

Sistemas de Processamento Digital de Sinais (SPDSina)

Digital Representation – Fixed Point

II — Consider the fractional numbers x and y both in Q_{14} format and their processing in a 16 bit fixed-point processor.

- a) Determine the format that should be used to represent $z = 2 \cdot (x^2 - y^2)$ with 16 bit such that z is always correct for all possible values of x and y . Justify.
- b) Consider now that $x = 1.25$ and $y = -0.75$. Write the most efficient and accurate fixed-point C code (32 bit processing when possible) which computes z in the format determined in a), including variable declarations and initializations, using only 16 bit words.
- c) Compute the true value of z and the value obtained with the computation in b). Explain the result.

$$z = 2 \cdot (x^2 - y^2), \quad x, y \text{ em } Q_{14} \rightarrow \text{Regras!}$$

a) formato de z ?

$$\begin{matrix} x^2 \rightarrow Q_{13} \\ y^2 \rightarrow \end{matrix} \quad x^2 - y^2 \rightarrow Q_{12} \quad 2(x^2 - y^2) \rightarrow Q_{11}$$

$$-(2 \cdot 2^{-14}) \leq x, y \leq 2 \cdot 2^{-14}$$

$$x^2 - y^2 \rightarrow \text{máximo quando } y=0, x^2 \text{ } Q_{13}$$

$$\rightarrow \text{mínimo quando } x=0, -y^2 \text{ } Q_{13}$$

$\Rightarrow z$ pode ser representado em Q_{12} ! e não $Q_{11} \Rightarrow$ 1 bit = mais! sem custo

b)

$$x = 1.25, \quad x_{Q_{14}} = \text{round}(1.25 \times 2^{14}) = 20480$$

$$y = -0.75, \quad y_{Q_{14}} = \text{round}(-0.75 \times 2^{14}) = -12288$$

$$\text{Int1: } x = 20480, \quad y = -12288, \quad z;$$

ficam MSBs

$$z = (((\text{Int32})x * x - (\text{Int32})y * y) \ll 1) \ll 1 \gg 1 \gg 16; \quad Q_{28} \rightarrow Q_{12}$$

$$x^2 - y^2 \equiv Q_{13}$$

$$\times 2 \quad Q_{13} - Q_{12}$$

ou deixar logo z em Q_{12}

sign

$$z = (((\text{Int32})x * x - (\text{Int32})y * y) \ll 1) \gg 16;$$

$$Q_{12}$$

$$\times 2$$

$$\text{MSB's}$$

$$z = 2(1.25^2 - 0.75^2) = z$$

$$z_{\text{real}} = (20480^2 - (-12288)^2) \times 2 \times \underbrace{\frac{1}{2^{16}}}_{\text{MSB's}} \times \underbrace{\frac{1}{2^{12}}}_{Q_{12}} = z$$

$z = z_{\text{real}} \Rightarrow \text{erro} = 0$ porque o resultado é representável por uma soma de potências (positivas e negativas) de 2 com 16 bit