

1 Determine the single precision 32-bit representation of the following decimal numbers:

a) 2^{-30}

$$\begin{aligned} & [2^{-30}]_{10} \\ &= 1 \times 2^{-30} \\ & \quad \text{*bias*} \\ & \quad 1 \times 2^{97} \\ &= [1.0 \times 10^{01100001}] \end{aligned}$$

0 01100001 000000000000000000000000000000 ✓

b) 64.015625

$$64_{10} = 0100\ 0000_2$$

$$\begin{aligned} & \left[\begin{array}{l|l} 0.015625 & \\ 0.03125 & 0 \\ 0.0625 & 0 \\ 0.125 & 0 \\ 0.25 & 0 \\ 0.5 & 0 \\ 1.0 & 1 \end{array} \right] \\ & \rightarrow = 0.000001_2 \end{aligned}$$

$$\begin{aligned} & [64.015625]_{10} \\ &= [01000000.000001]_2 \\ &= 1.000000000001 \times 10^{0110} \\ & \quad \text{*bias*} \\ &= 1.000000000001 \times 10^{10000101} \end{aligned}$$

0 10000101 0000000000000100000000000000 ✓

c) -8×2^{-24}

$$\begin{aligned} & [-8 \times 2^{-24}]_{10} \\ & \quad \text{*bias*} \\ & \quad -8 \times 2^{103} \\ &= [-1000.0 \times 10^{01101010}]_2 \\ &= -1.0 \times 10^{01100110} \end{aligned}$$

1 01101010 000000000000000000000000000000 ✓

d) 0.5

$$\begin{aligned} & [0.5]_{10} \\ &= [0.1]_2 \\ &= 1.0 \times 10^{-1} \\ & \quad \text{*bias*} \\ &= 1.0 \times 10^{01111110} \end{aligned}$$

0 01111110 000000000000000000000000000000 ✓

e) 42.424242

5 sig. bits after 1

$42_{10} = 00101010$

0.424242	
0.848484	0
1.696968	1
1.393936	1
0.787872	0
1.575744	1
1.151488	1
0.302976	0
0.605952	0
1.211904	1
0.423808	0
0.847616	0
1.695232	1
1.390464	1
0.780928	0
1.561856	1
1.123712	1
0.247424	0
0.494848	0

need 18 (23-5)
more bits
for mantissa

$\rightarrow 0.0110111001001101100$

$[42.424242]_{10} = [00101010.011011001001101100]_2$

$= 1.01010011011001001101100 \times 10^{0101}$

bias

$1.01010011011001001101100 \times 10^{10000100}$

$0 \mid 10000100 \mid 01010011011001001101100 \checkmark$

f) $76.234567 \times 10^{15} \rightarrow [0.0000000000000076234567]$

Goal: Rewrite in scientific notation w/ base 2

76.234567×10^{15}
 $= 7.6234567 \times 10^{14}$

Find Exponent

$2^K \leq 7.6234567 \times 10^{14} < 2^{K+1}$
 $\Rightarrow K = \lfloor \log_2(7.6234567 \times 10^{14}) \rfloor$
 $= \lfloor 43.57... \rfloor$
 $= -44$

$\therefore 2^{-44} \leq 7.6234567 \times 10^{14} < 2^{-43}$
 $5.68...e-14 \leq 7.62...e-14 < 1.13...e-13$

$|2^{-44} - 7.62...e-14| < |2^{-43} - 7.62...e-14|$
 $\therefore 7.6234567 \times 10^{14} = C \times 2^{-44}$

Find Coefficient

$C = 2^{44} \times (7.6234567 \times 10^{14})$
 $= 1.341132685679496527872$

Coefficient to Binary

0.34...	
0.68...	0
1.36...	1
0.72...	0
1.45...	1
0.91...	0
1.83...	1
1.66...	1
1.32...	1
0.65...	0
1.31...	1
0.63...	0
1.27...	1
0.55...	0
1.11...	1
0.23...	0
0.47...	0
0.94...	0
1.88...	1
1.77...	1
1.54...	1
1.09...	1
0.18...	0
0.37...	0

bias

$[1.341132685679496527872 \times 2^{83}]_{10}$
 $= 1.01010111010101000111100 \times 10^{01010011}$

$0 \mid 01010011 \mid 01010111010101000111100 \checkmark$

g) 1.4345678

0 sig bits

$$1.0_{10} = 1.0_2$$

need 23 mo'

0.4345678	
0.8691356	0
1.7382712	1
1.4765424	1
0.9530848	0
1.9061696	1
1.8123392	1
1.6246784	1
1.2493568	1
0.4987136	0
0.9974272	0
1.9948544	1
1.9897088	1
1.9794176	1
1.9588352	1
1.9176704	1
1.8353408	1
1.6706816	1
1.3413632	1
0.6827264	0
1.3654528	1
0.7309056	0
1.4618112	1
0.9236224	0

$$\rightarrow = 0.0110111001111101010$$

$$1.0110111001111101010 \times 10^{00000000}$$

bias

$$1.0110111001111101010 \times 10^{01111111}$$

$$0 \quad 01111111 \quad 0110111001111101010 \quad \checkmark$$

↑ 1 bit off. why?

