page) should be consulted.

But it may be useful to remind the student of some elementary facts about the positional number system. When we say that we are using 3 digits (a binary digit is called a “bit”), a base  $\beta$ , and an exponent  $e$ , then a number

$$0 . d_1 d_2 d_3 \cdot \beta^e$$

means (by the positional number system, the exponent of the base is related to the number of places shifted past the “decimal” point)

$$\left[ (d_1 \cdot \beta^{-1}) + (d_2 \cdot \beta^{-2}) + (d_3 \cdot \beta^{-3}) \right] \cdot \beta^e$$

and therefore 0.100 with  $e = -1$  and  $\beta = 2$  represents the number

$$\left[ 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 0 \cdot 2^{-3} \right] \cdot 2^{-1} = \frac{1}{4}$$

while 0.101 with  $e = 2$  then represents

$$\left[ 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} \right] \cdot 2^2 = \frac{5}{2}.$$

This is a straightforward extension of the (possibly more commonly encountered in elementary school) use of the positional number system for representing numbers with digits to the left of the “decimal” point.

## 1 An excursion into hexadecimal

If we choose the base  $\beta = 16$ , then we need more symbols for digits than we have in the decimal system. The usual symbols are 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, and F. We use A to mean the number represented as 10 in decimal notation; B for 11, and so on up to F for 15. We therefore count in hex 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F, 10, 11, 12, ..., 19, 1A, 1B, 1C, 1D, 1E, 1F, 20, 21, 22, ..., 29, 2A, 2B, 2C, 2D, 2E, 2F, 30, and so on. This is an increasingly attractive number system for human ages: my mother just turned 50, for example (in hexadecimal), while in 5 years I will be 30.

Similarly, 0.1F2A means  $1/16 + 15/16^2 + 2/16^3 + 10/16^4$ . The idea should be clear now.