

Analysis of finite Word length Effect in fixed-point DSP Systems:-

Finite word length Effects in Digital filters :-

Introduction :- Fundamental operations in various Computational Procedures like Convolution spectral estimation etc., in DSP are Multiplication and Addition. These operations are performed using sample of I/p sequence, sample of impulse, response and coefficients of differential equation governing system.

→ The Registers are basic storage device in digital systems. The Maximum size of binary information (d) data that can be stored in register is called Register word length.

Eg:- When register stores an 8-bit data then its word length is 8-bit.

For storing I/p data in register they have quantized and coded in binary. Where Quantization & Coding depends on register word length.

Eg:- If register word length is 8-bit we can generate $2^8 = 256$ binary codes, so we have 256 quantized levels.

The quantization and coding will introduce error in I/p data, bcoz the analog data has infinite precision but the digital equivalent has finite precision.

→ While performing computation size of result may be exceeding the size of register used storing the result.

Eg:- If result of addition of two eight bit data may be 8(8)9 bits and result of multiplication of two 8-bit data may be 16-bit.

→ In case of register used to store the result 8-bit the result has to be truncated (d) rounded to accommodate in the register. This makes system is non-linear and leads to limit cycle behaviour.

→ The effect of truncation or rounding can be represented in terms of an additive error signal, which is called round off noise.

* Finite word length effects in digital filters:-

- ① Error due to quantization of input data by A/D converter
- ② " " " " " filtered co-efficients
- ③ " " " " sounding the product in multiplication.
- ④ " " " " overflow in addition.
- ⑤ Limit cycles.

Q) Representation of Numbers in Digital Systems:-

In a radix representation the no. can be represented by a summation relation as

$$\text{Ex:- } 178.25_{10} = (1 \times 10^2) + (7 \times 10^1) + (8 \times 10^0) + (2 \times 10^{-1}) + (5 \times 10^{-2}) \\ = (d_2 \times 10^2) + (d_1 \times 10^1) + (d_0 \times 10^0) + (d_{-1} \times 10^{-1}) + (d_{-2} \times 10^{-2})$$

$$\text{Ex:- } 111.11_2 = (1 \times 2^2) + (1 \times 2^1) + (1 \times 2^0) + (1 \times 2^{-1}) + (1 \times 2^{-2}) \\ = (d_2 \times 2^2) + (d_1 \times 2^1) + (d_0 \times 2^0) + (d_{-1} \times 2^{-1}) + (d_{-2} \times 2^{-2})$$

i.e generally any number be represented as

$$N = \sum_{i=-A}^{B} d_i \gamma^{-i} \rightarrow ①$$

A = no. of integer digits

B = no. of fraction digits

γ = radix (or) Base

d_i = i th digit of the number.

where d_i ranges $0 \leq d_i \leq (\gamma - 1)$

$$\text{Ex:- if } \gamma = 2; d_i = 0, 1 \quad (8) 1$$

$$\text{where } \gamma = 10; d_i = 0, 1, 2, 3, 4, 5, 6, 7, 8, 9.$$

→ The radix $\delta = 2$ can be represented from binary number. (3)

$$\text{Binary number } N = \sum_{i=-A}^B d_i 2^i \rightarrow (2)$$

In binary digits d_{-A} is called Most Significant Bit (MSB) and binary digit d_B is called Least Significant Bit (LSB) of binary number 'N'.

The binary point (.) the digits d_0 and d_1 does not exists physically in digital systems.

→ In various computation procedures of Dsp all the fraction format bcz mixed no:- (no:- with integer and fraction parts) are difficult to multiply and the number of bits representing an integer cannot be reduced by truncation (0).
So for fraction format of binary no:- is modified as.

$$\text{Binary fraction no:- } N = \pm \sum_{i=1}^B d_i 2^{-i}$$

$$\text{Binary fraction number } N = \sum_{i=0}^B d_i 2^{-i} \rightarrow (3)$$

where d_0 is used to represent the sign of number.

There are 2-methods representing binary number.

(a) Fixed Point representation

(b) Floating Point "

(a) Fixed Point representation :- 3-Different formats

(i) Sign-magnitude format

(ii) Ones Complement "

(iii) Two's Complement "

In fixed point representation is only unique way of representing positive binary fraction no:-

$$\begin{aligned} \text{Positive binary fraction number } N_p &= 0.d_1 d_2 \dots d_B \\ &= (0 \times 2^0) + (d_1 \times 2^{-1}) + (d_2 \times 2^{-2}) + \dots + (d_B \times 2^{-B}) \\ &= \sum_{i=0}^B d_i 2^{-i} \text{ where } d_0 = 0 \rightarrow (4) \end{aligned}$$

where d_0 is set to zero to represent the positive sign.

(i) Sign magnitude format:-

The negative value of given no. differ only in sign bit (i.e. bit d_0)

The sign bit d_0 is zero for +ve and one for -ve number.

$$\therefore \text{Positive binary fraction no. } N_p = (0 \times 2^0) + \sum_{i=1}^B d_i 2^{-i} \rightarrow 5$$

$$\text{Negative " " " } N_n = (1 \times 2^0) + \sum_{i=1}^B d_i 2^{-i} \rightarrow 6$$

Eg:- $+0.125_{10} = 0.001_2 \text{ (or) } 0001_2 ; -0.125_{10} = 1.001_2 \text{ (or) } 1001_2$

Note:- The binary point b/w d_0 and d_1 is not necessarily below

(ii) One's Complement format:-

Here the no. is same as sign magnitude format and
-ve no is given by obtaining by bit by bit complement.

