

[6 points] Perform the following calculations supposing that the binary values represent floating point numbers according to the IEEE 754 half-precision format (1 sign bit, 5-bit biased exponent, and 10-bit fraction). *Hint: check whether each binary value represents a special number before performing the calculations.*

- (a) $0111110000000000 + 0111110000000000$
- (b) $0100010000000010 - 1011110000000100$
- (c) $0011110000000010 \times 000100000000100$
- (d) $1111111111111111 \div 0100011000000001$

(a) $0111110000000000 + 0111110000000000$

- i) Check for special cases
 - Both numbers represent $+\infty$
 - Result is $+\infty$

Result = 0 1111 0000000000

(b) $0100010000000010 - 1011110000000100$

- i) Check for special cases
 - no special cases
- ii) Transform subtraction to addition and negate second number
 - $0\ 10001\ 0000000010 + 0\ 01111\ 0000000100$
- iii) Align
 - Second number has smaller exponent
 - Add 2 to its exponent and shift its fraction to the right twice
 - Exponent of second number = 10001
 - Significand of second number = 0.0100000001
- iv) Add significands (taking signs into consideration)
 - Significand of result = $1.0000000010 + 0.0100000001 = 1.0100000011$
 - Sign bit of result = 0
- v) Normalize
 - Significand is already normalized
 - Exponent of result = 10001
 - Fraction of result = 0100000011
- vi) Round
 - Not needed
 - Fraction of result = 0100000011

Result = 0 10001 0100000011

(c) $00111100000000010 \times 0001000000000100$

- i) Check for special cases
→ no special cases
- ii) Add exponents (and subtract the bias)
→ Exponent of result = $01111 + 00100 - 01111 = 00100$
- iii) Calculate sign of result
→ Sign bit of result = 0
- iv) Multiply significands
→ Significand of result = $1.0000000010 * 1.0000000100$
= $00.00000001000000001000 + 01.00000000100000000000$
= 01.00000001100000001000
- v) Normalize
→ Significand is already normalized
→ Exponent of result = 00100
→ Fraction of result = 00000001100000001000
- vi) Round
→ Suppose "rounding to nearest" is used
→ Fraction of result = ~~00000001100000001000~~

Result = 0 00100 0000000110

(d) $1111111111111111 \div 0100011000000001$

- i) Check for special cases
→ First number is NaN
→ Result is NaN

Result = 1 11111 1111111111

(e) Perform the following calculations by interpreting the binary values according to the 12-bit floating-point format and rounding the results to the nearest representable number if needed:

i. $0\ 0000\ 0000000 \div 0\ 1111\ 0000000$

ii. $0\ 1011\ 1000001 \times 0\ 1001\ 1000000$

i. $0\ 0000\ 0000000 \div 0\ 1111\ 0000000$

* Check special cases

→ zero divided by ∞

→ result = zero = $0\ 0000\ 0000000$

ii. $0\ 1011\ 1000001 \times 0\ 1001\ 1000000$

* Check special cases

→ no special case

* Add exponents (and subtract bias)

→ Exponent of result = $1011 + 1001 - 0111 = 1101$

* Calculate sign of result

→ Sign of result = $\boxed{0}$

+ve * +ve

* Multiply significands

$$\begin{array}{r}
 1.1000001 \\
 1.1000000 \\
 \hline
 1100000100000000 \\
 1100000100000000 \\
 \hline
 1001000011000000
 \end{array}$$

product

★ Multiply significands

$$\begin{array}{r}
 1.1000001 \\
 1.1000000 \\
 \hline
 1100000100000000 \\
 1100000100000000 \\
 \hline
 100100001100000000
 \end{array}$$

product

★ Normalize

→ Shift product right one-bit position and increment exponent by 1

→ Product = 1.0010000110000000

→ Exponent = 1101 + 1 = 1110

★ Round to nearest

→ Fraction = $001000\boxed{0}1000000$
 $001000\boxed{1}1000000$ greater than 1000000

Result = 0 1110 0010001

[9 points] Consider the IEEE 754 half-precision format in which floating point numbers are represented using 16 bits: 1 sign bit, 5-bit biased exponent, and 10-bit fraction.

(a) Convert the following numbers to their IEEE half-precision counterparts:

i. 109.6875

ii. -6.103515625E-5

(a)
 9) i. $109.6875 = 1.7138671875 \times 2^6$

$\rightarrow 0 \overset{15+6}{10101} 1011011011$

ii. $-6.103515625 \times 10^{-5} = 1.0 \times 2^{-14}$

$\rightarrow 1 \overset{15}{00001} 0000000000$