

UNIT-IV FINITE WORD LENGTH EFFECTS

Representation of numbers - ADC Quantization

noise - Coefficient Quantization error, Product

Quantization error - truncation & rounding errors,

Limit cycle due to product round-off error-

Roundoff noise power.

Introduction:

When digital systems are implemented either in hardware or in software, the filter coefficients are stored in binary registers. These registers can accommodate only a finite number of bits and hence the filter coefficients have to be truncated or rounded off in order to fit into these registers.

Truncation or rounding of the data results in degradation of system performance.

Also, in digital processing systems, a continuous time input signal is sampled and quantised

in order to get the digital signal. The process of quantisation introduces an error in the signal.

Following are errors due to finite word length effects in digital filters:

1. ADC quantisation error (due to quantization of input data by A/D convertor)
2. Coefficient quantisation error (due to quantization of filter coefficients)
3. Product quantisation error (due to rounding the product in multiplication)
4. Limit cycles due to overflow in addition.
5. limit cycles due to zero input.

Representation of numbers:

Rounding or truncation introduces an error whose magnitude depends on the number of bits truncated or rounded-off. Also, the characteristics of the error depends on the form of binary number representation.

3.

We can represent a number 'N' to any desired accuracy by a finite series,

$$N = \sum_{i=n_1}^{n_2} d_i r^i$$

where $r \rightarrow$ radix or Base

$d_i \rightarrow$ i^{th} digit of the number.

$n_2 \rightarrow$ number of integer digits.

$n_1 \rightarrow$ number of fraction digits.

When $\boxed{r=10}$, the representation is known as decimal rep. having numbers from $\underbrace{0 \text{ to } 9}_{d_i}$

$$\text{eg: } 30.285 = \sum_{i=-3}^1 d_i 10^i$$

$$= 3 \times 10^1 + 0 \times 10^0 + 2 \times 10^{-1} + 8 \times 10^{-2} + 5 \times 10^{-3}$$

When $\boxed{r=2}$, the representation is known as binary representation with two numbers $\underbrace{0 \text{ and } 1}_{d_i}$

$$\text{eg: } 110.010 = \sum_{i=-3}^2 d_i 2^i$$

$$d_i = 0 \text{ or } 1$$

$$= 1 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 0 \times 2^{-3}$$

Similarly,

when $\boxed{r=8}$, \Rightarrow octal representation

$\boxed{r=16}$, \Rightarrow hexadecimal representation.

4.

Pbm Convert the decimal number 30.275 to binary form.

Soln: Integer part:

$$\begin{array}{r}
 2 \overline{)30} \\
 2 \overline{)15} - 0 \\
 2 \overline{)7} - 1 \\
 2 \overline{)3} - 1 \\
 1 - 1
 \end{array}
 \quad (30)_{10} \Rightarrow (11110)_2$$

Fractional part:

$$0.275 \times 2 = 0.550 \rightarrow 0$$

$$0.550 \times 2 = 1.10 \rightarrow 1$$

$$0.10 \times 2 = 0.2 \rightarrow 0$$

$$0.2 \times 2 = 0.4 \rightarrow 0$$

$$0.4 \times 2 = 0.8 \rightarrow 0$$

$$0.8 \times 2 = 1.6 \rightarrow 1$$

$$0.6 \times 2 = 1.2 \rightarrow 1$$

$$0.2 \times 2 = 0.4 \rightarrow 0$$

$$(0.275)_{10} \Rightarrow (01000110\ldots)_2$$

$$\text{Ans: } (30.275)_{10} \Rightarrow (11110.01000110\ldots)_2$$

Types of Number Representation:

1. Fixed Point Representation
2. Floating " "
3. Block floating " "

Fixed point Representation:

In fixed point arithmetic, the position of the binary point is fixed. The bit to the right represent the fractional part of the number and those to the left represent the integer part.

for eg: 01.1100 has the value 1.75 in decimal.

there is only one unique way of representing binary fraction positive number,

$$\text{i.e., } N_p = \overset{d_0}{\underset{\downarrow}{0}}.d_1 d_2 d_3 \dots d_{n-1} \\ = (0 \times 2^0) + (d_1 \times 2^{-1}) + (d_2 \times 2^{-2}) + \dots + (d_{n-1} \times 2^{-n+1})$$

where $d_0 = 0 \Rightarrow$ sign bit (=0 for positive number)
(MSB Bit)

there are 3 different formats for representing negative binary numbers.

