

# Unit 5

## Quantization Errors in DSP

## FINITE WORDLENGTH EFFECTS

- Registers are basic storage device in digital system. The maximum size of the binary information that can be stored in a register is called **Register word length**.
- For storing the input data in registers they have to be quantized, to be coded in binary and it depends on register word length.
- This **quantization and coding** will introduce error in the input data as the analog data (**infinite precision**) is converted into digital data (**Finite Precision**).
- While performing computation, the results may exceed the size of the register used for storing the result. Then results should be truncated or rounded off to accommodate which makes the system non-linear and leads to limit cycle behaviour. The effects of **Truncation & rounding** can be represented in terms of additive error signal, which is called as **round-off noise**.

## FINITE WORDLENGTH EFFECTS

Effects due to finite precision representation of number in a digital system are commonly referred to as **Finite Word Length effects**. Some of the finite word length effects in digital filters are

- Errors due to quantization of input data by A/D converter
- Errors due to quantization of filter coefficients.
- Errors due to rounding the product in multiplication.
- Errors due to overflow in addition.
- Limit Cycles.

- \* Digital Signal Processing algorithms are mainly used in Digital Systems like digital computers.
- \* In Digital Systems, the numbers and coefficients are stored in finite length registers.
- \* Therefore, the numbers are quantized by
  - \* truncation (or)
  - \* rounding off

## Coefficient Quantization error

**Truncation**- Simply cutoff the remaining bits

0.5485

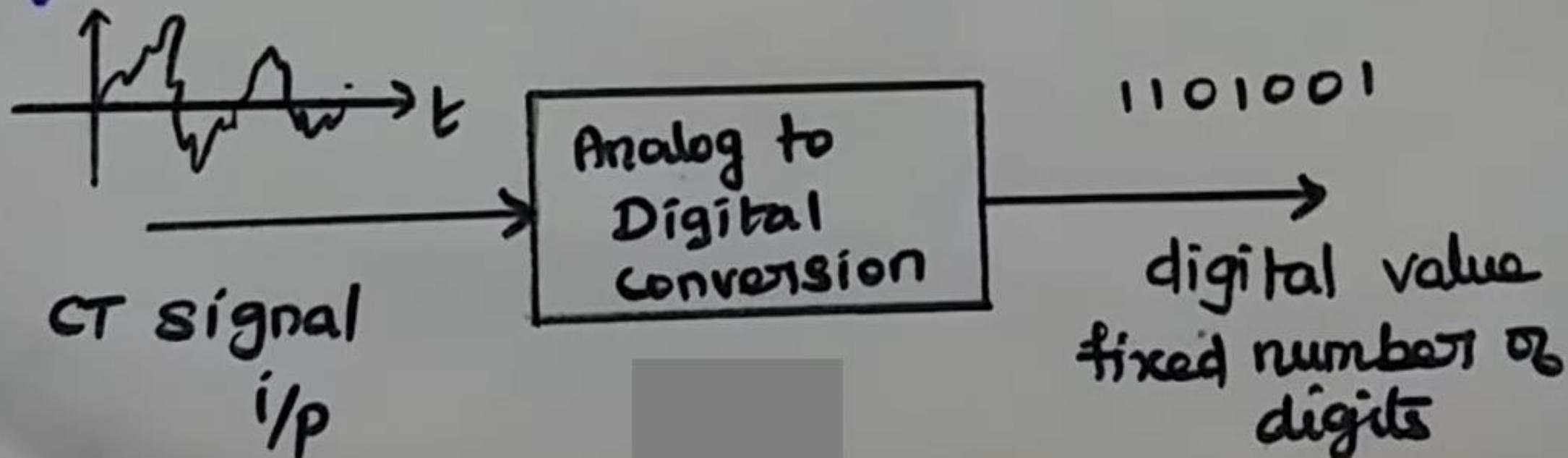
Rounding off -0.7878  
-0.788

## Effects / Errors due to quantization :-

- \* Input quantization errors
- \* Product quantization errors
- \* Coefficient quantization errors
- \* Round off noise &
- \* Limit cycle oscillations

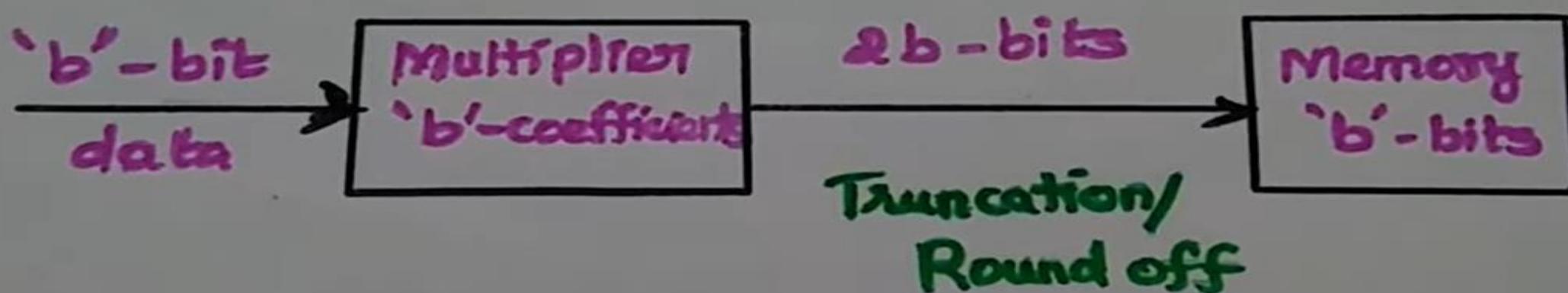
## Input quantization error :-

→ Input quantization error occurs when the continuous time signal is converted into digital value in A/D conversion process.



## \* Product Quantization Error

→ Product Quantization error occurs at the output of the multiplier.



- Since 'b'- bit register is used , the multiplier output must be rounded (or) truncated to 'b' bits .
- It produces a product error .

## \* Coefficient Quantization Error :-

→ Coefficient quantization error arises due to the process of quantizing the filter coefficients.

→ If the filter coefficients are quantized, the frequency response of the resulting filter may differ from the desired response.

## \* Limit cycle oscillations :-

→ It is a low level oscillation which arises due to the non-linearity associated with rounding off the internal filter coefficients.

\* Zero input limit cycle oscillation

\* Overflow limit cycle oscillation

## Number Representation



- In DSP, a number 'N' can be represented to any desired format using Number system.
- To represent the numbers in any digital hardware.

## Types of Number Representation:-

### (i) Fixed Point Representation

- \* Sign magnitude form
- \* 1's complement form
- \* 2's Complement form

### (ii) Floating Point Representation

### (iii) Block Floating Point Representation

## Fixed Point Representation :-

→ The position of the binary point is fixed.

Ex:-

In Binary,

11.01011  
Integer      ↑      Fractional Part  
Part              Point

In Decimal,

3.34375

# Representation of Negative Numbers

## in Fixed - point Arithmetic

- \* Sign - Magnitude form
- \* One's complement form
- \* Two's complement form

## Sign-Magnitude form:-

→ The most significant bit (MSB) is set to '1' → to represent the negative sign.  
'0' → to represent the positive sign

Ex:-

$$-1.25 \Rightarrow 11.01$$

$$+1.25 \Rightarrow 01.01$$

## One's Complement form :-

→ In this method, the negative number is obtained by complementing all the bits.

Ex:-

$$(0.875)_{10} \Rightarrow (0.111000)_2$$

$$(-0.875)_{10} \Rightarrow (1.000111)_2$$

## Two's Complement form :-

→ The negative number can be obtained by complementing all the bits and adding one to the least significant bit.

Ex:-

$$(0.875)_{10} \Rightarrow (0.111000)_2$$

↓ ↓↓↓↓↓

complementing  
each  
bit

1.000111

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

Adding  
one bit

$$(-0.875)_{10} \Rightarrow \underline{\underline{1.0010\ 00}}$$

## \* Floating Point Representation :-

→ A positive number is represented as,

$$F = 2^C \cdot M$$

where,

M → Mantissa →  $\frac{1}{2} \leq M < 1$

C → Exponent → either positive or negative

Ex:-

$$2.25 = 2^2 \times 0.25$$

$$F = 2^{010} \times 0.01$$

$$0.125 = 2^0 \times 0.001$$

$$F = 2^{000} \times 0.001$$

→ Negative floating point number is represented by considering the mantissa as a fixed point number.

