

Assignment 6

Sonntag, 22. Juli 2018

20:32

Q format	Data type from <stdint.h>	Number of sign bits	Number of integer bits	Number of fractional bits	Total amount of bits	Minimum value	Maximum value
Q1.13	<i>int16_t</i>	3	0	13	16	-1	0.9998
U Q2.14	<i>int16_t</i>	0	2	14	16	0	3.9999
Signed Q0.15	<i>int16_t</i>	1	0	15	16	-1	0.9999
Q4.4 + Q4.4	<i>int16_t</i>	8	4	4	16	-16	15.9375
Q7 * Q15	<i>int32_t</i>	10	0	22	32	-1	0.9999
UQ1.2 + Q2.2	<i>int8_t</i>	4	2	2	8	-4	3.7500

- Notation
 - Q_n (implicit) or $Q_{m,n}$ (explicit)
 - n ... Number of fractional bits
 - m ... Number of integer bits
 - inclusive (or exclusive) sign bit
 - U $Q_{m,n}$... unsigned designator
- Value range
 - Unsigned: $[0, (2^m - 2^{-n})]$
 - $m \geq 0$
 - Signed: $[-(2^{m-1}), (2^{m-1} - 2^{-n})]$
 - $m \geq 1$ (sign bit included)
- Adding
 - $Q_{m_1, n_1} + Q_{m_2, n_2} = Q_{m,n}$
 - Fraction n : $\max(n_1, n_2)$
 - Integer m : $\max(m_1, m_2) + 1$
- Multiplying
 - $Q_{m_1, n_1} * Q_{m_2, n_2} = Q_{m,n}$
 - Fraction n : $n_1 + n_2$
 - Integer m : $m_1 + m_2$
 - Signed * Signed: $m_1 + m_2 - 1$ (extra sign bit)
- Convert from float
 - $q = f \cdot 2^n$
 - Problem
 - Truncation or rounding?
- Convert to float
 - $f = q \cdot 2^{(-n)}$

• $Q4.4 + Q4.4 = Q5.4$ • $UQ1.2 + Q2.2 = Q3.2$

• $Q7 * Q15 = Q22$

- Calculate the result by hand

- $2.125 [Q3.3] + (-1.5) [Q1.2]$
- $-0.125 [Q1.4] * 0.725 [Q1.4]$

• $2.125 [Q3.3] + (-1.5) [Q1.2]$

$Q3.3 + Q1.2 = Q4.3 \Rightarrow SIIIFFF$

$$\begin{array}{r|l} 2.125 & \rightarrow 0.125 \cdot 2 = 0.250 \quad | 0 \\ & 0.250 \cdot 2 = 0.500 \quad | 0 \\ & 0.500 \cdot 2 = 1.000 \quad | 1 \\ \hline & 10b \end{array}$$

$= 0010.001$

$$\begin{array}{r} 0010.001 \\ 1110.100 \\ \hline 0000.101 \end{array}$$

$2^{-3} \cdot 5 = 0.625$

$$\begin{array}{r|l} -1.5 & \rightarrow 0.500 \cdot 2 = 1 \quad | 1 \\ & \downarrow \\ & 01b \end{array}$$

$= 0001.100$

$$\begin{array}{r} 0001.100 \\ 1110.011 \\ \hline 1110.100 \end{array}$$

• $-0.125 [Q1.4] \cdot 0.725 [Q1.4]$

$Q1.4 \cdot Q1.4 = Q1.8 \Rightarrow SFFFFFFF$

$$\begin{array}{r|l} -0.125 & \rightarrow 0.125 \cdot 2 = 0.250 \quad | 0 \\ & 0.250 \cdot 2 = 0.500 \quad | 0 \\ & 0.500 \cdot 2 = 1.000 \quad | 1 \\ \hline & 0b \end{array}$$

$$\begin{array}{r|l} 0.725 & \rightarrow 0.725 \cdot 2 = 1.450 \quad | 1 \\ & 0.450 \cdot 2 = 0.900 \quad | 0 \\ & 0.900 \cdot 2 = 1.800 \quad | 1 \\ \hline & 0b \end{array}$$

$$\begin{array}{rcl} -0.125 & \rightarrow & 0.125 \cdot 2 = 0.250 \quad | \quad 0 \\ & & 0.250 \cdot 2 = 0.500 \quad | \quad 0 \\ & & 0.500 \cdot 2 = 1.000 \quad | \quad 1 \\ \hline & & 0b \end{array}$$

$$= 0.0010$$

$$1.1101$$

$$+1$$

$$1.1110$$

$$1.1110 \cdot 0.1011$$

$$0.0000$$

$$1.11110$$

$$0.000000$$

$$1.1111110$$

$$1.11111110$$

$$1.1111$$

$$1.111111$$

$$1.11101010$$

$$\begin{array}{rcl} 0.725 & \rightarrow & 0.725 \cdot 2 = 1.450 \quad | \quad 1 \\ & & 0.450 \cdot 2 = 0.900 \quad | \quad 0 \\ & & 0.900 \cdot 2 = 1.800 \quad | \quad 1 \\ & & 0.800 \cdot 2 = 1.600 \quad | \quad 1 \\ \hline & & 0b \end{array}$$

$$= 0.1011$$

- What is the best minimum Q-format and resulting error to store

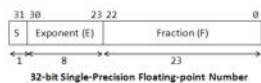
- 3.141592

Leider nicht berechnet

- Proof by hand that for IEEE 754 binary32 floating point format
 - the largest positive number representable is roughly $3.4e+38$
 - the smallest positive number representable is roughly $1.4e-45$