1. Sign-magnitude format
2. One's Complement "
3. Two's Complement "

$$N_p = (0 \times 2^0) + \sum_{i=1}^{n-1} d_i \cdot 2^{-i}$$

6.

Sign-Magnitude format:

In this format, the Most Significant Bit (MSB) bit is set to '1' to represent the negative sign.

For eg: $+1.75 \Rightarrow (01.1100)_2$
 $-1.75 \Rightarrow (11.1100)_2$

In this format, the number '0' has two representations i.e., 00.000000_2 (or) 10.000000_2 .

with 'b' bits, only $(2^b - 1)$ numbers can be represented.

$$N_n = (1 \times 2^0) + \sum_{i=1}^n d_i 2^{-i}$$

Pbms: Express the fraction $+7/8$, $-7/8$, $+6/8$, $-6/8$ in sign-magnitude format.

Ans: $(7/8)_{10} = (0.875)_{10}$

$$\begin{aligned} (+0.875)_{10} &\Rightarrow 0.875 \times 2 = \underbrace{1.750}_{\downarrow} \rightarrow 1 \\ &0.750 \times 2 = \underbrace{1.500}_{\downarrow} \rightarrow 1 \\ &0.5 \times 2 = \underbrace{1.0}_{\downarrow} \rightarrow 1 \end{aligned} \quad \Rightarrow (0.111)_2$$

*↓
Sign bit*

$$(-0.875)_{10} \Rightarrow (1.111)_2$$

*↓
Sign bit*

Sign bit = $\begin{cases} 0 & \text{for +ve no.} \\ 1 & \text{for -ve no.} \end{cases}$

$$(6/8) \Rightarrow (0.75)_{10}$$

$$\begin{aligned} (+6/8)_{10} &\Rightarrow (0.75)_{10} = 0.75 \times 2 = \underbrace{1.500}_{\downarrow} \rightarrow 1 \\ &0.5 \times 2 = \underbrace{1.0}_{\downarrow} \rightarrow 1 \end{aligned} \quad \Rightarrow (0.110)_2$$

$$(-6/8)_{10} \Rightarrow (1.110)_2$$

7.

One's complement format:

In this format, positive number is represented as in the sign-magnitude notation. But the negative number is obtained by complementing all the bits of the positive number.

Ex: $(0.875)_{10} = (0.111000)_2$

$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$

$$(-0.875)_{10} = (1.000111)_2$$

$$N_n = (1 \times 2^0) + \sum_{i=1}^n (1-d_i) 2^{-i}$$

here the number '0' can be represented as 0.000000 and 1.111111 .

with 'b' bits, $(2^b - 1)$ numbers can be represented exactly.

Ex. Pblm: Express the fraction $+7/8$, $-7/8$, $+6/8$, $-6/8$ in one's magnitude format.

Soln: $(+7/8)_{10} = (0.875)_{10} = (0.111)_2$

$$(-7/8)_{10} = (1.000)_2 \leftarrow \text{complement of all bits}$$

$$(+6/8)_{10} = (0.110)_2$$

$$(-6/8)_{10} = (1.001)_2$$

8.

Two's Complement format:

In this format, positive numbers are represented as in sign-magnitude and one's complement form.

The negative number is obtained by complementing all the bits of the positive number and adding one to LSB.

$$N_n = (1 \times 2^0) + \sum_{i=1}^n (1 - d_i) 2^{-i} + (1 \times 2^{-n})$$

for e.g: $(+0.875)_{10} \Rightarrow (0.111\ 000)_2$

$$\begin{array}{r} 1.000\ 111 \\ \hline 1 \\ \hline \end{array} \leftarrow \text{ones complement}$$

$$\leftarrow \text{add 1 to LSB}$$

$$(-0.875)_{10} \leftarrow \underline{\underline{(1.001\ 000)}_2}$$

Pbm
ef: Express the fraction $+7/8, -7/8, +6/8, -6/8$ in 2's complement format with 4 bit representation including sign bit.