→ The first bit of Mantissa represents the sign of the floating point number.

## \* Block Floating point representation:-

- the combination of fixed point and floating point representations.
- The set of signals is divided into blocks.
- The arithmetic operations within the block uses fixed point arithmetic and only one exponent per block is stored.

## Fixed Point Representation

\* Fast & inexpensive implementation

\* Limited Dynamic range

\* Round off errors occur only for addition

## Floating Point Representation

Slow & Expensive implementation.

Large Dynamic range

Round off errors can occur with both addition & multiplication

\* Overflow occurs in addition process

Overflow does not occur.

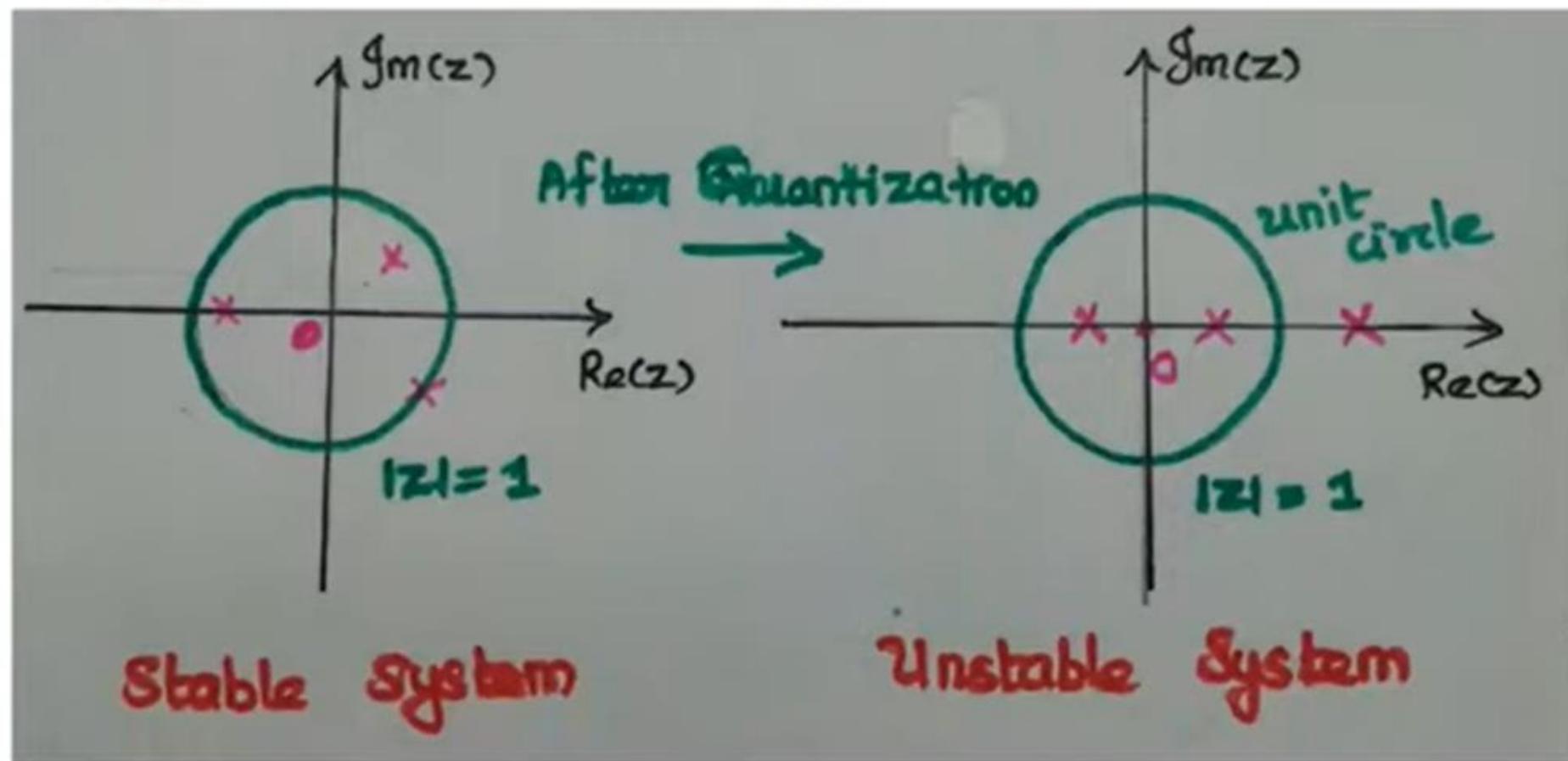
\* Low power consumption

High power consumption

\* Less flexible

More flexible

# Coefficient Quantization



**It changes the pole locations**

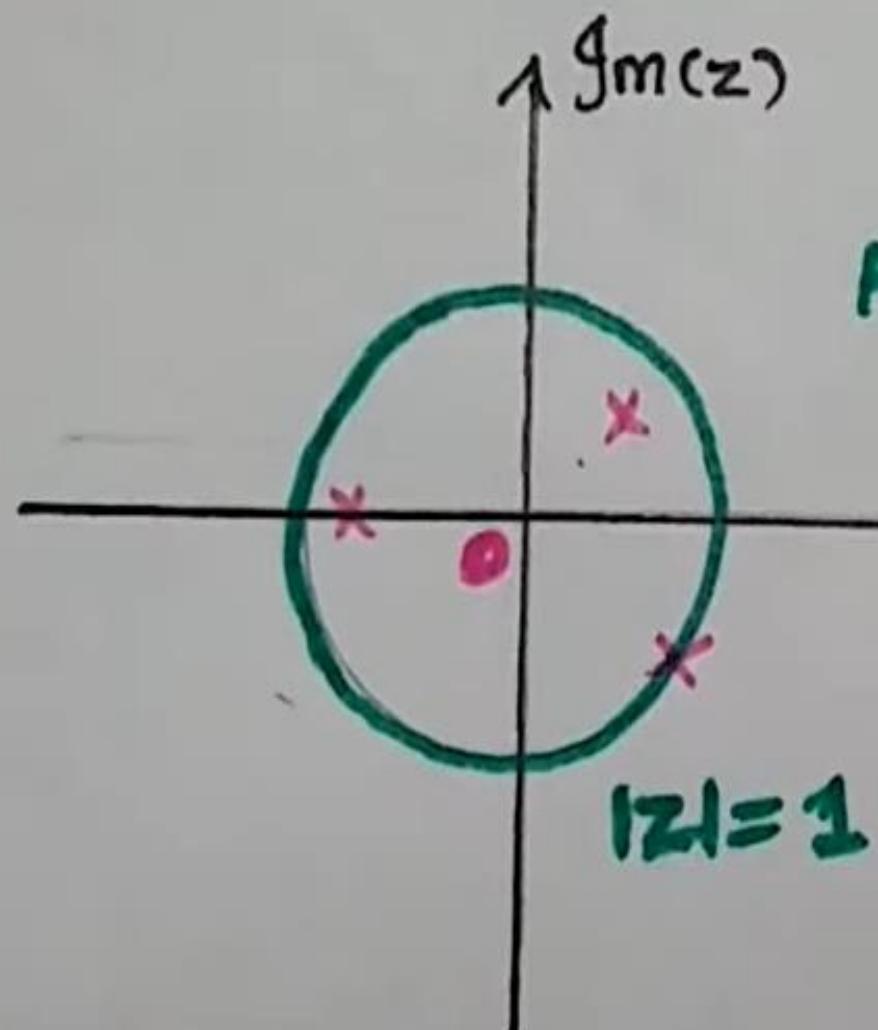
## \* Explanation :-

- Theoretically, the filter coefficients are computed to infinite precision.
- But in digital system, the filter coefficients are represented in binary and are stored in registers.
- Hence the coefficients must be quantized using rounding / truncation method

