

Lecture 4 Activity Solution

Model 1: What is floating point?

Please note that some answers represent one interpretation and there are other valid approaches.

1. 1.5213×10^4
2. One possible representation: 8 digits. 15213104. But the answer may vary, depending how you want to represent the number.
3. 18213104. 8 digits
4. 18213107. 8 digits
5. 10001100, which is 1.0001.
6. 99999999, which is 9.9999×10^9 or 9.9999×99^9 or 9×9^{9999} .
7. No.

Model 2: Binary Scientific Notation

1. 1
2. 1.0111×2^4 , 1.0111×2^2 , 1.0111×2^1 , 1.0111
3. 1

Model 3: IEEE Representation

1. Sign bit. The number is negative.
2. 0111
3. 1
4. With no bias, it would be 2, which is greater than 1.
5. 0b10000000
6. $exp = 13 + 127 = 140$. $15213_{10} = 0b01000110011011011011010000000000$, or 0x466db400
7. From -1022 to 1023

Model 4: Extreme Exponents

1. 1.0000
2. No.
3. Two, one positive, one negative.
4. 0.0001
5. $+inf$. Yes, 2^{127} is 0111111100000 and $+inf$ is 0111111110000.
6. Largest denormalized number has all 0 for exponent bits and all 1 for fraction bits. Smallest normalized number has all 0 except the lowest exponent bit to be one and all 0 for fraction bits.

Model 5: Addition and Multiplication

1. 1.0011×2^4
2. 4
3. 1, 1.0, 1, 2.0
4. 1.00011, 1.00, 1; 1.00101, 1.01, 1.25; 1.111, 10.0, 2; 1.101, 1.10, 1.5
5. 1.0011×2^4 , since the number is normalized. However, if we include the leading 1, rounding to even yields 1.010×2^4 .
6. 2048
7. 2^{11}

Model 6: Simple Floating-point

1. 15.5 (01101111), 1/64 (00000001)
2. 0, this is a cancellation error, it happens because larger floating point numbers are farther apart than smaller ones and there is a trade off between range and precision.
3. 7, 111
4. 0x63
5. +Inf (0x70)

Model 7: Review

1. Yes it will, at 2^{24} . Some large numbers will have precision that cannot be represented exactly in float. The rounding rules for adding 1.0f round down at this value.
2. It won't terminate
3. No. It implies that some ints cannot be casted into an equivalent float value.
4. No. Double uses 53 fractional bits, so all 32-bit ints can have equivalent doubles.