h) 3.141592653589793238462643383279502884

(my calculator only has 15 sig digits)

$$[3]_{10} = [0011]_2$$

0.14...	
0.28...	0
0.56...	0
1.13...	1
0.26...	0
0.53...	0
1.06...	1
0.12...	0
0.24...	0
0.49...	0
0.99...	0
1.98...	1
1.96...	1
1.92...	1
1.85...	1
1.70...	1
1.41...	1
0.83...	0
1.66...	1
1.32...	1
0.65...	0
1.31...	1
0.63...	0

$$[3.14...]_{10}$$

$$= [11.0010010000111101010]_2$$

$$= 1.10010010000111101010 \times 10$$

$$0 \quad 10000000 \quad 10010010000111101010 \quad \checkmark$$

↑ 1 bit off. why?

i) 3.14

$$\begin{bmatrix} 3 \\ 11 \end{bmatrix}_{10} = \begin{bmatrix} 11 \\ 2 \end{bmatrix}_2$$

.14	
.28	0
.56	0
1.12	1
.24	0
.48	0
.96	0
1.92	1
1.84	1
1.68	1
1.36	1
.72	0
1.44	1
.88	0
1.76	1
1.52	1
1.04	1
.08	0
.16	0
.32	0
.64	0
1.28	1
0.56	0
1.12	1

$$\begin{aligned} & [3.14]_{10} \\ & = [11.001000111101011000010]_2 \\ & = 1.1001000111101011000010 \times 10 \end{aligned}$$

0 10000000 1001000111101011000010 ✓

↑ 1 bit off. why?

j) 17/31

17/31	
34/31	1
6/31	0
12/31	0
24/31	0
48/31	1
34/31	1

repeats (5 times fit)

$$\begin{aligned} & [17/31]_{10} \\ & = [0.1000110001100011000110001]_2 \\ & = 1.000110001100011000110001 \times 10^0 \\ & \text{bias} \\ & = 1.000110001100011000110001 \times 10^{0111110} \end{aligned}$$

0 0111110 00011000110001100011000 ✓

↑ 1 bit off. why?

2 Identify the floating point numbers corresponding to the following single precision machine binaries:

a) 0 00000000 000000000000000000000000 |

0 ✓

b) 1 1111111 000000000000000000000000 |

-Infinity ✓

c) 0 10000001 011000000000000000000000 |

$$\begin{aligned}
 &+ \quad 129_{10} \quad 1.011_2 \\
 &\quad -127 \\
 &= [1.011 \times 10^{10}]_2 \\
 &= [2.395 \times 2^2]_{10} \\
 &= 9.25 \quad \checkmark
 \end{aligned}$$

d) 0 0111111 000000000000000000000000 |

$$\begin{aligned}
 &+ \quad 127 \quad 1.0 \\
 &\quad -127 \\
 &= [1.0 \times 10^{00}]_2 \\
 &= [2 \times 2^0]_{10} \\
 &= 2 \quad \checkmark
 \end{aligned}$$

e) 0 0111011 100110011001100110001100 |

$$\begin{aligned}
 &+ \quad 123 \quad 1.1001...100 \quad \begin{array}{l} 8388608 \\ 4194304 \end{array} \\
 &\quad -127 \\
 &= [1.10011001100110011001100 \times 10^{-0100}]_2 \\
 &= 0.00011001100110011001100 \\
 &= \left[\frac{1}{16} + \frac{1}{32} + \frac{1}{256} + \frac{1}{512} + \frac{1}{4096} + \frac{1}{8192} + \frac{1}{65536} + \frac{1}{131072} + \frac{1}{1048576} + \frac{1}{2097152} + \frac{1}{16777216} + \frac{1}{33554432} \right]_{10} \\
 &= \frac{3}{32} + \frac{3}{512} + \frac{3}{8192} + \frac{3}{131072} + \frac{3}{2097152} + \frac{3}{33554432} \\
 &= 3 \left(\frac{2^{20} + 2^{16} + 2^{12} + 2^8 + 2^4 + 1}{2^{26}} \right) \quad \checkmark
 \end{aligned}$$

3

Consider a decimal machine in which floating-point numbers are represented with a precision of 12 decimal places. Compute the relative error for the following numbers (assume numbers are rounded correctly):

a) $\underset{\uparrow}{x} \ 1.5673456545567890621$

$$\hat{x} = 1.56734565456$$

$$\varepsilon = \left| \frac{x - \hat{x}}{x} \right|$$

$$= 2.0480485531 \times 10^{-13}$$

b) $\underset{\uparrow}{x} \ 1 - 2^{-7}$

$$\hat{x} = 1 - 2^{-7}$$

$$\varepsilon = \left| \frac{x - \hat{x}}{x} \right|$$

$$= 0$$

c) $\underset{\uparrow}{x} \ 1 - 10^{-13}$

$$x = 0.99999999999999$$

$$\hat{x} = 1$$

$$\varepsilon = \left| \frac{x - \hat{x}}{x} \right|$$

$$= 1.00000000000001$$

d) $\underset{\uparrow}{x} \ 2.3456 - 0.00000456789011234$

$$\hat{x} = 2.34559543210988966$$

$$\hat{x} = 2.34559543211$$

$$\varepsilon = \left| \frac{x - \hat{x}}{x} \right|$$

$$\approx 4.9900234011 \times 10^{-14}$$

e) $\underset{\uparrow}{x} \ 3.14562345678912 - 3.145623451233476$

$$\hat{x} = 5.555644 \times 10^{-9}$$

$$\hat{x} = 5.556 \times 10^{-9}$$

$$\varepsilon = \left| \frac{x - \hat{x}}{x} \right|$$

$$\approx 0.0000640789798626$$