## \* Effects of coefficient Quantization :-

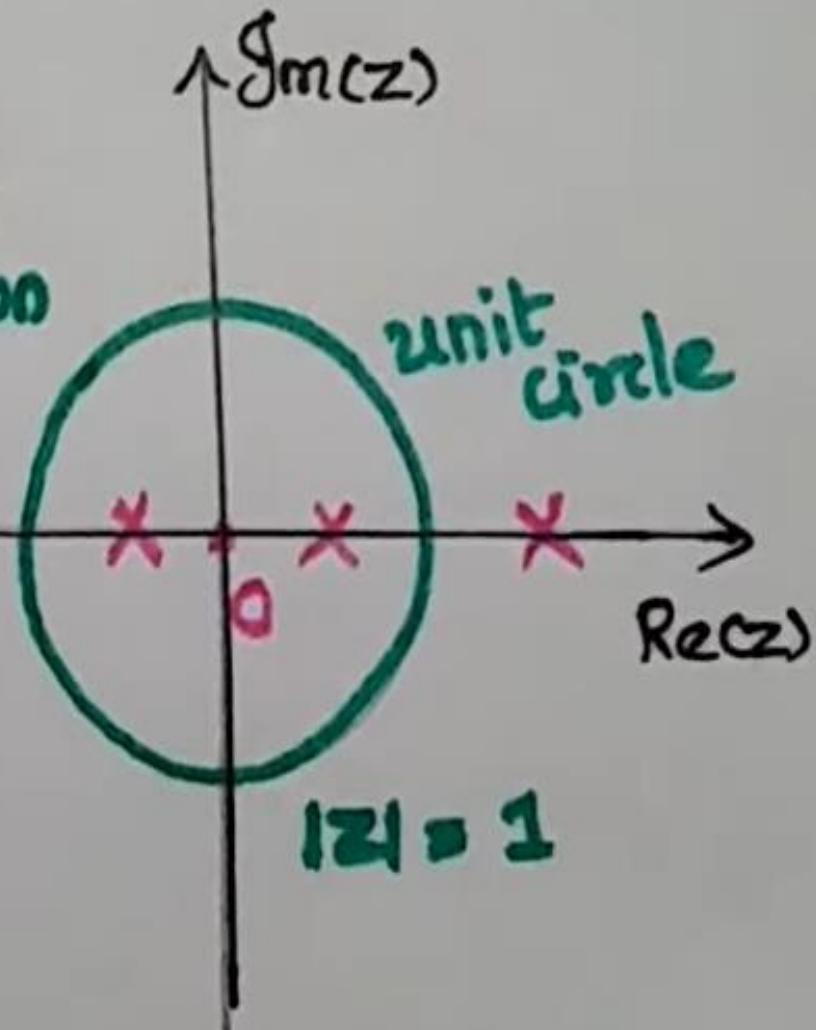
---

- Due to the coefficient quantization ,  
the frequency response of the filter  
deviates from the desired response of  
the system.
- Filter may fail to meet the desired  
specifications .



Stable System

**After Quantization**



Unstable System

- \* For stable system, poles are inside the unit circle.
- \* Quantized coefficients may lie outside the unit circle of the z-plane leading to instability.
- \* Method used to reduce coefficient Quantization
  - A high -order filter as a cascade of second -order sections.

Pbm:- Find the effect of coefficient quantization on pole locations for the given second order IIR filter using Direct form - I

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

Solution:-

Take 3 bits for truncation.

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

$$= \frac{1}{1-0.45z^{-1}-0.5z^{-1}+0.225z^{-2}}$$

$$H(z) = \frac{1}{1-0.95z^{-1}+0.225z^{-2}}$$

## Coefficient Quantization:-

$$H(z) = \frac{1}{1 - Q[0.95z^{-1}] + Q[0.225z^2]}$$

$Q[0.95]$  :-

$$\frac{0.95}{2}$$

$$\frac{\boxed{0} \cdot 9 \ 0}{2} - 1$$

$$\frac{\boxed{0} \cdot 8 \ 0}{2} - 1$$

$$\frac{1.6 \ 0}{2} - 1$$

$$\frac{2}{1.2 \ 0} \rightarrow 1$$

$Q[0.95]_{10} \Rightarrow Q[0.1111\cdots]_2$   
 $\Rightarrow (0.111)_2$

$Q[0.95]_{10} = (0.875)_{10}$

$Q[0.225]_{10}$

$$\begin{array}{r} 0.225 \\ \times 2 \\ \hline 0.450 \rightarrow 0 \\ \hline 0.900 \rightarrow 0 \\ \hline 1.800 \rightarrow 1 \\ \hline 0.600 \rightarrow 1 \\ \hline 1.200 \rightarrow 1 \end{array}$$

$Q[0.225]_{10}$

$$\Rightarrow [0.\underbrace{001}_{\frac{1}{8}}\underbrace{11\dots}_2]_2$$

$$\Rightarrow [0.001]_2$$

$$0.001_2 = \frac{1}{8} = 0.125$$

$$Q[0.225]_{10} \Rightarrow (0.001)_2$$

$$= (0.125)_{10}$$

$$Q[0.225]_{10} \Rightarrow (0.125)_{10}$$

$$H(z) = \frac{1}{1 - 0.875 z^{-1} + 0.125 z^{-2}}$$

Pole locations are changed to.

$$0.95 \rightarrow 0.875$$

$$0.225 \rightarrow 0.125$$

## Quantization of filter coefficients in cascade method

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

Solution :-

No. of bits = 3 ; Truncation method.

$Q(0.5)$  :-

$$\frac{0.5}{2} \\ \frac{2}{1.0} \rightarrow 1$$

$$Q[0.5]_{10} = Q[0.100]_2$$

$$Q[0.5] = 0.5$$

Q[0.45] :-

0.45

2  
0.90 → 0

2

1.80 → 1

2

1.60 → 1

2

1.20 → 1

$$Q[0.45]_{10} = Q[0.\underline{01}\underline{11}\dots]_2$$

$$= (0.011)_2$$

$$\Rightarrow 0 \cdot \left(0 + \frac{1}{4} + \frac{1}{8}\right)$$

$$\Rightarrow 0 \cdot \left(\frac{3}{8}\right)$$

$Q[0.45] = 0.375$

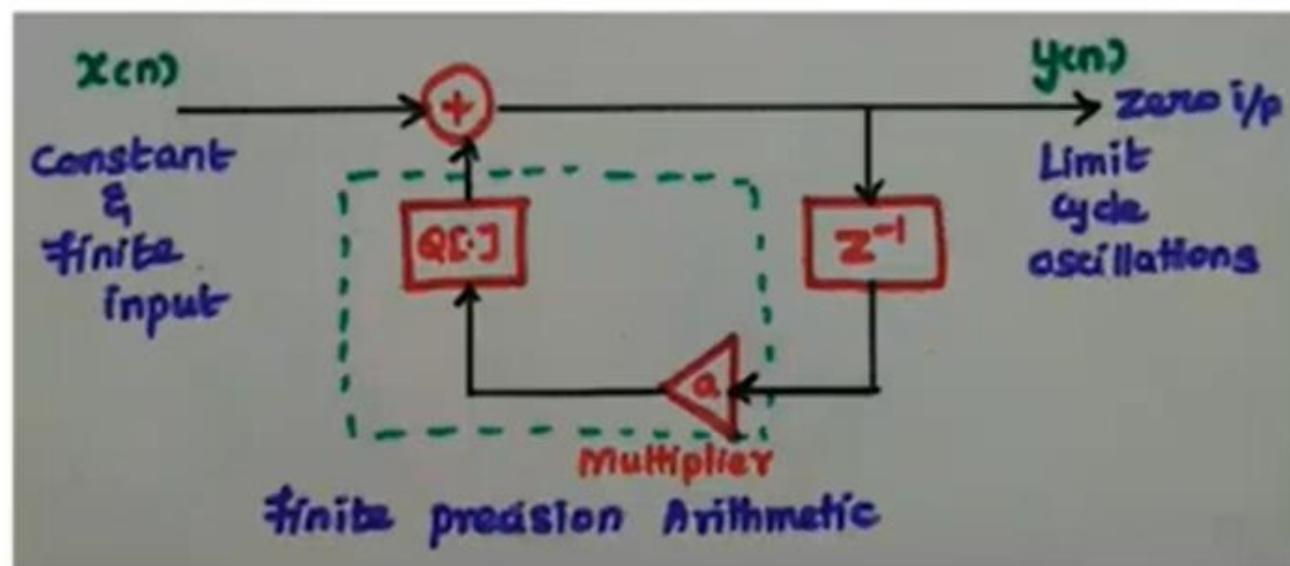
**Pole locations are changed from**

**0.5- 0.5**

**0.45-0.375**

# *Limit Cycle Oscillations*

\* *Zero - Input Limit Cycle Oscillations  
Due to Product Quantization*



$$\text{Dead Band} = \frac{\frac{1}{2}a^2\alpha^b}{1-|\alpha|}$$

# Limit Cycle Oscillations

\* Defn:-

→ It is a low-level oscillation due to the product quantization effects in a stable IIR digital filters [Recursive filters]