Soln: $(+7/8)_{10} = (0.875)_{10} = (0.111)_2$

$$\begin{array}{r} 1.000 \\ \hline 1 \\ \hline \end{array}$$

$$(-7/8)_{10} = (1.001)_2$$

$$(+6/8)_{10} = (0.75)_{10} = (0.110)_2$$

$$\begin{array}{r} 1.001 \\ \hline 1 \\ \hline \end{array}$$

$$(-6/8) = (1.010)_2$$

$$\begin{array}{r} 1 \\ \hline 1.010 \\ \hline \end{array}$$

9.

Comparison of 3 formats for 4-bit word length (including sign-bit)

Decimal Equivalent	Binary no. in fixed point representation	1's complement	2's complement
	Sign-magnitude		
0	[0.000]	[0.000]	0.000
1/8	0.001	0.001	0.001
2/8	0.010	0.010	0.010
3/8	0.011	0.011	0.011
4/8	0.100	0.100	0.100
5/8	0.101	0.101	0.101
6/8	0.110	0.110	0.110
7/8	0.111	0.111	0.111
-0	[1.000]	[1.111]	X
-1/8	1.001	1.110	1.111
-2/8	1.010	1.101	1.110
-3/8	1.011	1.100	1.101
-4/8	1.100	1.011	1.100
-5/8	1.101	1.010	1.011
-6/8	1.110	1.001	1.010
-7/8	1.111	1.000	1.001
-8/8			1.000

b=4

$$\text{total: } 2^4 - 1$$

$$= 15 \text{ dup.}$$

$$= 2^4 - 1$$

$$= 15 \text{ rep}$$

$= 2^4$ representations

Addition of two fixed point numbers:

two numbers are added bit by bit starting from right, with carry bit being added to next bit.

eg: $(0.5)_{10} + (0.125)_{10} = ?$ (Assume 4 bits including sign bit)

$$(0.5)_{10} = 0.100_2$$

$$\begin{array}{r} (0.125)_{10} = 0.001_2 \\ \hline 0.101_2 \end{array} \Rightarrow (0.625)_{10}$$

↑
Sign bit

Binary to decimal conversion:

$$\begin{array}{r} 0.101 \\ \downarrow \quad \downarrow \quad \downarrow \\ 1 \times 2^3 = 0.125 \\ 0 \times 2^{-2} = 0.000 \\ 1 \times 2^{-1} = 0.500 \\ 0 \times 2^0 = 0.000 \\ \hline 0.625 \\ \downarrow \\ \text{Sign bit} \end{array}$$

When two numbers are added, the sum cannot be represented by 'b' bits, an overflow is said to occur.

for eg: $(0.5)_{10} + (0.625)_{10} = ?$

$$(0.5)_{10} = 0.100_2$$

$$\begin{array}{r} (0.625)_{10} = 0.101_2 \\ \hline 1.001_2 \end{array} \Rightarrow (-0.125)_{10} \text{ in sign magnitude format}$$

↑
Sign bit

Binary in decimal:

$$\begin{array}{r} 1.001_2 \Rightarrow 1 \times 2^{-3} = 0.125 \\ 0 \times 2^{-2} = 0.000 \\ 0 \times 2^{-1} = 0.000 \\ 1 \times 2^0 = 1.000 \\ \hline 1.125_{10} \end{array}$$

overflow occurs. it cannot be represented in 'b' bit number system.

Subtraction of 2 fixed point numbers:

Subtraction is done by addition both in one's and 2's complement format.

procedure for 1's complement addition:

- Steps:
1. find 1's complement of negative number
 2. Add positive number with 1's complement of negative number.
 3. If carry is generated, it is added to the LSB of addition result, If no carry, Answer will be in 1's complement form.

Ex: $(0.625)_{10} - (0.375)_{10} = ?$

$$(0.625)_{10} + (-0.375)_{10} = ?$$

$$(0.625)_{10} \Rightarrow 0.101_2$$

$$(0.375)_{10} \Rightarrow 0.011_2$$

$$\text{1's complement of } (0.375)_{10} \Rightarrow 1.100_2$$

$$(0.625)_{10} \underline{\quad} 0.101_2 (+)$$

$$\text{carry} \leftarrow \boxed{1} \dots \rightarrow 1$$

$$\underline{\quad 0.10_2} \Rightarrow (0.250)_{10}$$

Note: Subtraction of smaller number from larger number results in 'carry = 1', \therefore add carry to LSB to get original result.