\therefore Complement of digit d_p can be obtained by subtracting the digit
from one.

$$\text{i.e. } d_p = \bar{d}_p = (1 - d_p) \rightarrow 7$$

$$\therefore \text{Negative binary fraction } \left\{ \begin{array}{l} N_{1c} = (1 \times 2^0) + \sum_{i=1}^B (1 - d_i) 2^{-i} \\ \text{no. is one's complement} \end{array} \right. \rightarrow 8$$

Eg:-

$$+0.125_{10} = 0.001_2 \text{ (or) } 0001_2 ;$$

$$-0.125_{10} = 1.110_2 \text{ (or) } 1110_2 .$$

(iii) Two's Complement:-

The the no. is same as sign magnitude format
-ve no. is obtained by taking ones complement of positive
representation and adding one to LSB.

$$\therefore \text{Negative binary fraction } \left\{ \begin{array}{l} N_{2c} = (1 \times 2^0) + \sum_{i=1}^B (1 - d_i) 2^{-i} + (1 \times 2^{-B}) \\ \text{no. is two's complement} \end{array} \right. \rightarrow 9$$

Eg:- $-0.125_{10} = 1.111_2 \text{ (or) } 1111_2 .$

$$+0.125_{10} = 0.001_2 \text{ (or) } 0001_2 ;$$

Decimal equivalent of 4-bit binary no:- in fixed Point representation:- (5)

Binary number in fixed Point representation.	Sign magnitude	One's complement	Two's complement	Decimal Equivalent
0.000	0.000	0.000	0.000	0
0.001	0.001	0.001	0.001	1/8
0.010	0.010	0.010	0.010	2/8
0.011	0.011	0.011	0.011	3/8
0.100	0.100	0.100	0.100	4/8
0.101	0.101	0.101	0.101	5/8
0.110	0.110	0.110	0.110	6/8
0.111	0.111	0.111	0.111	7/8
1.000	1.000	1.000	1.000	-0
1.001	1.001	1.001	1.001	-1/8
1.010	1.010	1.010	1.010	-2/8
1.011	1.011	1.011	1.011	-3/8
1.100	1.100	1.011	1.100	-4/8
1.101	1.101	1.010	1.011	-5/8
1.110	1.110	1.001	1.010	-6/8
1.111	1.111	1.000	1.001	-7/8
			1.000	-8/8

(b) Floating Point representation:-

Is employed to represent larger range of no:- in a given binary word size

$$\text{Floating Point no:- } N_f = M \times 2^E$$

M - mantissa ; E = exponent

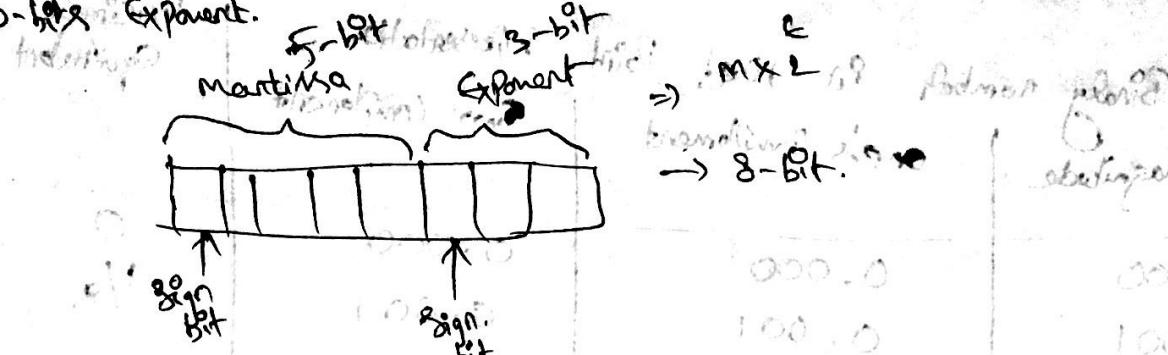
The value of M - ranges $0.5 \leq M \leq 1$ & M is binary frac format

E is called exponent and it is either positive (P) negative (N) integer.

LSB of mantissa/exponent represents sign (either +ve/-ve).

+
0
-

in floating point representation explained by 5-bit mantissa and 3-bit exponent.



The range of floating point represented is

$$\pm \left(\frac{1}{2} \times 2^{-3} \right) \text{ to } \pm \left(2 - \frac{1}{2} \right) \times 2^{3-1}, \text{ i.e. } (\pm 7.8125 \times 10^{-3} \text{ to } \pm 15.5)$$

Ex:-

$+5_{10} = +101_2 = 0.1010 \times 2^{+3} = 0.1010 \times 2 = \boxed{0|1|0|1|0|0|1|1}$

$-5_{10} = -101_2 = 1.1010 \times 2^{+3} = 1.1010 \times 2^{+11_2} = \boxed{1|1|0|1|0|0|1|1}$

$+0.125_{10} = +0.001_2 = 0.0010 \times 2^0 = 0.1000 \times 2^{-2} = 0.1000 \times 2^{-10_2}$

$= \boxed{0|0|1|0|0|0|1|1|0}$

Complement

$-0.125_{10} = -0.001_2 = 1.0010 \times 2^0 = 1.0010 \times 2^{-11_2} = 1.1000 \times 2^{-2}$

$= 1.1000 \times 2^{-10_2} = \boxed{1|1|0|0|0|1|1|0}$

2) Types of Arithmetic in Digital Systems :-

classified broad classes

fixed point Arithmetic floating point Arithmetic.

3-types

sign magnitude Arithmetic
one's complement Arithmetic
two's complement Arithmetic

(avoided in digital system due to difficulties in handling during addition)

1's Complement Addition:-

(7)

i- Add $+0.375$ and -0.625 by 1's Complement Addition.