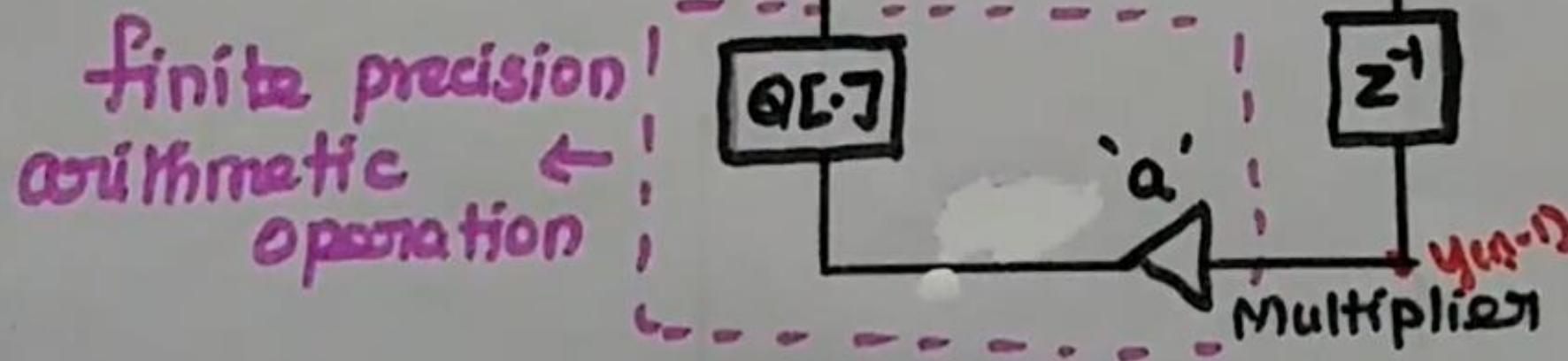
## \* Zero - Input Limit cycle oscillation:-

- When a stable IIR filter is excited by a constant and finite input signal, the output must decay to zero.
- But it does not occur due to the product quantization process which produces the non-linearity's in the system operation.

\* Ex:

Consider a first order IIR filter,

$$y(n) = x(n) + a y(n-1)$$



$$y(n) = x(n) + Q[a y(n-1)]$$

\* The non-linearities due to this finite precision arithmetic operations cause Periodic oscillations in the output

→ Such oscillations in recursive systems are called zero - input limit cycle oscillations.

Analysis of Limit cycle oscillations:-

[Important Problem in Limit cycle oscillation]

The first order recursive filter is given as ,

$$y(n) = x(n) + a y(n-1)$$

$$y(n) = x(n) + Q[a y(n-1)]$$

where,

$a = -\frac{1}{2}$  and the input signal is,

$$x(n) = \begin{cases} 0.875 & ; n=0 \\ 0 & ; \text{otherwise} \end{cases}$$

Assume that the register size is  $b = 3+1$

$b \Rightarrow 3 \text{ bits} + 1 \text{ sign bit}$

$n \Rightarrow \text{no. of samples in the sequence}$

n=0 :-

$$y(n) = x(n) + \alpha [a y(n-1)]$$

$$y(0) = x(0) + \alpha [a y(0-1)]$$

$$y(0) = 0.875 + 0$$

$$\boxed{y(0) = 0.875}$$

\*  $D = 1$  :-

$$y(1) = x(1) + Q \left[ \frac{1}{2} \times y(1-1) \right]$$

$$= 0 + Q \left[ \frac{1}{2} y(0) \right]$$

$$y(1) = Q \left[ \frac{1}{2} \times 0.875 \right] = Q[0.4375]$$

$$\begin{array}{r}
 0.4375 \\
 \times 2 \\
 \hline
 0.8750 \rightarrow 0 \\
 \hline
 0.7500 \rightarrow 1 \\
 \hline
 0.5000 \rightarrow 1 \\
 \hline
 1.0000 \rightarrow 1
 \end{array}$$

$$= Q[0.01\textcolor{red}{1}\textcolor{red}{1}]$$

[Round off to 3 bits]

$$\Rightarrow \begin{array}{r}
 0.011 \\
 \hline
 0.100
 \end{array}_{10}$$

$$y_{(1)} = Q[0.0111] = (0.100)_2$$

$$y_{(1)} = 0.5$$

$$\begin{array}{r} 0.100 \\ \downarrow 1 \times \frac{1}{2} \Rightarrow 0.5 \end{array}$$

For n=2 :-

$$y(2) = x(2) + Q \left[ \frac{1}{2} \times y(2-1) \right]$$

$$= 0 + Q \left[ \frac{1}{2} \times y(1) \right] = Q \left[ \frac{1}{2} \times 0.5 \right]$$

$$y(2) = Q [0.25]$$

$$= Q [0.010]$$

$$\begin{array}{r} 0.25 \\ \times 2 \\ \hline 0.50 \end{array} \rightarrow 0$$
  
$$\begin{array}{r} 0.50 \\ \times 2 \\ \hline 1.00 \end{array} \rightarrow 1$$

$$y(2) = (0.010)_2 = 0.25$$

y(2) = 0.25

\*  $n=3$  :-

$$y(3) = x(3) + Q[y_2 y(3-1)]$$

$$= 0 + Q[\frac{1}{2} \times 0.25]$$

$$= Q[0.125] = Q[0.001]$$

$$\begin{array}{r} 0.125 \\ \times 2 \\ \hline 0.250 \rightarrow 0 \\ \hline 2 \\ \hline 0.500 \rightarrow 0 \\ \hline 2 \\ \hline 1.000 \rightarrow 1 \end{array}$$

$$y(3) = (0.001)_2$$

$$y(3) = 0.125$$

\*  $n=4$  :-

$$\begin{aligned}
 y(4) &= 0 + Q \left[ \frac{1}{2} y(3) \right] \\
 &= Q \left[ \frac{1}{2} \times 0.125 \right] \\
 &= Q [0.0625]
 \end{aligned}$$

