

# Hardware - Tutorial 1 - Correction

## Two's Complement and Floating point representation

### Exercise 1: Signed Numbers

Encode the following numbers into 8 bits signed binary.

- $-1 = 1111\ 1111$
- $-29 = 1110\ 0011$
- $-42 = 1101\ 0110$
- $-127 = 1000\ 0001$
- $-128 = 1000\ 0000$
- $-175 = \text{Impossible in 8 bits binary}$

### Exercise 2: Signed Operations

Perform the following 8 bits binary operations. Give the result in 8 bits binary and convert it to decimal given the context is signed or unsigned. If an overflow occurs, write down 'ERROR' in the corresponding cell.

Operation	Binary Result	Decimal Value	
		Unsigned	Signed
$1111\ 0101 + 1111\ 1010$	1110 1111	ERROR	-17
$1110\ 1000 - 1100\ 0110$	0010 0010	34	34
$0101\ 1110 - 1001\ 1110$	1100 0000	ERROR	ERROR
$0111\ 1110 + 0000\ 0101$	1000 0011	131	ERROR
$1100\ 1011 - 0001\ 1010$	1011 0001	177	-79
$1000\ 0000 + 1111\ 1010$	0111 1010	ERROR	ERROR
$1000\ 0011 - 0000\ 1010$	0111 1001	121	ERROR

## Exercice 3: Decimal to float

Convert the following decimal numbers to their binary **single-precision** floating point representation:

- 128
- $-32.75$
- 18.125
- 0.0625

**Solution:**

1)  $N = 128$

- $S = 0$
- $m = |128_{10}| = 128_{10} = 10000000_2$
- $\mathbf{M} = 0$
- $e = 7$
- $\mathbf{E} = e + \text{bias} = 134_{10}$
- $\mathbf{E} = 10000110_2$
- $128_{10} = 0\ 10000110\ 000000000000000000000000_f$

2)  $N = -32.75$

- $S = 1$
- $m = |-32.75_{10}| = 32.75_{10} = 100000.11_2$
- $\mathbf{M} = 0000011$
- $e = 5$
- $\mathbf{E} = e + \text{bias} = 132_{10}$
- $\mathbf{E} = 10000100_2$
- $-32.75_{10} = 1\ 10000100\ 000001100000000000000000_f$

3)  $N = 18.125$

- $S = 0$
- $m = |18.125_{10}| = 18.125_{10} = 10010.001_2$
- $\mathbf{M} = 0010001$
- $e = 4$
- $\mathbf{E} = e + \text{bias} = 131_{10}$
- $\mathbf{E} = 10000011_2$
- $18.125_{10} = 0\ 10000011\ 001000100000000000000000_f$

4)  $N = 0.0625$

- $S = 0$
- $m = |0.0625_{10}| = 0.0625_{10} = 0.0001_2$
- $\mathbf{M} = 0$
- $e = -4$
- $\mathbf{E} = e + \text{bias} = 123_{10}$
- $\mathbf{E} = 01111011_2$
- $0.0625_{10} = 0\ 01111011\ 000000000000000000000000_f$

## Exercise 4: Float to decimal

Convert the following **single-precision** floating point numbers to decimal representation:

- 1011 1101 0100 0000 0000 0000 0000 0000
- 0101 0101 0110 0000 0000 0000 0000 0000
- 1100 0001 1111 0000 0000 0000 0000 0000
- 1111 1111 1000 0000 0000 0000 0000 0000
- 0000 0000 0100 0000 0000 0000 0000 0000

**Solution:**

1)  $F = 1011\ 1101\ 0100\ 0000\ 0000\ 0000\ 0000\ 0000$

- $S = 1$
- $E = 01111010$
- $M = 1$
- $m = 1.M = 1.1$
- $e = E - \text{bias} = 122 - 127 = -5$
- $N = -1 \times 1.1_2 \times 2^{-5} = -1 \times 3_{10} \times 2^{-6}$
- $N = -0.046875$

2)  $F = 0101\ 0101\ 0110\ 0000\ 0000\ 0000\ 0000\ 0000$

- $S = 0$
- $E = 10101010$
- $M = 11$
- $m = 1.M = 1.11$
- $e = E - \text{bias} = 170 - 127 = 43$
- $N = 1 \times 1.11_2 \times 2^{43} = 1 \times 111_2 \times 2^{41}$
- $N = 7 \times 2^{41} \approx 1.53 \times 10^{13}$

3)  $F = 1100\ 0001\ 1111\ 0000\ 0000\ 0000\ 0000\ 0000$

- $S = 1$
- $E = 10000011$
- $M = 111$
- $m = 1.M = 1.111$
- $e = E - \text{bias} = 131 - 127 = 4$
- $N = -1 \times 1.111_2 \times 2^4 = -1 \times 1111_2 \times 2^1$
- $N = -30$

4)  $F = 1111\ 1111\ 1000\ 0000\ 0000\ 0000\ 0000\ 0000$

- $S = 1$
- $E = 11111111$
- $M = 0$
- $E = 255, M = 0 \rightarrow \text{Infinity}$
- $N = -\infty$

5)  $F = 0000\ 0000\ 0100\ 0000\ 0000\ 0000\ 0000\ 0000$

- $S = 0$
- $E = 0$
- $M = 1$
- $E = 0, M \neq 0 \rightarrow \text{Denormalized Mantissa}$
- $m = 0.M = 0.1$
- $e = 1 - \text{bias} = -126$
- $N = 1 \times 0.1_2 \times 2^{-126} = -1 \times 1_2 \times 2^{-127}$
- $N = 2^{-127} \approx 5.88 \times 10^{39}$