Convert Decimal to binary

$$\begin{array}{r}
 0.375 \\
 \times 2 \\
 \hline
 0.750 \\
 \times 2 \\
 \hline
 1.500 \\
 \times 2 \\
 \hline
 0.000
 \end{array}$$

0.011_2

$$\begin{array}{r}
 0.625 \\
 \times 2 \\
 \hline
 1.250 \\
 \times 2 \\
 \hline
 0.500 \\
 \times 2 \\
 \hline
 0.000
 \end{array}$$

0.101_2

Step - 2:- Representation is given as

$$\begin{array}{c}
 +0.375_{10} \xrightarrow{\substack{\text{Convert to} \\ \text{binary}}} +0.011_2 \xrightarrow{\substack{\text{add sign} \\ \text{bit}}} 0.011_2 \\
 -0.625_{10} \xrightarrow{\substack{\text{Convert to} \\ \text{binary}}} -0.101_2 \xrightarrow{\substack{\text{add sign} \\ \text{bit}}} 1.101_2 \xrightarrow{\substack{\text{Complement} \\ \text{Fraction Part}}} 1.010_2
 \end{array}$$

Step - 3 Addition

$$\begin{array}{r}
 (+0.375)_{10} = 0.011_2 \\
 + (-0.625)_{10} = +1.010_2 \\
 \hline
 \boxed{0.110_2 \leftarrow \text{sum}}
 \end{array}$$

Carry \leftarrow

∴ Carry is zero so the sum is -ve so the sum is converted to

decimal as

$$\begin{array}{r}
 1.101_2 \xrightarrow{\substack{\text{Extract} \\ \text{Sign Bit}}} -0.101_2 \xrightarrow{\substack{\text{Complement} \\ \text{Fraction Part}}} -0.10_2 \xrightarrow{\substack{\text{Convert to} \\ \text{decimal}}} -0.25_{10}
 \end{array}$$

$$\text{Binary to decimal } (0 \times 2^{-1}) + (1 \times 2^{-2}) + (0 \times 2^{-3}) = 0.25_{10}$$

$$.10_2 = (0 \times 2^{-1}) + (1 \times 2^{-2}) + (0 \times 2^{-3}) = 0.25_{10}$$

$$\therefore (+0.375_{10}) + (-0.625_{10}) = (-0.25_{10})$$

Eq 2:- Add $+0.625$ and -0.375 by 1's Complement addition.

Step 1:- Convert decimal to binary we get

$$+0.625 = 0.101_2$$

$$-0.375 = 0.011_2$$

Step 2 Representation:

$$+0.625_{10} \xrightarrow[\text{Binary}]{\text{Convert to}} +0.101_2 \xrightarrow[\text{bit}]{\text{add sign}} 0.101_2$$

$$-0.375_{10} \xrightarrow[\text{Binary}]{\text{Convert to}} -0.011_2 \xrightarrow[\text{bit}]{\text{add sign}} 1.011_2 \xrightarrow[\text{fraction part}]{\text{Complete the}} 1.100_2$$

Step 3:- now

$$\begin{array}{r}
 & 0.101 \\
 + & 0.100 \\
 \hline
 & \text{Sum} \\
 \text{Carry} & \leftarrow 1.001 \\
 & \text{,1 add Carry to LSB} \\
 \hline
 & 0.010 \leftarrow \text{final sum}
 \end{array}$$

\because Carry is one, the sum is positive.

now final sum converted to decimal.

$$0.010 \xrightarrow[\text{sign bit}]{\text{extract}} +.010 \xrightarrow{\text{convert to decimal}} +0.25_{10}$$

$$\therefore (+0.625) + (-0.375) = (+0.25)_{10}.$$

Eq:-3 Add $+0.375$ and -0.625 by two Complement addition

Step 1:- Representation are:-

$$+0.375_{10} \xrightarrow[\text{Binary}]{\text{Convert to}} +0.011_2 \xrightarrow[\text{bit}]{\text{add sign}} 0.011_2$$

$$-0.625_{10} \xrightarrow[\text{Binary}]{\text{Convert to}} -1.01_2 \xrightarrow[\text{bit}]{\text{add sign}} 1.101_2 \xrightarrow[\text{fraction part}]{\text{Complete the}} 1.010_2 \xrightarrow{\text{add one}} 1.011_2$$

$$\begin{array}{r}
 0.011 \\
 1.011 \\
 \hline
 0.101 \leftarrow \text{sum}
 \end{array}$$

Carry.

The Carry is zero the sum is 've'

9

convert to decimal

$$1.110_2 \xrightarrow{\substack{\text{Extract} \\ \text{sign bit}}} -0.110_2 \xrightarrow{\substack{\text{Complete the} \\ \text{fraction part}}} -0.001_2 \xrightarrow{\substack{\text{add one} \\ \text{to LSB}}} -0.010_2 \xrightarrow{\substack{\text{Convert to} \\ \text{Binary}}} -0.25_{10}$$

$$\begin{array}{r} +0.375_{10} \\ -0.625_{10} \\ \hline 1.110_2 \end{array} \Rightarrow \begin{array}{l} 0.011_2 \\ 1.011_2 \\ \hline -0.25_{10} \end{array}$$

Ex:- Add $+0.625_{10}$ and -0.375_{10} by two's Complement Addition.

The 2's Complement representation is

$$\begin{array}{r} +0.625_{10} \xrightarrow{\substack{\text{Convert to} \\ \text{Binary}}} +0.101_2 \xrightarrow{\substack{\text{add sign} \\ \text{bit}}} 0.101_2 \\ -0.375_{10} \xrightarrow{\substack{\text{Convert to} \\ \text{Binary}}} -0.111_2 \xrightarrow{\substack{\text{add sign} \\ \text{bit}}} 1.011_2 \xrightarrow{\substack{\text{Complete the} \\ \text{fraction part}}} 1.100_2 \xrightarrow{\substack{\text{add one} \\ \text{to LSB}}} 1.101_2 \end{array}$$

$$\Rightarrow \begin{array}{r} 0.101_2 \\ +1.101_2 \\ \hline \end{array}$$

$$\cancel{\text{Carry}} \leftarrow \underline{0.010} \leftarrow \text{sum}$$

∴ Carry is '1' the sum is Positive. Then carry is discarded.