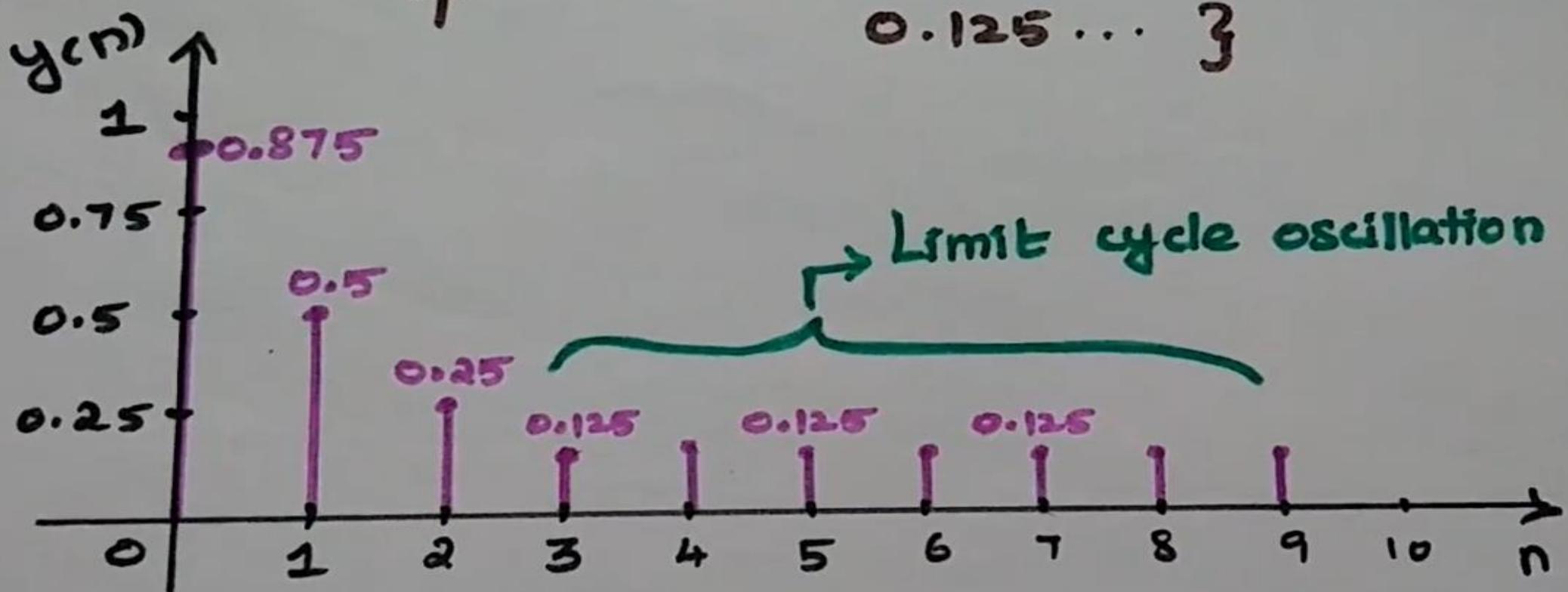
$$\begin{array}{r}
 0.0625 \\
 \times 2 \\
 \hline
 0.12500 \\
 \times 2 \\
 \hline
 0.25000 \rightarrow 0 \\
 \times 2 \\
 \hline
 0.50000 \rightarrow 0 \\
 \times 2 \\
 \hline
 1.00000 \rightarrow 1
 \end{array}$$

$$\begin{aligned}
 y(4) &= Q [0.0625] \\
 &\quad \Downarrow \\
 &= Q [0.000\overset{1}{\cancel{0}}\overset{1}{\cancel{1}}] \\
 &\quad \overset{0.0001}{\cancel{+}} \\
 &\Rightarrow \underline{\underline{0.001}}
 \end{aligned}$$

$$y(4) = (0.001)_2$$

$$y(4) = 0.125$$

$$y(n) = \{ 0.875, 0.5, 0.25, 0.125, 0.125, \\ 0.125 \dots \}$$



When  $n \geq 4$ , the value of  $y(n)$  is equal to 0.125. This is known as limit cycle oscillations.

## \* Dead Band :-



- The limit cycles occur as a result of the quantization effects in multiplications.
- The amplitudes of the output during a limit cycle are confined to a range of values.
- This range of values is known as dead band.

\* For the first order filter, the dead band is given as,

$$\text{Dead Band} = \frac{\frac{1}{2} a^2 b}{1 - |a|}$$

Note :-

→ Limit cycle oscillation does not exist in FIR filters since there is no feed back network in the structure.

*Zero - Input  
Limit Cycle Oscillation*

*An Important Problem solved*

$$y(n) = x(n) + 0.95 y(n-1)$$

Pb10

Find the characteristics of a limit cycle oscillation of the given first order digital system.

$$y(n) = 0.95 y(n-1) + x(n)$$

where,

$$x(n) = \begin{cases} 0.875 & ; n=0 \\ 0 & ; \text{otherwise} \end{cases}$$

Also determine the dead band of the filter.

Solution:-

Given :

$$a = \alpha = 0.95 ;$$

$$x(0) = 0.875$$

Assume that the number of bits :

$b \Rightarrow 4$  bits + 1 sign bit .

The output with rounding value is given as ,

$$y(n) = x(n) + Q [0.95 y(n-1)]$$

$n=0$  :-

$$\begin{aligned}y(0) &= x(0) + \alpha [0.95 y(0-1)] \\&= 0.875 + \alpha [0]\end{aligned}$$

$$y(0) = 0.875$$

\*  $n=1$  :-

$$y^{(1)} = x^{(1)} + \alpha [0.95 y^{(1-1)}]$$

$$= 0 + \alpha [0.95 y^{(0)}]$$

$$= \alpha [0.95 \times 0.875]$$

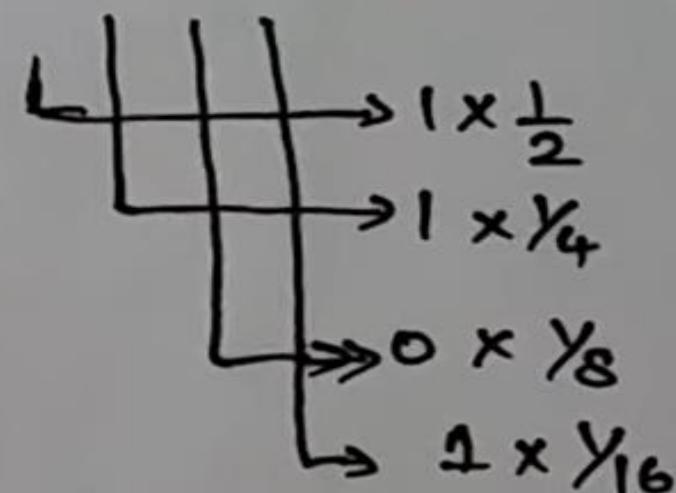
$$y^{(1)} = \alpha [0.83125]$$

$$\Rightarrow Q[0.11010\ldots]_2$$

$$\begin{array}{r}
 0.83125 \\
 \times 2 \\
 \hline
 1.66250 \rightarrow 1 \\
 \times 2 \\
 \hline
 1.32500 \rightarrow 1 \\
 \times 2 \\
 \hline
 0.65000 \rightarrow 0 \\
 \times 2 \\
 \hline
 1.30000 \rightarrow 1 \\
 \times 2 \\
 \hline
 0.60000 \rightarrow 0 \downarrow
 \end{array}$$

$$\Rightarrow Q[0.\underline{\underline{11010}}\dots]_2$$

$$y(1) \Rightarrow (0.1101)_2$$



$$\Rightarrow \frac{1}{2} + \frac{1}{4} + \frac{1}{16} \Rightarrow \frac{8+4+1}{16} = \frac{13}{16}$$

$$y(1) = Q[0.11010\dots] = (0.1101)_2$$

$y(1) = 0.8125$

\*  $n = 2$  :-

$$\begin{aligned}y(2) &= x(2) + \alpha [0.95 y(2-1)] \\&= 0 + \alpha [0.95 \times 0.8125] \\&= \alpha [0.771875]_{10} \\&= \alpha [0.110001\cdots]_2\end{aligned}$$

$$y(2) = \alpha [0.110001]_2$$

$$y(2) = (0.1100)_2 = (0.75)_{10}$$

$$y(2) = 0.75$$

\*  $n=3$  :-

$$\begin{aligned}y(3) &= x(3) + [0.95 y(3-1)] \\&= 0 + 0.95 \times 0.15 \\&= 0.1425\end{aligned}$$

$$= 0.101101\dots_2 \Rightarrow (0.1011)_2$$

$$y(3) = (0.1011)_2 = (0.6875)_{10}$$

$y(3) = 0.6875$

\*  $n=4$  :-

$$y(4) = xc(4) + Q[0.95 y(4-1)] \\ = 0 + Q[0.95 \times 0.6875]$$

$$= Q[0.653125]_{10}$$

$$= Q[0.101001\ldots]_2 = [0.1010]_2$$

$$y(4) = (0.1010)_2 = (0.625)_{10}$$

$$\boxed{y(4) = 0.625}$$

\*  $n = 5$  :-

$$\begin{aligned}y(5) &= x(5) + Q[0.95 y(5-1)] \\&= 0 + Q[0.95 \times 0.625] \\&= Q[0.59375]_{10} \\&= Q[0.10011]_2 \quad \begin{array}{r} 0.1001 \\ \hline 1010 \end{array} \\&= [0.1010]_2\end{aligned}$$

$$y(5) = (0.1010)_2 = (0.625)_{10}$$

$$y(5) = 0.625$$

$$y_{cn}) = \{ 0.875, 0.8125, 0.75, 0.6875, \\ 0.625, 0.625, 0.625 \dots \} \\ y_{(0)} \quad y_{(1)} \quad y_{(2)} \quad y_{(3)} \\ y_{(n)} \quad y_{(5)} \dots \dots \}$$

For  $n \geq 5$ , the output remains constant at 0.625 causing limit cycle oscillations.