Ex 12.

Ex: $(0.375)_{10} - (0.625)_{10} = ?$

$(0.375)_{10} + (-0.625)_{10} = ?$

\Downarrow
 0.011_2

\Downarrow
 0.101_2 complement
 1.010_2

$$(0.375)_{10} = 0.011_2$$

$$(-0.625)_{10} = \underline{1.010_2}$$

$$\begin{array}{r} & 1.101_2 \\ \underline{-} & 1.010_2 \\ \hline & 0.011_2 \end{array} \xrightarrow{\text{result in } 2\text{'s complement form.}}$$

$$\text{no carry} \quad 0.011_2 \Rightarrow (-0.250)_{10} \leftarrow \text{original result.}$$

Carry = 0.

Note: Subtraction of larger number from smaller number results no carry. Original result is obtained by complementing the addition result except sign bit.

Procedure for 2's complement addition:

- Steps: 1. find 2's complement of negative number.
2. Add the positive number with 2's complement of negative number.
3. If ~~yes~~ carry is generated, omit the carry to get original result. If no carry, answer will be in 2's complement form. Hence find 2's complement of addition result to get original result except sign bit.

Ex: $(0.625)_{10} - (0.375)_{10} = ?$

$(0.625)_{10} + (-0.375)_{10} = ?$

$$\begin{array}{r} 0.101_2 \\ \Downarrow \\ \boxed{0.010} \end{array}$$

$$0.011$$

$1 \cdot 100 \Rightarrow 1's \text{ complement}$

$1 \Rightarrow \text{add } 1 \text{ to LSB}$

carry omitted

Answer $\Rightarrow \text{rest } (0.010)_2 \Rightarrow (0.250)_{10}$

Ex: $(0.375)_{10} - (0.625)_{10} = ?$

$(0.375)_{10} + (-0.625)_{10} = ?$

$$\begin{array}{r} 0.011 \\ \Downarrow \\ \boxed{1.110} \end{array}$$

$$0.101$$

$$\begin{array}{r} 1.010 \\ \Downarrow \\ 1.011_2 \end{array}$$

no carry $001 \Rightarrow 1's \text{ complement}$

$1 \Rightarrow \text{add } 1 \text{ to LSB}$

$1.010 \Rightarrow \text{original result: } (1.010)_2 \Rightarrow (-0.250)_{10}$

Subtraction of large number from smaller number results in "no carry"

Multiplication in fixed point arithmetic:

In multiplication of 2 fixed point numbers,

first sign and magnitude components are separated.

The magnitude of the numbers are multiplied first, then sign of the product is determined and applied to the result. When 'b' bit number is

14.

multiplied with another 'b' bit number, the product may contain ' $2b$ ' bits.

 $b=2$

$$\text{eg: } (11)_2 \times (11)_2 = (1001)_2$$

$$\begin{array}{r}
 & 1 & 1 \\
 & 1 & 1 \\
 \hline
 & 1 & 1 \\
 & 1 & 1 \\
 \hline
 1 & 0 & 0 & 1 \leftarrow 2b \text{ bits}
 \end{array}$$

Note: In fixed point arithmetic, multiplication of 2 fractions result in a fraction.

For multiplication with fractions overflow can never occur.

$$\text{eg: } 0.1001 \times 0.0011 = 0.00011011$$

4 bits 4 bits 8 bits.

$$\begin{array}{r}
 1 & 0 & 0 & 1 \\
 0 & 0 & 1 & 1 \\
 \hline
 1 & 0 & 0 & 1 \\
 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 \\
 \hline
 0 & 0 & 1 & 1 & 0 & 1 & 1
 \end{array}$$

Floating Point Representation:

In floating point representation, a positive number is represented as

$$F = M \times 2^e$$

where $M \rightarrow$ mantissa, is a fraction $[0.5 \leq M \leq 1.0]$

$e \rightarrow$ exponent can be either +ve or -ve.