2's Complement addition.

$$(i.e. 0 \times 2^1 + 0 \times 2^2 + 0 \times 2^3) = 0.25_{10}$$

$$0.010_2 \xrightarrow{\substack{\text{Extract} \\ \text{sign bit}}} +0.010_2 \xrightarrow{\substack{\text{Convert to} \\ \text{Decimal}}} +0.25_{10}$$

$$\begin{array}{r} +0.625_{10} \\ -0.375_{10} \\ \hline \end{array} \Rightarrow \begin{array}{r} 0.101_2 \\ 1.101_2 \\ \hline \end{array}$$

$$(+0.625_{10}) + (-0.375_{10}) \Rightarrow \underline{0.010_2} \Rightarrow +0.25_{10}$$

Floating Point addition:-

(10)

Addition can be performed only when Exponent of both no: are equal.

i.e. Exponent of smaller no: is changed to Equal Exponent of larger no:

Consider 10-bit floating point format with

7-bit mantissa
3-bit exponent

Let LSB in mantissa and exponent be sign bit.

$$\text{let: } +5_{10} = 101_2 = 0.101000 \times 2^{+3}_{10} = 0.101000 \times 2^{+11_2}_{011_2}$$

$$+0.25_{10} \Rightarrow \cancel{+0.25_2} \quad \begin{array}{r} 0.25 \\ \times 2 \\ \hline 0.50 \\ \times 2 \\ \hline 1.00 \\ \hline .01 \end{array}$$

$$\Rightarrow +0.25 \Rightarrow \cancel{+0.25_2} + .01_2 = 0.100000 \times 2^{-1}_{10} = 0.100000 \times 2^{+10_2}$$

whose exponents are not equal. so

The small no: is $+0.25$ and so it is unnormalized to make 9 bits.

Exponent equal to that $+5$:

$$+0.25_{10} = 0.100000 \times 2^{+10_2} = 0.100000 \times 2^{\cancel{-1 \text{ unnormalize}}}_{\cancel{10_2}} = 0.000010 \times 2^3$$

now normalized mantissa of $+0.25_{10}$ is added to mantissa exponent $+5$.

$$+5_{10} \Rightarrow 0.101000 \times 2^{+11_2}$$

$$+0.25_{10} \Rightarrow \cancel{0.000010 \times 2^{+10_2}} \quad \begin{array}{r} 0.101010 \times 2^{+11_2} \\ \hline .01 \end{array}$$

$$\Rightarrow 101.01_2 = 5.25_{10}$$

dp =

3) Truncation & Rounding:

If finite word length of registered word have infinite word length of data then how to fit it in to a register

If 8-bit word and register

fit
into
register

number
is
10-bits

solved
by
using
2-methods

① Truncation

② Rounding

Method of shortening number what is effect on design parameters of DSP

If using 2-methods what is effect

10.201562387542

Neglected

↳ displayed in calculator

what is happened when neglected in number then we

Study the effect in truncation.

Then this type of truncating number stop the number required is called truncation

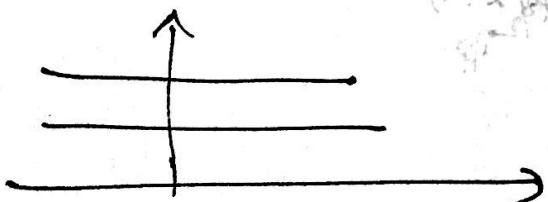
Rounding :- Eg:- 10.201562

↳ 2-ways of rounding if we keep '3' as it is then it is truncation.

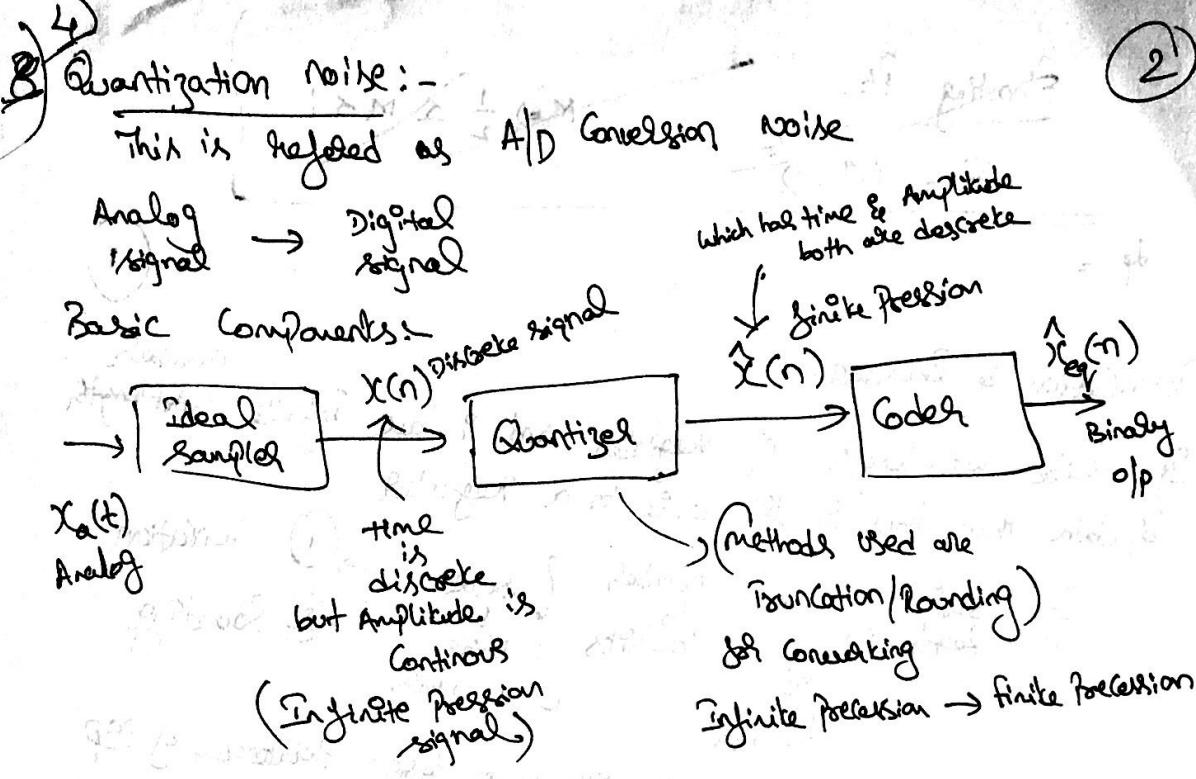
Rule :- If the number is high above '5' then it will round off to 5. If the number is less than or equal to 5 then it will round off to 4. Next digit to next higher one. 10.201562387542