Dead Band :-

$$\text{Dead band} = \frac{\frac{1}{2} \cdot 2^{-b}}{1 - |\alpha|}$$

$$\Rightarrow \frac{\frac{1}{2} \cdot 2^{-4}}{1 - 0.95} \Rightarrow \frac{\frac{1}{2} \cdot \frac{1}{2^4}}{0.05}$$

$$\Rightarrow \frac{Y_{3a}}{0.05} = \frac{0.03125}{0.05}$$

Dead band = 0.625



Convert the following decimal numbers in floating point with five bits for mantissa and three bits for exponent.

(i) 4.5

first convert the decimal number into binary.

The binary equivalent of 4 is 100

The binary equivalent of .5 is .1

$$0.5 \times 2 = 1$$

$$4.5_{10} = .100 \cdot 1_2$$

we have to represent by using 5 bit (4 + 1 sign bit)

If the binary is shifted three bits left we will get  
 point shifted three bits left so +3.  
 exponent.

$$.1001 \times 2$$

$$0.1001 \times 2$$

$\hookrightarrow$  sign bit for mantissa (0 for positive  
 1 for negative)

Floating point representation of  $4.5_{10} = 0.1001 \times 2^{011}$

(ii)  $1.5_{10}$

first convert 1.5 to binary

$$1.01_2$$

To bring this into floating point format the binary point is shifted 1 bit left.

$$0.11 \times 2^{001}$$

Since we have to represent by using five bit mantissa the above one is modified as

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

$$1.5_{10} = 0.1100 \times 2^{001}$$

(iii)  $6.5_{10}$

$$6.5_{10} = 110.1_2$$

Shifting the binary point three bit left.

$$0.1101 \times 2^{011}$$

$$6.5_{10} = 0.1101 \times 2^{011}$$

(iv)  $7_{10}$

convert  $7_{10}$  to binary

$111_2$

We have to represent the binary in floating point format

The floating point format is

$2^C \times M$

$C \rightarrow$  exponent (first bit represents sign)

$M \rightarrow$  Mantissa (first bit represents sign)  
positive number - 0  
negative numbers - 0

$111_2$  is  $111.0$

Shifting the binary point

three bits left

$0.1110 \times 2^{011}$

[Since the binary point is shifted three bits left in the exponent we are having 3].

$$7_{10} = 0.1110 \times 2^{011}$$

(v)  $-7_{10}$

convert  $7_{10}$  to binary

$111_2$ .

111.0

$$1.1110 \times 2^{011}$$

↳ sign bit (1 since we have -7)

$$\boxed{-7_{10} = 1.1110 \times 2^{011}}$$

(V<sup>i</sup>) 0.25<sub>10</sub>

Convert 0.25<sub>10</sub> to binary.

$$\begin{array}{r} 0.25 \times \\ \hline 2 \\ \hline 0.5 \times \\ \hline 2 \\ \hline 1.0 \end{array}$$

$$0.25_{10} = .01_2$$

$$.01 \times 2^{000}$$

$$\boxed{0.25_{10} = 0.0100 \times 2^{000}}$$

(Vii) -0.25<sub>10</sub>

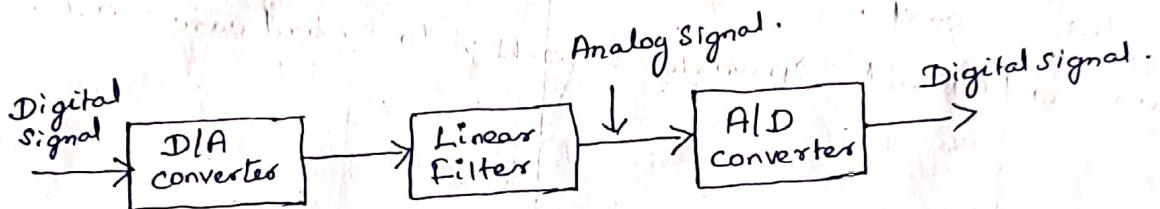
$$\boxed{0.25_{10} = 1.0100 \times 2^{000}}$$

## Multirate Signal Processing :-

When the digital signal processing system handles processing at multiple sampling rates, then it is called multirate digital signal processing.

There are two methods available for Sampling rate conversion.

1) First method - D/A conversion and resampling at required rate



- \* In this method, a digital signal is made to pass through a D/A converter then it is filtered if necessary, and then it is resampled at the desired sampling rate. The resampling of analog signal is carried out by using an A/D converter.

2). Second method - Sampling rate conversion in digital domain.

- \* In this method, Sampling rate conversion is carried out entirely in the digital domain.
- \* This method does not need a D/A and A/D converters.
- \* This method uses interpolator, or decimator or both depending upon the sampling rate conversion factor.

### First method

#### Advantages :-

(i) New sampling rate can have any value. It need not be related to old sampling rate.

(ii) The processing is straight forward and simplest.

#### Disadvantages :

(i) D/A and A/D conversion introduces additional distortion in the signal.

(ii) Two additional converters adds to hardware cost.

### Second Method :-

#### Advantages :-

(i) Distortion due to Sampling and quantization is reduced.

(ii) No hardware cost, since processing is totally in digital domain.

### Applications of Multirate DSP.

Multirate DSP find its application in

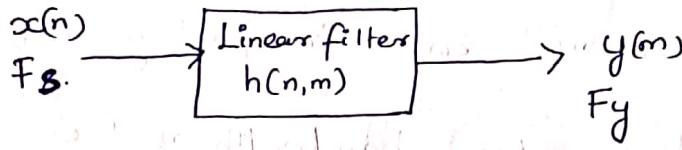
(i) Sub-band coding.

(ii) Voice privacy using analog phone lines.

(iii) Signal compression by Subsampling.

(iv) A/D, D/A Converters.

## Sampling Rate Conversion as Linear filtering.

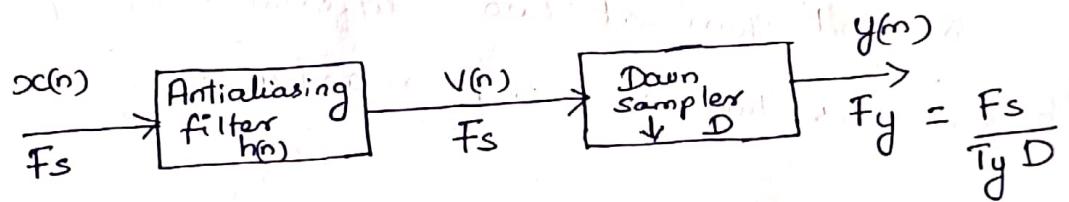


- \* The input signal has the sampling rate of  $F_s$ .
- \* The output signal has the sampling rate of  $F_y$ .
- \* The linear filter has impulse response  $h(n,m)$ , which relates  $x(n)$  and  $y(m)$  at two different sampling rates  $n$  and  $m$ .  
Hence  $h(n,m)$  has two variables.

## Decimation by a Factor D

- \* Decimation by a factor  $D$ , means to reduce the sampling rate by a factor  $D$ . It is also called down sampling by  $D$ .
- \* Let us consider a discrete-time signal  $x(n)$  with spectrum  $X(\omega)$  is to be down sampled by an integer factor  $D$ .
- \* The spectrum  $X(\omega)$  is assumed to be non-zero in the frequency interval  $0 \leq |\omega| \leq \pi$ .
- \* Let the Sampling frequency  $F_s$  and the maximum frequency  $F_{max}$  be related as  $F_{max} \leq \frac{F_s}{2D}$ .

