# Finite - Wordlength Effect in Digital filters

The basic fundamental operations involved in DSP algorithms like convolution, Correlation DFT, FFT etc. are multiplication and addition.

These operations are performed using input Sequence $x(n)$, impulse sequence $h(n)$ and the Co-efficients of the difference equation governing the system.

The input and output data are stored in register in a digital system or different memory location in the form of binary values.

The maximum size of the binary information that can be stored in a memory location or register is known as "Word-length".

* When a register stores an 8-bit data then its wordlength is 8-bit.

* Quantization & Coding depends on register wordlength.

* Quantization & Coding will introduce error in input data, because analog data has infinite precision but digital equivalent has finite precision.

There are some error occurred due to the
process of finite precision representation
of binary numbers in digital system. These
errors are commonly called as " finite word
length effects or error".

* finite word length effects in digital filters—

(1) Errors due to quantization of I/p data by
A/D converter.

(2) Error due to quantization of filter
co-efficients.

(3) Errors due to rounding of products in
multiplication.

(4) Error due to overflow in addition.

* Representation of Numbers —

(a) fixed point representation

(b) floating " " .

(a) fixed point Representation — Bits allowed for
integer & fractional
part and so position of binary point is fixed.

Drawback → Due to the fixed integer & fractional

part — too large and too small values cannot be represented.

Types — (i) Sign magnitude form

$$(2.75)_{10} = (010.1100)_2$$

$$(-2.75)_{10} = (110.1100)_2$$

(ii) One's Complement form    (for negative)

(iii) Two's Complement form    ( "   "   ).

* (b) **Floating Point Representation** :-

$$N_f = M \times 2^E$$

M = Mantissa → it will be in binary fraction format

$$0.5 \leq M < 1$$

E = Exponent → either a positive or negative integer.

(i)  $+7_{10} = +111_2 = 0.1110 \times 2^3 = 0.111 \times 2^{+11_2}$

$$= 0.1110 \quad 011$$
$$\quad\quad M \quad\quad E$$

(ii)  $-7_{10} = -111_2 = 1.111 \times 2^{+3_{10}} = 1.1110 \times 2^{+11_2}$

$$= 1.1110 \quad 011$$
$$\quad\quad M \quad\quad E$$

## Quantization: (or) Truncation

In fixed point or floating point arithmetic, the size of the result of the operation may exceed the size of binary used in the number system. In this case, the low order bits has to be eliminated in order to store the result.

Method → * Truncation

* Rounding.

* Truncation → 8 bits to 4 bits

$$0.01110011 \longrightarrow 0.0111$$

* Rounding → 0.101010 rounded to 4 bits

$$0.1010 \text{ or } 0.1011.$$

Rounding error ⇒ $e_r = N_r - N$.

         ↓              ↓ Unquantized
      Quantized      number
      or rounded
      number

Range:

$$-\frac{2^{-b}}{2} \leq e_r \leq \frac{2^{-b}}{2}$$

$b$ = no. of rounded bits.

i/p ——→ | Sampler | Quantizer | ——→ o/p
$x(t)$                                  $x_q(n)$

Block-diagram of
A/D Converter

Fold Here

Fold Here

Fold Here

Fold Here

Fold Here

Fold Here

# Quantization in filter- Coefficient

The filter co-efficients are quantized to word Sizes of the register used to store them either by truncation or by rounding.
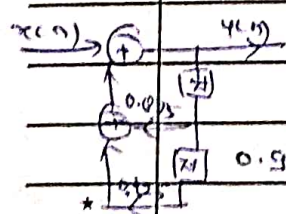
* Location of poles & zeros of digital filters directly depends on the value of filter coefficients. The quantization of the filter co-efficients will modify the value of poles and zeros and so the location of poles and zeros will be shifted from the desired location. This will create deviation in the frequency response of the system.

**SA:—**

$$H(z) = \frac{1}{(1-0.5\,z^{-1})(1-0.45\,x^{-1})}$$

3-bit- coefficient—

$$= \frac{x^2}{(z-0.5)(x-0.45)} = \frac{1}{(1-0.95\,x^{-1}+0.225\,x^{-2})}$$
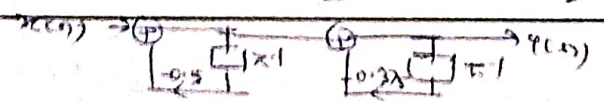


$0.95 \xrightarrow{D \to B} 0.1111_2 \xrightarrow{\text{trun to 3 bits}} 0.111_2 \xrightarrow{B \to D} 0.875_{10}$

$0.225 \xrightarrow{}_{10} 0.0011_2 \xrightarrow{} 0.001_2 \xrightarrow{B \to D} 0.125$

−for cascade :—  $H_1(z) = \frac{1}{1-0.5\,z^{-1}}$     $H_2(z) = \frac{1}{1-0.45\,x^{-1}}$

$0.5 \xrightarrow{D \to B} 0.1000_2 \xrightarrow{} 0.100_2 \xrightarrow{B \to D} 0.5_{10}$

$0.45 \xrightarrow{} 0.0111_2 \xrightarrow{} 0.011_2 \xrightarrow{} 0.375$

## Quantization in filter-Coefficient-

The filter co-efficients are quantized to word sizes of the register used to store them either by truncation or by rounding.

* Location of poles & zeros of digital filters directly depends on the value of filter coefficients. The quantization of the filter co-efficients will modify the value of poles and zeros and so the location of poles and zeros will be shifted from the desired location. This will create deviation in the frequency response of the system.
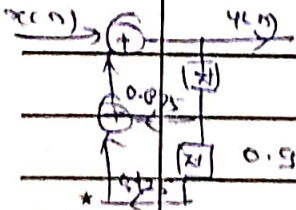
Ex:- $H(x) = \dfrac{1}{(1-0.5 x^{-1})(1-0.45 x^{-1})}$

B-bit-coefficient-

$$= \dfrac{x^2}{(x-0.5)(x-0.45)} = \dfrac{1}{(1-0.95 x^{-1}+0.225 x^{-2})}$$

$x(n) \rightarrow (+) \rightarrow y(n)$
$(Z)$
$0.95$
$(Z)$
$0.225$

$0.95 \xrightarrow{D \to B} 0.1111_2 \xrightarrow{trun.\ to\ 3\ bits} 0.111_2 \xrightarrow{B \to D} 0.075$
$\quad\quad\quad 16 \quad\quad\quad\quad\quad\quad\quad 2 \quad\quad\quad\quad 10$

$0.225 \xrightarrow{} 0.0011_2 \xrightarrow{} 0.001_2 \xrightarrow{B \to D} 0.125$
$\quad 16 \quad\quad\quad\quad 2 \quad\quad\quad\quad 2$

for cascade:- $H_1(x) = \dfrac{1}{1-0.5 x^{-1}}$   $H_2(z) = \dfrac{1}{1-0.45 x^{-1}}$

$0.5 \xrightarrow{D \to B} 0.1000_2 \xrightarrow{} 0.100_2 \xrightarrow{B \to D} 0.5_{10}$

$0.45 \xrightarrow{} 0.0111_2 \xrightarrow{} 0.011_2 \xrightarrow{} 0.375$

$x(n) \to (+) \boxed{z^{-1}}_{0.5} \to (+) \boxed{z^{-1}}_{0.25} \to y(n)$