- Single precision

- 1 sign bit (S)
- 8 exponent bits (E)
- 23 fraction bits (F)
- 22 effective fraction bits



- Normalized form

$$(-1)^S \cdot 2^{(E-127)} \cdot 1.F$$

- Denormalized form

$$(-1)^S \cdot 2^{(-126)} \cdot 0.F$$



S

EEEE	EEEE	FFFF	FFFF	FFFF	FFFF	FFFF	FFF
11	11	1110	1111	1111	1111	1111	1111

254 8388607

$$\begin{aligned} \max &= (-1)^S \cdot 2^{(E-127)} \cdot 1.F \Rightarrow \text{normalized Form} \\ &= (-1)^0 \cdot 2^{(254-127)} \cdot (1 + 8388607 \cdot 2^{-23}) \\ &= 3.4028 \cdot 10^{38} \end{aligned}$$

S

EEEE	EEEE	FFFF	FFFF	FFFF	FFFF	FFFF	FFF
0000	0000	0000	0000	0000	0000	0000	001

0 1

$$\begin{aligned} \max &= (-1)^S \cdot 2^{(-126)} \cdot 0.F \Rightarrow \text{denormalized Form} \\ &= (-1)^0 \cdot 2^{(-126)} \cdot (1 \cdot 2^{-23}) \\ &= 1.40129 \cdot 10^{-45} \end{aligned}$$

- $0.1f + 0.2f \neq 0.3f$

0.1f:

$0.1 \cdot 2 = 0.2$	0	} Shift um $n=4 \Rightarrow 2^{-4}$	$0.6 \cdot 2 = 1.2$	1	$2^{-n} = 2^{E-127}$ $2^{-4} = 2^{E-127} \quad \log_2$ $-4 = E-127 \quad +127$ <u>$E = 123$</u>
$0.2 \cdot 2 = 0.4$	0		$0.2 \cdot 2 = 0.4$	0	
$0.4 \cdot 2 = 0.8$	0		$0.4 \cdot 2 = 0.8$	0	
$0.8 \cdot 2 = 1.6$	1		$0.8 \cdot 2 = 1.6$	1	

S EEEE EEEE FFFF FFFF FFFF FFFF FFF
 0 0111 1011 1001 1001 1001 1001 101

123 5033165

$$\text{Value} = (-1)^S \cdot 2^{(E-127)} \cdot 1F$$

$$\text{Value} = (-1)^0 \cdot 2^{(123-127)} \cdot (1 + 5033165 \cdot 2^{-23}) = 0.1000000149$$

0.2f:

$0.2 \cdot 2 = 0.4$	0	} Shift um $n=3 \Rightarrow 2^{-3}$	$2^n \cdot 2^{E-127}$
$0.4 \cdot 2 = 0.8$	0		$2^{-3} = 2^{E-127} \quad \log_2$
$0.8 \cdot 2 = 1.6$	1		$-3 = E-127 \quad +127$
$0.6 \cdot 2 = 1.2$	1		<u>$E = 124$</u>

S EEEE EEEE FFFF FFFF FFFF FFFF FFF
 0 0111 1100 1001 1001 1001 1001 101

124 5033165

$$\text{Value} = (-1)^0 \cdot 2^{(124-127)} \cdot (1 + 5033165 \cdot 2^{-23}) = 0.2000000298$$

0.3f:

$0.3 \cdot 2 = 0.6$	0	} Shift um $n=2 \Rightarrow 2^{-2}$	$2^{-2} = 2^{(E-127)}$
$0.6 \cdot 2 = 1.2$	1		$E = 125$
$0.2 \cdot 2 = 0.4$	0		
$0.4 \cdot 2 = 0.8$	0		

S EEEE EEEE FFFF FFFF FFFF FFFF FFF
 0 0111 1101 0011 0011 0011 0011 010

125 1677722

$$\text{Value} = (-1)^0 \cdot 2^{(125-127)} \cdot (1 + 1677722 \cdot 2^{-23}) = 0.3000001192$$

$$0.1000000149 + 0.2000000298 \neq 0.3000001192 \quad \checkmark$$

$$0.1f + 0.2f \neq 0.3f \quad \checkmark$$