=>

10.2015624



$$\frac{2}{8} + \frac{1}{4}$$



Let us consider
3-bit Quantized ($\hat{x}(n)$) 3-bit Coder :-

$2^3 = 8$ levels i.e. $b=2$
 $b+1=3$

tree numbers

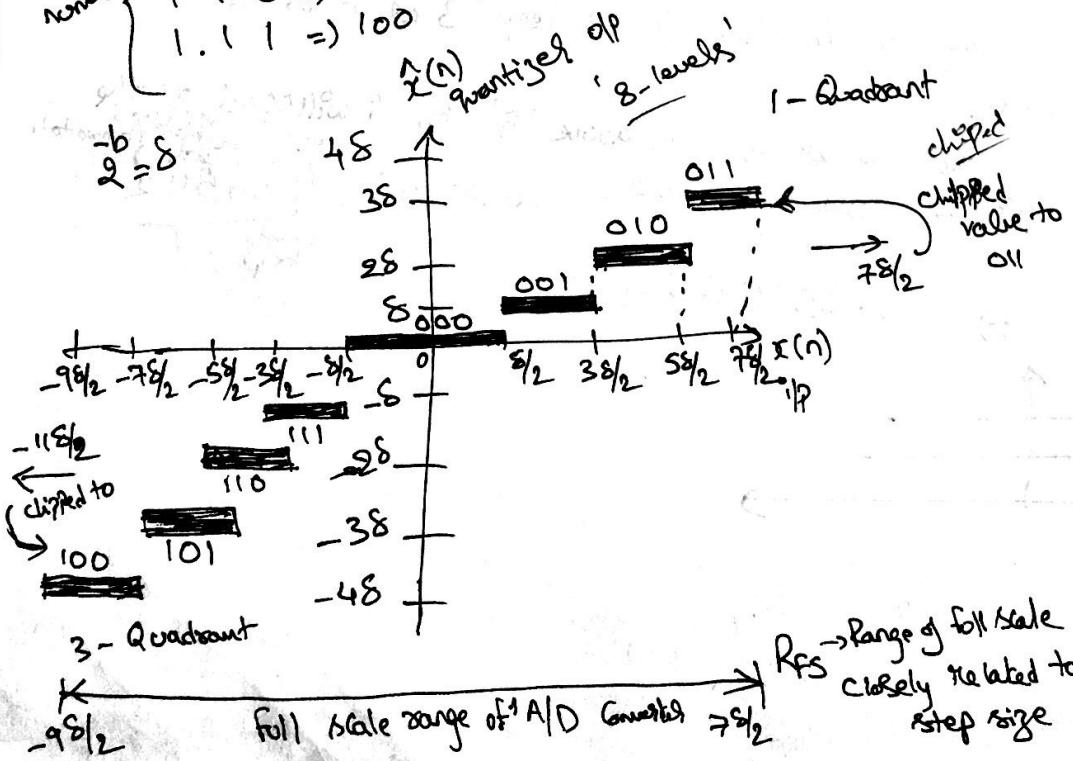
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1

jump $\frac{1}{4}$ $\frac{1}{2}$ $\frac{3}{4}$

$0.0 \rightarrow 0.25$
 $0.0 \rightarrow 0.5$ ($\frac{1}{4} + \frac{1}{2}$)
 $0.1 \rightarrow 0.75$ ($\frac{1}{4} + \frac{1}{2}$)

Complete 3-bit numbers

1.0	0.0	1.0	1.0
1.0	0.0	1.0	1.0
1.1	0	1.0	1.0
1.1	1	1.0	1.0



Using $b+1 = 3$ -bit, total no. of levels can be 8 but
for bipolar signals 3

Five no will be represented using magnitude
-ve " " " " " 2's complement

$$b = 2$$

$$2^{-5} = 2^{-2} = \frac{1}{4}$$

If %P is above $\pm 8\%$ then Quantized will clip the op to 0.11 (or) 0.1

$$\Rightarrow 0.75 \left(\frac{3}{4}\right) \Rightarrow 1 - \frac{1}{2^3} = 1 - \frac{1}{8} = \left(1 - \frac{1}{4}\right) = \frac{3}{4}$$

If α is below $-9\pi/2$ then quantized will clip off to $\int_{1/2}^{100}$ is complete

From the graph to define translation & bounding

If old value is in b/w $\frac{8}{2}$ to $\frac{38}{2}$ then the half value. $\frac{2 \cdot 2 \cdot 2^2}{2} =$

(3) middle value then decision goes
if truncation (000)

if it is down then it is subtraction
if it is up then it is rounding (00)

→ Quantization is a nonlinear process.

$$R_{fg} = \frac{b+1}{2} \cdot s \quad \leftarrow R_{fg} \text{ step size iteration}$$

$$\text{final } \hat{x}_{eq} = \frac{\hat{x}}{R_{FS}} \quad \text{i.e. } -1 < \hat{x}_{eq} < 1 \leftarrow \text{Codebook bound}$$

$\leftarrow R_{FS} \text{ Codebook relation}$

$$\text{Range is } -\left(\frac{-b}{2} + 1\right) \leq x_a(n) \leq \left(\frac{b+1}{2} - 1\right) \cdot \frac{8}{2}$$

The noise is

$$e(n) = \hat{x}(n) - x(n)$$

$$e(n) = Q(x(n)) - x(n)$$

Rounding noise is bounded by

$$-\frac{1}{2} \left(2^b - 2^{-b_1} \right) \leq F_R \leq \frac{1}{2} \left(2^{-b} - 2^{-b_1} \right)$$

where b_1 - before rounding bits

b_1 - before b - after " bits i.e. $b_1 \gg b$.