- \* If we reduce the sampling rate simply by selecting every  $D^{\text{th}}$  value of signal  $x(n)$ , the resulting signal would be an aliased version of  $x(n)$  with folding frequency  $\frac{f_s}{D}$ .
- \* To avoid aliasing, the bandwidth of the signal  $x(n)$  is reduced to  $F_{\max} \leq \frac{f_s}{2D}$  or  $\omega \leq \frac{\pi}{D}$



$$f_s = \frac{1}{T_s}$$

- \* The input signal or sequence  $x(n)$  is passed through a antialiasing filter ie. Low pass filter.
- \* The LPF is characterised by the impulse response  $h(n)$  and frequency response  $H_D(\omega)$ .

$$H_D(\omega) = \begin{cases} 1, & |\omega| \leq \frac{\pi}{D} \\ 0, & \text{elsewhere} \end{cases}$$

- \* The filter eliminates the spectrum of  $X(\omega)$  in the range  $\frac{\pi}{D} < \omega < \pi$ .

## Derivation of Decimation Equation:-

The output of the filter is a sequence  $v(n)$  which is given by

$$v(n) = \sum_{k=0}^{\infty} h(k) x(n-k) \quad \text{--- (1)}$$

when  $v(n)$  is downsampled by the factor  $D$ , we get

$$y(m) = v(mD)$$

$$y(m) = \sum_{k=0}^{\infty} h(k) x(mD-k) \quad \text{--- (2)}$$

- \* The downsampling operation is a time varying operation

## Relationship between the spectrums of $x(n)$ and $y(m)$

- \* The frequency-domain characteristics of the output sequence  $y(m)$  may be obtained by relating the spectrum of  $y(m)$  to the spectrum of the input sequence  $x(n)$ .

The sequence  $\tilde{v}(n)$  can be defined as follows.

$$\tilde{v}(n) = \begin{cases} v(n), & n = 0, \pm D, \pm 2D, \dots \\ 0, & \text{elsewhere} \end{cases} \quad \text{--- (3)}$$

- \* The sequence  $\tilde{v}(n)$  may be viewed as a sequence obtained by multiplying  $v(n)$  with a periodic train of impulses  $p(n)$ , with period  $D$ .

$$\tilde{V}(n) = V(n) \cdot P(n) \quad \text{--- } ④$$

The discrete Fourier series representation of

$P(n)$  is given by

$$P(n) = \frac{1}{D} \sum_{k=0}^{D-1} e^{\frac{j2\pi kn}{D}} \quad \text{--- } ⑤$$

$$y(m) = V(n) \cdot P(n)$$

$$y(m) = \tilde{V}(mD) \quad \text{--- } ⑥$$

$$y(m) = V(mD) P(mD)$$

Taking z-transform of the output sequence  $y(m)$

$$y(m) = \tilde{V}(mD)$$

$$Y(z) = \sum_{m=-\infty}^{\infty} \tilde{V}(mD) z^{-m} \quad \text{--- } ⑦$$

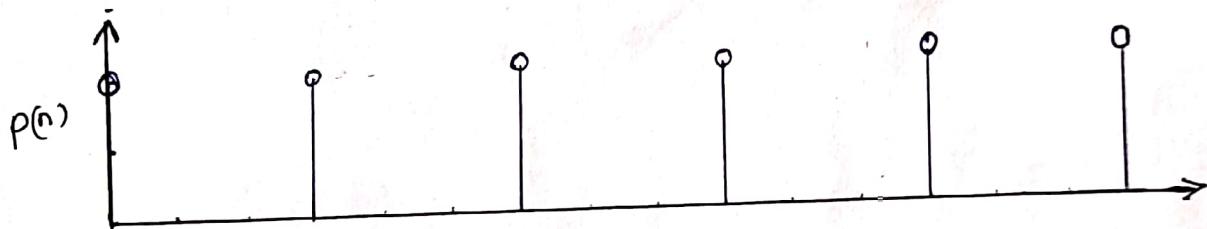
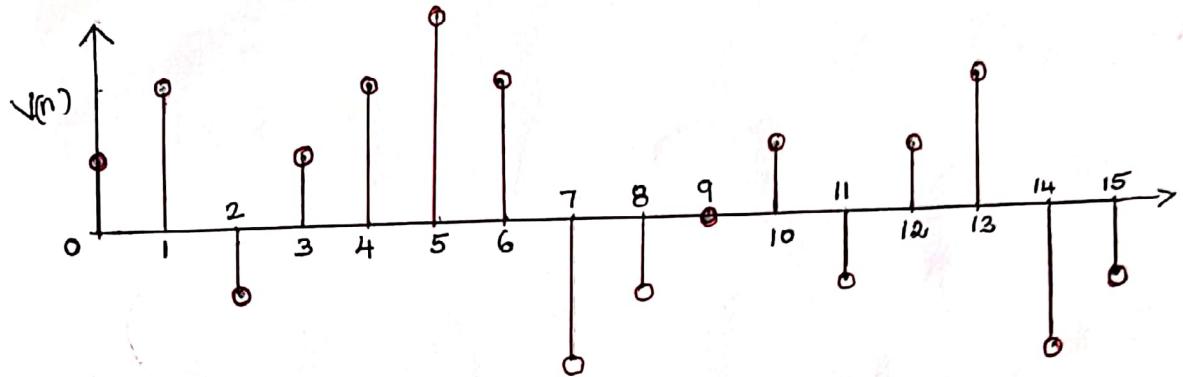
we know that  $\tilde{V}(m) = 0$ , except multiples of  $D$ .

Hence the variable 'm' in above equation can be manipulated as follows.

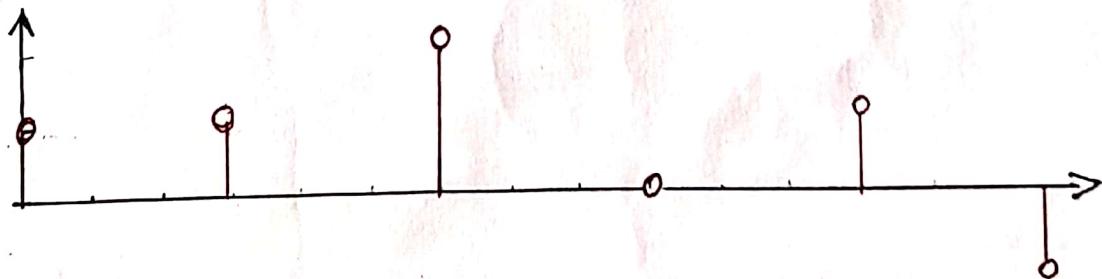
$$Y(z) = \sum_{m=-\infty}^{\infty} \tilde{V}(m) z^{-\frac{m}{D}}$$

$$Y(z) = \sum_{m=-\infty}^{\infty} V(m) P(m) z^{-\frac{m}{D}} \quad \text{--- } ⑧$$

Down Sampling operation. Observed as multiplication of two sequences.



$$y(m) = v(n) \cdot p(n)$$



$$Y(z) = \sum_{m=-\infty}^{\infty} v(m) \cdot \left[ \frac{1}{D} \sum_{k=0}^{D-1} e^{j \frac{2\pi m k}{D}} \right] z^{-\frac{m}{D}}$$

$$Y(z) = \sum_{m=-\infty}^{\infty} v(m) \left[ \frac{1}{D} \sum_{k=0}^{D-1} e^{-j \frac{2\pi m k}{D}}, z^{\frac{+1}{D}} \right]^{-m}$$

$$Y(z) = \frac{1}{D} \sum_{k=0}^{D-1} \sum_{m=-\infty}^{\infty} v(m) \cdot \left[ e^{-j \frac{2\pi k}{D}}, z^{\frac{+1}{D}} \right]^{-m}$$

$$Y(z) = \frac{1}{D} \sum_{k=0}^{D-1} V \left[ e^{j \frac{2\pi k}{D}}, z^{\frac{1}{D}} \right] \quad \text{--- (9)}$$

where  $V \left[ e^{-j \frac{2\pi k}{D}}, z^{\frac{1}{D}} \right] = \sum_{m=-\infty}^{\infty} v(m) \left[ e^{-j \frac{2\pi k}{D}}, z^{\frac{1}{D}} \right]^{-m}$

But  $v(n) = h(n) * x(n)$

Hence their  $z$  transform will multiply.

i.e  $V(z) = H(z) * X(z)$

$$Y(z) = \frac{1}{D} \sum_{k=0}^{D-1} H \left( e^{-j \frac{2\pi k}{D}}, z^{\frac{1}{D}} \right) \times \left( e^{-j \frac{2\pi k}{D}}, z^{\frac{1}{D}} \right) \quad \text{--- (10)}$$

- If  $Y(z)$  is evaluated on unit circle, then we get spectrum of  $y(m)$ . For this we have to put  $z = e^{j\omega y}$ .