## Exercise 5: Decimal to double

Convert the following decimal numbers into their binary **double-precision** floating point representation:

- 1
- -64
- 12.06640625
- 0.2734375

**Solution:**

1)  $N = 1$

- $S = 0$
- $m = |1_{10}| = 1_{10} = 1_2$
- $M = 0$
- $e = 0$
- $E = e + bias = 1023_{10}$
- $E = 0111111111_2$
- $1_{10} = 0\ 0111111111\ 00\dots 0_d$

2)  $N = -64$

- $S = 1$
- $m = |-64_{10}| = 64_{10} = 1000000_2$
- $M = 0$
- $e = 6$
- $E = e + bias = 1029_{10}$
- $E = 10000000101_2$
- $-64_{10} = 1\ 10000000101\ 00\dots 0_d$

3)  $N = 12.06640625$

- $S = 0$
- $m = |12.06640625_{10}| = 12.06640625_{10} = 1100.00010001_2$
- $M = 10000010001$
- $e = 3$
- $E = e + bias = 1026_{10}$
- $E = 10000000010_2$
- $12.06640625_{10} = 0\ 10000000010\ 100000100010\dots 0_d$

4)  $N = 0.2734375$

- $S = 0$
- $m = |0.2734375_{10}| = 0.2734375_{10} = 0.0100011_2$
- $M = 00011$
- $e = -2$
- $E = e + bias = 1021_{10}$
- $E = 0111111101_2$
- $0.2734375_{10} = 0\ 0111111101\ 000110\dots 0_d$

## Exercise 6: Double to decimal

Convert the following **double-precision** floating point numbers to decimal representation:

- 403D 4800 0000 0000
- C040 0000 0000 0000
- BFC0 0000 0000 0000
- 8000 0000 0000 0000
- FFF0 0001 0000 0000

**Solution:**

1)  $D = 403D\ 4800\ 0000\ 0000$ 

- $D_2 = 0100\ 0000\ 0011\ 1101\ 0100\ 1000\dots 0$
- $\mathbf{S} = 0$
- $\mathbf{E} = 100\ 0000\ 0011$
- $\mathbf{M} = 1101\ 0100\ 1000\dots$
- $m = 1.M = 1.1101\ 0100\ 1000$
- $e = E - \text{biais} = 1027 - 1023 = 4$
- $\mathbf{N} = 1 \times 1.1101\ 0100\ 1000_2 \times 2^4 = 1110101001_2 \times 2^{-5}$
- $\boxed{\mathbf{N} = 937 \times 2^{-5}}$

2)  $D = C040\ 0000\ 0000\ 0000$ 

- $D_2 = 1100000000100\dots 0$
- $\mathbf{S} = 1$
- $\mathbf{E} = 100\ 0000\ 0100$
- $\mathbf{M} = 0\dots$
- $m = 1.M = 1.0$
- $e = E - \text{biais} = 1028 - 1023 = 5$
- $\mathbf{N} = 1 \times 1.0_2 \times 2^5$
- $\boxed{\mathbf{N} = -32}$

3)  $D = BFC0\ 0000\ 0000\ 0000$ 

- $D_2 = 1011\ 1111\ 1100\dots 0$
- $\mathbf{S} = 1$
- $\mathbf{E} = 011\ 1111\ 1100$
- $\mathbf{M} = 0$
- $m = 1.M = 1.0$
- $e = E - \text{biais} = 1020 - 1023 = -3$
- $\mathbf{N} = 1 \times 1.0_2 \times 2^{-3}$
- $\boxed{\mathbf{N} = -0.125}$

4)  $D = 8000\ 0000\ 0000\ 0000$ 

- $D_2 = 1000\dots 0$
- $\mathbf{S} = 1$
- $\mathbf{E} = 0$
- $\mathbf{M} = 0$
- $E = 0, M = 0 \rightarrow \text{Zero}$
- $\boxed{\mathbf{N} = -0}$

5)  $D = FFF0\ 0001\ 0000\ 0000$ 

- $D_2 = 1111\ 1111\ 1111\ 0000\ 0000\ 0000\ 0000\ 0001\dots 0$
- $\mathbf{S} = 1$
- $\mathbf{E} = 111\ 1111\ 1111$
- $\mathbf{M} = 0000\ 0000\ 0000\ 0000\ 0001$
- $E = 2047, M \neq 0 \rightarrow NaN$
- $\boxed{\mathbf{N} = NaN}$



- Therefore, the variable  $f1$  must be less than  $2^{28}$  so that the addition will affect  $f3$
- Wich gives:  
 $f1 < 2^{28}$   
 $10^n < 2^{28}$   
 $n < \text{Log}(2^{28})$   
 $n < 8.42$   

$n_{max} = 8$

3) Assuming that  $f1$ ,  $f2$ ,  $f3$  and  $r$  are declared as double, what is the largest value of  $n$  that still gives a correct value of  $r$  ?

**Solution:**

With the same line of reasoning:

$$\begin{aligned} f1 &< 2^{5+52} \\ f1 &< 2^{57} \\ 10^n &< 2^{57} \\ n &< \text{Log}(2^{57}) \\ n &< 17.15 \\ \div n_{max} &= 17 \end{aligned}$$