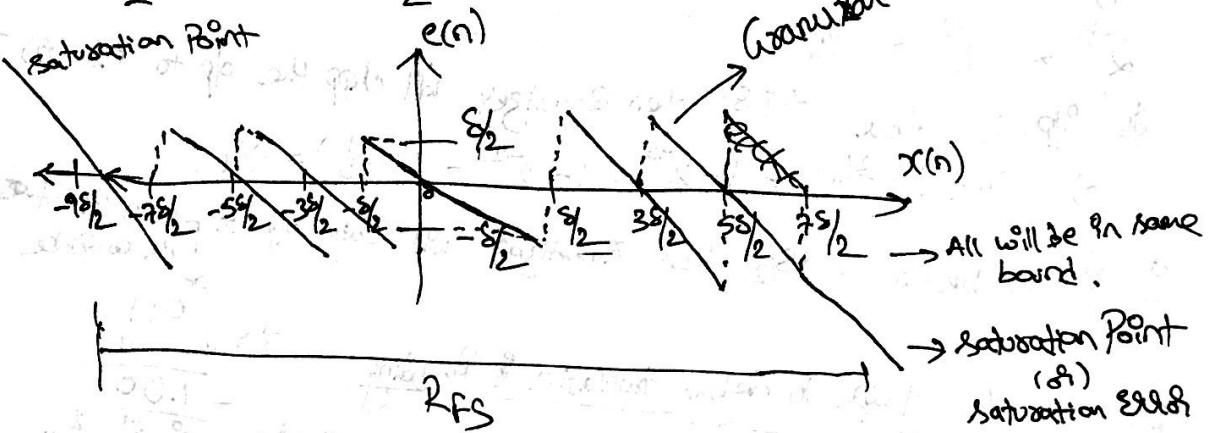
∴ The error is bounded to about $\pm \frac{\delta}{2}$ for Robot $\Rightarrow \text{Eqn. 2}$

$$-\frac{1}{2} \cdot \frac{\delta}{2} \leq e(n) \leq \frac{1}{2} \cdot \frac{\delta}{2}$$

where $\frac{\delta}{2} = \text{Span of } \Delta$

$$-\frac{\delta}{2} \leq e(n) \leq \frac{\delta}{2}$$

Saturation Point



If error is swinging b/w $-\frac{\delta}{2}$ to $\frac{\delta}{2}$ then it is guaranteed error.

If the error is going beyond the saturation point it is saturation error.

* Statistical model of error:-

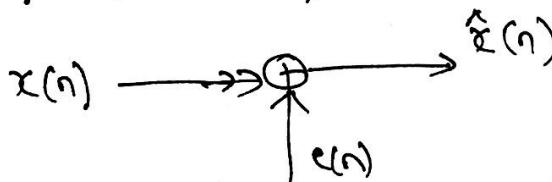
Consider error signal is random process

① $e(n) \rightarrow$ WSS white noise process

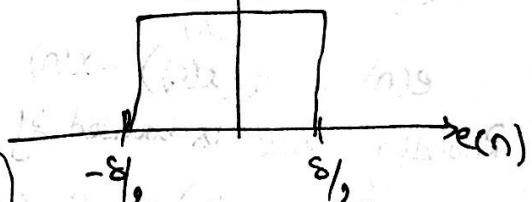
white sense stationary process
if it is uniformly distributed over 1-sample period

② $e(n) \rightarrow$ are uncorrelated to op. (i.e. error remains same)

③ $x(n) \rightarrow$ Stationary Random Process



$P(e) \rightarrow$ Probability density

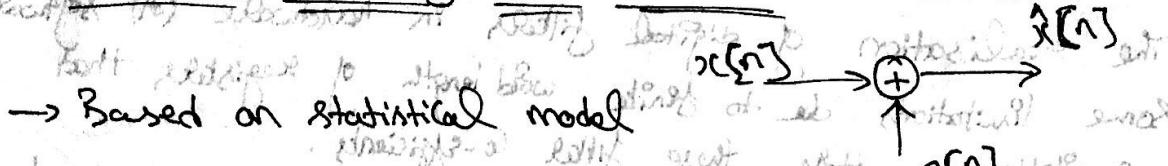


$$\text{Mean of error } (m_e) = -\frac{\delta/2 + \delta/2}{2} = 0$$

$$\text{Variance } \sigma_e^2 = \frac{((+\frac{\delta}{2}) - (-\frac{\delta}{2}))^2}{12} = \frac{\delta^2}{12} \Rightarrow \text{Power of noise.}$$

* Signal to Quantization Noise Ratio :-

(5)



→ Based on statistical model

We evaluate the effect of additive quantization noise $e[n]$ on A/D Quantizer

opposite signals $x[n]$ by computing Signal to noise Ratio (SNR) for

to Quantization noise ratio in dB defined by

$$\text{SNR}_{\text{A/D}} = 10 \log_{10} \left(\frac{\sigma_x^2}{\sigma_e^2} \right) \cdot \text{dB} \rightarrow (a)$$

σ_x^2 is signal variance representing signal power

σ_e^2 is noise variance representing quantization noise power

In case of bipolar $(b+1)$ bit A/D converted

$$S = \frac{-(b+1)}{2} R_{fs}$$

$$2^{-2b} (R_{fs})^2$$

→ (b)

$$\text{SNR}_{\text{A/D}} = 10 \log_{10} \left(\frac{48 \sigma_x^2}{2^{-2b} (R_{fs})^2} \right) \rightarrow (c)$$

Substituting eq (b) in eq (a) we get

$$\text{SNR}_{\text{A/D}} = 10 \log_{10} \left(\frac{48 \sigma_x^2}{2^{-2b} (R_{fs})^2} \right) \rightarrow (d)$$

$$\therefore \text{SNR}_{\text{A/D}} = 6.02b + 16.81 - 20 \log_{10} \left(\frac{R_{fs}}{\sigma_x} \right) \cdot \text{dB} \rightarrow (d)$$

∴ eq (d) is used to determine the minimum noise needed to meet a specified signal to quantization noise ratio

$$(\text{SNR})_{\text{A/D}} = 6.02b + 16.81 - 20 \log_{10} (K) \rightarrow (e)$$

$$\text{where } K = \frac{R_{fs}}{\sigma_x}$$

5) Finite word length effect in IIR digital filters:- (6)

The realisation of digital filters in hardware (d) software has some limitations due to finite word length of registers that are available to store these filter Co-efficients.

