



TÉCNICO
LISBOA

Instituto Superior Técnico Sistemas de Processamento Digital de Sinais (SPDSina)

Fixed-point numerical representation

Consider the real numbers $x = 17.35$ and $y = 0.15$ and its processing using fixed point arithmetic. Assume a symmetric two's complement representation: the number w is represented in the interval $-(2^i - 2^{-m}) \leq w \leq 2^i - 2^{-m}$ where i is the number of integer bits and m the number of fractional bits (total number of bits is $n=1+i+m$).

- Determine the arithmetic formats which allow the most accurate representation of x and y with $n=16$ bit and $z = x \cdot y$ with 32 bit.
- Determine the values of x and y in these formats and the resulting value of z , z_{real} . Compute the relative absolute error of z , $\varepsilon_{\text{rel}}(z) = \left| \frac{z_{\text{ideal}} - z_{\text{real}}}{z_{\text{ideal}}} \right|$. How could this error be made smaller? Compute the same error if z is represented with 16 bit instead of 32.
- Assume you know *a priori* the values of the operands x and therefore of z . Determine the arithmetic formats which allow the most accurate representation of z and determine the relative error.

Exercise: $x = 17.35$, $y = 0.15$, fixed point 16 bit

$$z = x \cdot y = 2.6025 = z_{true}$$

$y_{Q_{15}}$, $|x| > 16$, $< 32 \rightarrow$ need 5 integer bit $\rightarrow Q_{10}$

a) Assume $y_{Q_{15}}$, $x_{Q_{10}}$ but actual values are unknown.

$$z = x \cdot y \rightarrow \text{needs } 5+0 = 5 \text{ integer bits} \Rightarrow Q_{10} \text{ or } Q_{25} (32 \text{ bit})$$

$$x = \text{round}(2^{10} \times 17.35) = \text{round}(17766.4) = 17766$$

$$y = \text{round}(2^{15} \times 0.15) = \text{round}(4915.2) = 4915$$

b) $z = 4915 \times 17766 = 87319890$ but is in Q_{25} because of extra sign bit

$$z = 2 \times 87319890 = 174639780$$

$$z_{real} = \frac{174639780}{2^{26}} = 2.602335512638...$$

$$E_n = \left| 1 - \frac{z_{real}}{2.6025} \right| = 6.32 \times 10^{-5} \approx 10^{-4.2} \quad (4.2 \text{ decimal places})$$

Code: $\text{Int16 } x = 17766, y = 4915, z_{16};$ $\log_{10} 2^{-15} \approx$

$\text{Int32 } z$ (z_{Int32})

$$z = (x * y) \ll 1; \quad (Q_{26}) \quad (\text{Int32})$$

$$z_{16} = (z \gg 16); \quad \text{or} \quad z_{16} = ((x * y) \ll 1) \gg 16;$$

c) In this case we know the values of the operands and the result which is $2.6025 \Rightarrow$ in fact need only 2 integer bits because $|z_{true}| < 4 \Rightarrow$ can store in Q_{29} (or Q_{13})

$$z = ((x * y) \ll 1) \ll 3;$$

$$\begin{aligned} z_{real} &= 174639780 \gg 13 \\ &= 21318 \\ &= \frac{21318}{2^{13}} \\ &= 2.6022949 \end{aligned}$$

$$E = 7.88 \times 10^{-5} \rightarrow \text{almost equal to 32 bit}$$

$$= 10^{-4.1} \quad \text{Note: symmetric interval:}$$

$$-(2^i - 2^{-m}) \leq w \leq 2^i - 2^{-m}$$

With 16 bit (z_{16})
 $z_{real} = 174639780 \gg 16$
 $= 2664 \quad (Q_{16})$

$$z_{real} = \frac{2664}{2^{10}} = 2.6015625$$

$$E = 3.6 \times 10^{-4} = 10^{-3.44}$$