$$\text{eg: } (4.5)_{10} \Rightarrow 2 \begin{array}{r} 4 \\ \hline 2 \end{array} \begin{array}{l} 2-0 \\ \uparrow \\ 1-0 \end{array} \quad 0.5 \times 2 = \begin{array}{r} 1.0 \\ \uparrow \end{array} \begin{array}{l} 1 \\ 0 \end{array} \downarrow$$

$$\Rightarrow (100.10\ldots)_2$$

$$\Rightarrow 0.1001 \times 2^3 = \boxed{0.1001 \times 2^{011}}$$

$$(0.625)_{10} \Rightarrow 0.1010 \times 2^0 = \boxed{0.1010 \times 2^{000}}$$

$$(6.5)_{10} \Rightarrow 2 \begin{array}{r} 6 \\ \hline 2 \end{array} \begin{array}{l} 2-0 \\ \uparrow \\ 1-1 \end{array} \quad 0.5 \times 2 = \begin{array}{r} 1.0 \\ \uparrow \end{array} \begin{array}{l} 1 \\ 0 \end{array} \downarrow$$

$$= (110.10\ldots)_2$$

$$= 0.1101 \times 2^3 = \boxed{0.1101 \times 2^{011}}$$

$$(1.5)_{10} = (1.100\ldots)_2 = 0.1100 \times 2^1$$

$$= \boxed{0.1100 \times 2^{001}}$$

Negative floating point numbers are generalised represented by considering mantissa as a fixed point number. The sign of the floating point number is obtained from the first bit of mantissa.

Multiplication of 2 floating point numbers:

$$F_1 = 2^{e_1} \times M_1, \quad F_2 = 2^{e_2} \times M_2$$

$$F_3 = F_1 \times F_2 = (\underbrace{M_1 \times M_2}_2) 2^{(e_1+e_2)}$$

must be: 0.25 to 1.0

In decimal

$$425.3$$

$$= 0.4253 \times 10^3$$

$$= 42530 \times 10^{-2}$$

16.

$$\text{Q: } (1.5)_{10} \times (1.25)_{10} = ?$$

$$(1.5)_{10} = 0.75 \times 2^1 = 0.1100 \times 2^{001}$$

$$(1.25)_{10} = 0.625 \times 2^1 = 0.1010 \times 2^{001}$$

$$0.1100 \times 0.1010 = 0.\underbrace{0111000}_{8 \text{ bits}}$$

$$2^{001} \times 2^{001} \Rightarrow 2^{010} \leftarrow (001+001)$$

$$\text{Ans: } 0.0111000 \times 2^{010}$$

$$\begin{array}{r}
 1010 \\
 1100 \\
 \hline
 0000 \\
 0000 \\
 1010 \\
 \hline
 1010 \\
 \hline
 1111000
 \end{array}$$

Addition and subtraction of 2 floating point nos:

⇒ more difficult than that of fixed point nos.

To carry out addition, first adjust the exponent of the smaller number until it matches the exponent of larger number. The mantissa are then added or subtracted. Finally, the resulting representation is rescaled so that its mantissa lies in the range 0.5 to 1.0

$$\text{Q: } (3.0)_{10} + (0.125)_{10} = ?$$

\downarrow

$$\begin{aligned}
 & 0.110000 \times 2^{010} \\
 & \quad \downarrow \\
 & 0.001000 \times 2^{000} \\
 & = 0.0000100 \times 2^{010} \\
 & \underline{0.1100000 \times 2^{010}} \\
 & \hline
 & 0.1100100 \times 2^{010}
 \end{aligned}$$

Answer:

Comparison of fixed and floating point arithmetic:

fixed	floating
1. Fast operation	1. slow operation
2. small dynamic range	2. large dynamic range.
3. Round off errors occur only for addition	3. occurs for both addition & multiplication
4. overflow occurs in addition	4. overflow does not arise.
5. hardware implementation is cheaper	5. costlier
6. preferred for real time applications	6. non-real time applications
7. used in small computers	7. used in larger, general purpose computers.

Block floating point numbers:

A compromise between fixed and floating point systems is the block-floating point arithmetic. Here, the set of signals to be handled is divided into blocks. Each block have the same value for exponent.

The arithmetic operations within the block uses fixed point arithmetic and only one exponent per block is stored, thus saving memory. This type of representation is most suitable in certain FFT flowgraphs and in digital audio applications.