Co-efficients stored in these registers are either truncated/Rounded off. The system may have different frequency response not accurate the system may have different frequency response than the desired one.

There are 2 approaches for analysis and synthesis of digital filters with quantized Co-efficients.

1st approach the digital filter is assumed to be an ideal one and represented by a transfer function in parallel with ideal filter. 2nd approach, each filter is studied separately and the quantized Co-efficients can be optimised to minimise the maximum weighted difference b/w the ideal and actual frequency response.

Let the transfer function of Co-efficients quantised digital filter being realised

$$H(z) = \frac{\sum_{k=0}^M b_k z^{-k}}{\sum_{k=0}^N a_k z^{-k}} \quad \rightarrow (1)$$

If \bar{a}_k & \bar{b}_k are the Co-efficients of unquantized filter and a_k and b_k are the error quantities then the Co-efficients of quantised filter can be written as

$$a_k = \bar{a}_k + \epsilon_k$$

$$b_k = (\bar{b}_k + \beta_k) - 18.21 + 480.2 = (502)$$

7

That effect due to quantisation of filter co-efficients is

$$e(n) = y'(n) - y(n)$$

where $y(n) = \text{op of ideal filter}$

$y'(n) = \text{"actual"}$

When 2-bit numbers are multiplied the product is 23 bit long

This Product must be rounded to B bits in all digital Processing

application of finite word length multiplied can be expressed as

$$Q\{\alpha_i x(n)\} = \alpha_i x(n) + e(n) \rightarrow ②$$

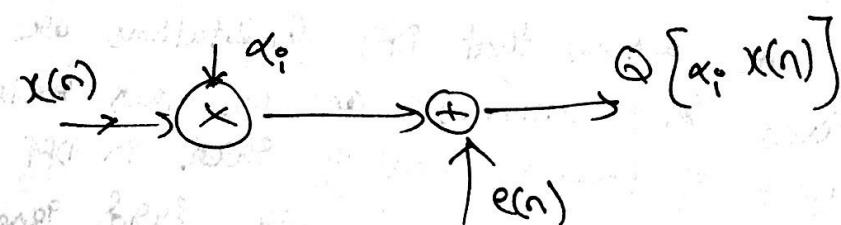


fig :- Product Quantisation noise model

where $\alpha_i x(n)$ is the product which is 2B bit long & $e(n)$ is

the error resulting from rounding the product to B-bits

The analysis of Product round off noise is similar to A/D converter

The Avg Power (Variance) is given by

$$\text{Avg Power}^2 = \frac{2^{2B}}{12} \rightarrow ③$$

6) Finite word length effect in FFT algorithms:-

(or)

Roundoff errors in FFT algorithms:-

since FFT is often employed in a number of digital signal processing applications.

and it is interest to analyze the effect of finite word lengths in FFT wordlengths in FFT computations.

The most critical error in the computation is that due to arithmetic round off errors.

As earlier sections we assume that DFT Computations are being carried out using fixed point arithmetic and we thus restrict our analysis to the effect of product round off error in DFT computation via FFT algorithm and to compute them with error generated in DFT computation.

→ If DFT computation requires K-complex multiplications, there are $4K$ sources of quantization errors in the computational structure.

Assumptions :-

- (a) All 4K errors are uncorrelated with each other and uncorrelated with input sequence.
- (b) The quantization errors are random variables uniformly distributed with a variance $\sigma_0^2 = \frac{2^{-2b}}{12}$ assuming a signed b-bit fractional fixed point arithmetic.

(a) Direct form DFT Computation:-

Recall that the N -Point DFT $X[k]$ of length N Complex sequence $x[n]$ is

$$X[k] = \sum_{n=0}^{N-1} x[n] \cdot w_N^{nk}, \quad 0 \leq k \leq N-1 \rightarrow @$$

The computation of a single DFT sample requires N -Complex multiplications and hence total no.: of real multiplications for the computation of a single DFT sample is $4N$.

The $2N$ Quantization error source. The variance of the error in computation of one DFT sample is

$$\therefore \sigma_e^2 = 4N \sigma_0^2 = \frac{2N}{3} \rightarrow b)$$

indicating op count off error is proportional to the DFT length.
Now the op sequence $x[n]$ must be scaled to avoid overflow in computation of $X[k]$

$$|X[k]| \leq \sum_{n=0}^{N-1} |x(n)| \leq N$$

assuming that op range samples satisfy the dynamic range constraint

$$|x[n]| \leq 1 \text{ to prevent overflow we need to scale}$$

$$|X[k]| < 1$$

which can be guaranteed by dividing each sample $x[n]$ in the op sequence by N .

