

Floating Point Reference Sheet

*(using an 8-bit representation like pg.86 of CS324 book) – look at the examples there as well

***Bolded words** are key differences between Denormalized and Normalized

Denormalized

<i>Sign</i>	<i>Exponent (exp)</i>				<i>Frac</i>		
1/0	0	0	0	0	1/0	1/0	1/0

$e = 0$ (since exp always has 0's)

$\text{Bias} = (2^{(\# \text{ of exp bits} - 1)}) - 1 = 2^{(4-1)} - 1 = 7$

$E = 1 - \text{bias} = 1 - 7 = -6$

$f = (\text{Value of Frac As Int}) / (2^{(\# \text{ of frac Bits})}) = (\text{Value of Frac as Int}) / (2^3)$

$M = f$

(sign) = -1 if Sign=1

(sign) = +1 if Sign=0

$V = M \times (2^E) \times (\text{sign}) = M \times (2^{(-6)}) \times (\text{sign})$

Normalized

<i>Sign</i>	<i>Exponent (exp)</i>				<i>Frac</i>		
1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0

Exp field needs at least one '1' in it and at least one '0' in it (i.e. Range 0001 to 1110 in binary).

$e = \text{int value of exp field}$

$\text{Bias} = (2^{(\# \text{ of exp bits} - 1)}) - 1 = 2^{(3)} - 1 = 7$

$E = e - \text{bias}$

$f = (\text{Value of Frac As Int}) / (2^{(\# \text{ of frac Bits})}) = (\text{Value of Frac as Int}) / (2^3)$

$M = f + 1 = ((\text{Value of frac as Int}) / (2^3)) + (2^3) / (2^3)$

(sign) = -1 if Sign=1

(sign) = +1 if Sign=0

$V = M \times (2^E) \times (\text{sign})$

Infinity (+oo / - oo)

<i>Sign</i>	<i>Exponent (exp)</i>				<i>Frac</i>		
1/0	1	1	1	1	0	0	0

NaN (Not a Number)

<i>Sign</i>	<i>Exponent (exp)</i>				<i>Frac</i>		
1/0	1	1	1	1	<-----	1 anywhere	----->

Frac contains at least one '1' in it

pg. 90 Default is to Round to Even i.e. On tie round to nearest even number

<i>Mode</i>	<i>\$1.40</i>	<i>\$1.60</i>	<i>\$1.50</i>	<i>\$2.50</i>	<i>\$-1.50</i>
Round-to-even	\$1	\$2	\$2	\$2	\$-2

Example Denormalized

0 0000 010

$e = 0$

$E = -6$

$f = (0b010 / (2^3)) = 2 / (2^3) = (2/8)$

$M = 2/8$

$\text{Value} = (2/8) \times (2^{-6}) = (2/(2^3)) \times (1/(2^6)) = (2/(2^9)) = (2/512) = (1/256)$

Example Normalized

1 0111 001

$e = 0b0111 = 7$

$E = e - 7 = 7 - 7 = 0$

$f = (0b001 / (2^3)) = (1/8)$

$M = (1/8) + 1 = (1/8) + (8/8) = (9/8)$

$\text{Value} = (9/8) \times (2^0) \times (-1) = (-9/8)$

Side Notes (pg.83-84)

Single-precision floating-point format (a float in C)

$s = 1$ bit

$\text{exp}(k) = 8$ bits

$\text{frac}(n) = 23$ bits

32 bits total

Double-precision floating-point format (a double in C)

$s = 1$ bit

$\text{exp}(k) = 11$ bits

$\text{frac}(n) = 52$ bits

64 bits total