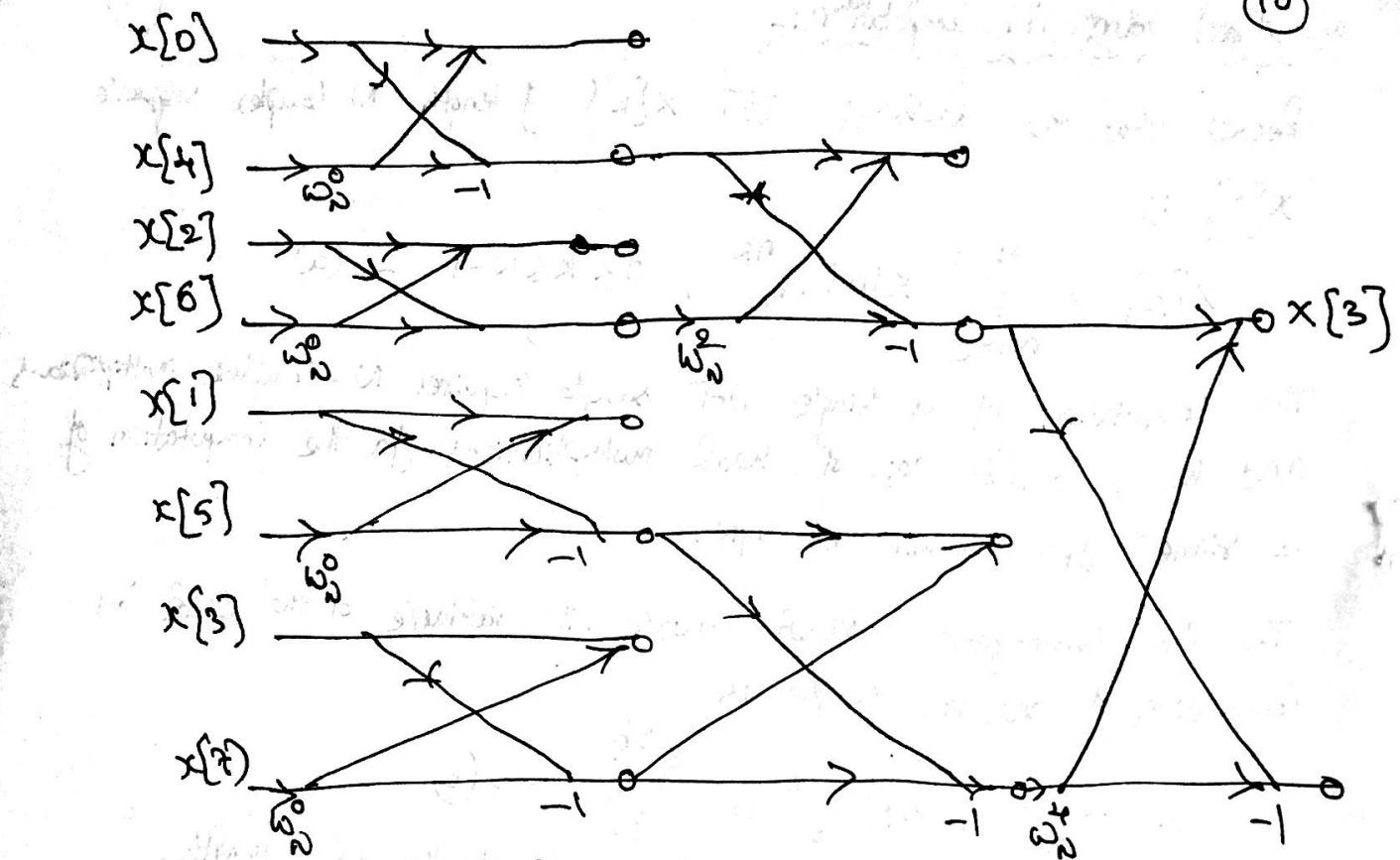
The op signal power is given by:

$$\sigma_x^2 = \frac{(2/N)^2}{12} = \frac{1}{3N^2}$$

The corresponding op signal power is

$$\sigma_x^2 = N \sigma_x^2 = \frac{1}{3N}$$

(10)



As a result signal to noise ratio is

$$SNR = \frac{\sigma_x^2}{\sigma_y^2} = \frac{2^N}{N}$$

\therefore Above expression indicates SNR has been reduced by a factor of N due to scaling and round off error.

④ DFT Computation via FFT algorithm:-

FFT is used in large no:- of applications. Hence it is necessary to analyze the effect due to finite word length in necessary to FFT algorithm.

Let σ_x^2 represents the variance of DFT Co-efficients $|X(k)|$

for N -point DFT σ_y^2 is given as

$$\sigma_x^2 = \frac{1}{3N} \quad \text{--- (1)}$$

for direct computation of DFT the variance of quantization error in multiplication is

$\sigma_q^2 = \frac{N \cdot A^2}{3}$ --- (2)

Hence σ_q^2 is variance of quantization error and

Δ is step size

$$\text{where } \Delta = 2^{-b} \rightarrow (3)$$

where b is no. of bits represents level.

∴ Eq ② becomes

$$\sigma_q^2 = \frac{N}{3} \cdot 2^{-2b} \rightarrow (4)$$

The signal to noise ratio at the op can be considered as measure of quantization errors.

The ratio of variance DFT coefficients (σ_x^2) to variance quantization error (σ_q^2) is

$$\text{Signal to noise ratio direct computation of DFT} = \left(\frac{\sigma_x^2}{\sigma_q^2} \right)_{\text{Direct DFT}}$$

from Eq ① & ② we have

$$\left(\frac{\sigma_x^2}{\sigma_q^2} \right)_{\text{direct DFT}} = \frac{\frac{1}{3N}}{\frac{N}{3} \cdot 2^{-2b}} = \frac{2^{2b}}{N^2} \rightarrow (5)$$

when DFT is computed using FFT algorithm the variance of signal remains same.

$$\sigma_x^2 = \frac{1}{3N} \text{ from Eq ①} \rightarrow (6)$$

But with algorithm the variance of the quantization error is

$$\sigma_q^2 = \frac{2}{3} \cdot 2^{-2b} \rightarrow (7)$$

Hence signal to noise ratio in FFT algorithm is

$$\left(\frac{\sigma_x^2}{\sigma_q^2} \right)_{\text{FFT}} = \frac{\frac{1}{3N}}{\frac{2}{3} \cdot 2^{-2b}} = \frac{2^{2b}}{2N} //$$