Here  $w_y$  is frequency of  $y(m)$ .

$$Y(e^{jw_y}) = \frac{1}{D} \sum_{k=0}^{D-1} H\left(e^{-j\frac{2\pi k}{D}}, e^{\frac{jw_y}{D}}\right) \times \left(e^{-j\frac{2\pi k}{D}}, e^{\frac{jw_y}{D}}\right)$$

$$Y(e^{jw_y}) = \frac{1}{D} \sum_{k=0}^{D-1} H\left(e^{\frac{jw_y - 2\pi k}{D}}\right) \times \left(e^{\frac{w_y - 2\pi k}{D}}\right)$$

We can write  $Y(e^{jw_y})$  as  $Y(w_y)$

$$Y(w_y) = \frac{1}{D} \sum_{k=0}^{D-1} H\left(\frac{w_y - 2\pi k}{D}\right) \times \left(\frac{w_y - 2\pi k}{D}\right)$$

From the above equation we observe that  $\times\left(\frac{w_y - 2\pi k}{D}\right)$  indicates replicas of  $\times\left(\frac{w_y}{D}\right)$  at  $0, 2\pi, 4\pi, 6\pi, \dots, 2\pi(D-1)$ . These replicas are generated due to down sampling operation.

If we consider only the first spectrum substitute  $k=0$  in the above equation. we will get :

$$Y(w_y) = \frac{1}{D} H\left(\frac{w_y}{D}\right) \times \left(\frac{w_y}{D}\right) \quad \text{--- (1)}$$

The antialiasing filter bandlimits the signal to  $\frac{\pi}{D}$ .

It has the magnitude response of

$$H(\omega_s) = \begin{cases} 1 & \text{for } \omega_s \leq \frac{\pi}{D}, \\ 0 & \text{elsewhere.} \end{cases}$$

— (12)

$$\omega_s = 2\pi \frac{f}{f_s}$$

$$\omega_y = 2\pi \frac{f}{f_y}$$

The two sampling rates are related as  $f_y = \frac{f_s}{D}$

$$\omega_y = 2\pi \frac{f}{\frac{f_s}{D}}$$

$$\omega_y = 2\pi f \cdot D$$

$$\boxed{\omega_y = \omega_s D}$$

or

$$\boxed{\omega_s = \frac{\omega_y}{D}}$$

$$H\left(\frac{\omega_y}{D}\right) = \begin{cases} 1 & \text{for } \frac{\omega_y}{D} \leq \frac{\pi}{D}, \\ 0 & \text{elsewhere.} \end{cases}$$

Putting  $H\left(\frac{\omega_y}{D}\right) = 1$  in equation.

$$Y(\omega_y) = \frac{1}{D} \times \left(\frac{\omega_y}{D}\right)$$

13'

The above equation relates the spectrums of  $x(n)$  and  $y(n)$ .

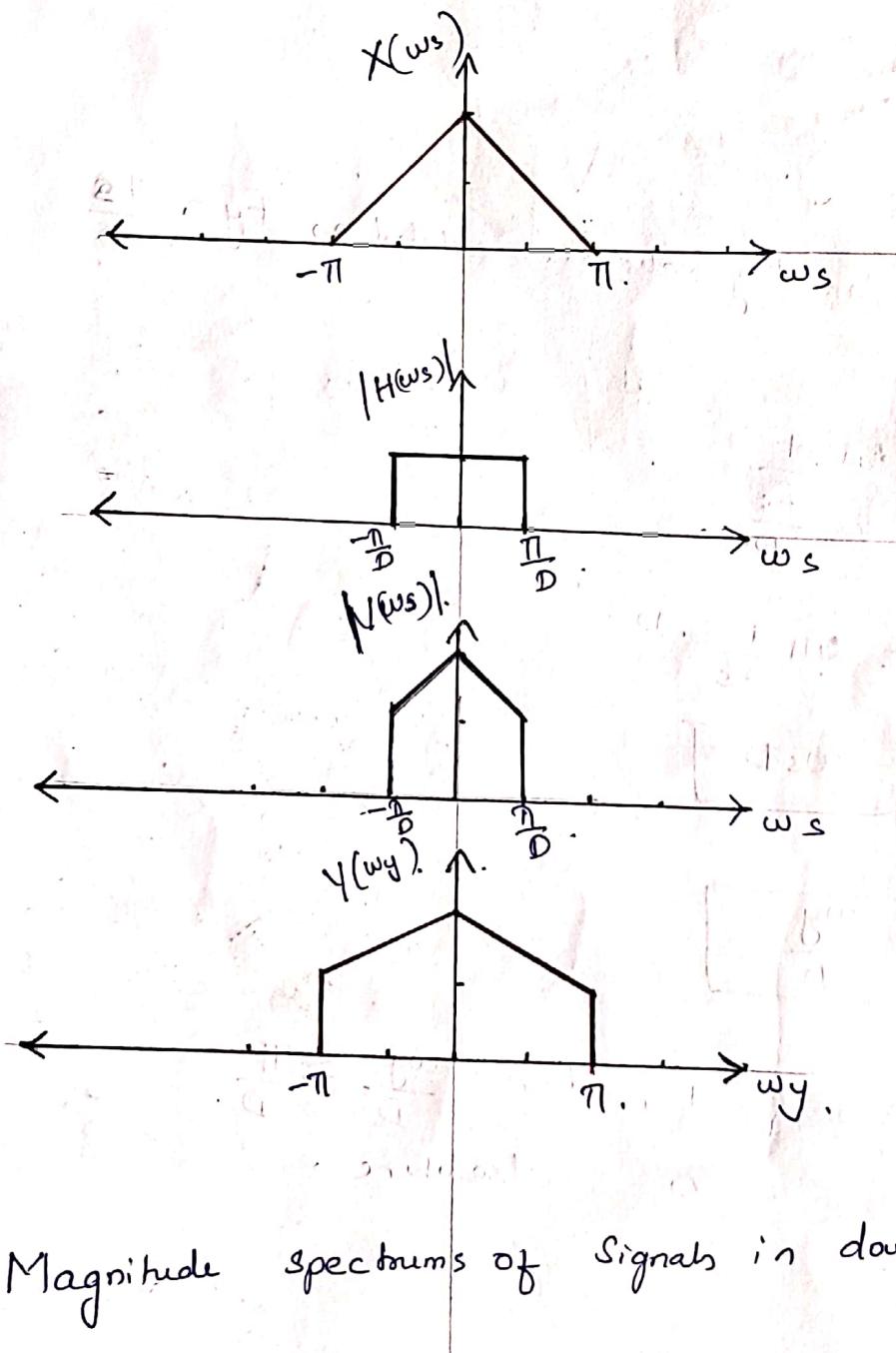


Fig. Magnitude spectrums of Signals in down sampling operation.

## Interpolation by factor I.

- \* Interpolation by a factor I, means to increase the sampling rate by a factor I.
- \* It is also called upsampling by I.
- \* If the Sampling frequency of input signal is  $f_s$ , then it is increased by I. Thus the Sampling frequency of the output signal is  $F_y = I f_s$ .
- \* The interpolator simply puts  $(I-1)$  zeros between successive samples of  $x(n)$ .
- \* The anti-imaging filter removes the image spectrums and interpolates the samples of  $V(m)$ .



- \* The output  $y(m)$  of the anti-imaging filter is given as the convolution of  $h(m)$  and  $V(m)$  i.e.,

$$y(m) = \sum_{k=-\infty}^{\infty} h(m-k) V(k)$$

But the value of  $V(k)$  is zero except at multiples of I.

Hence the product

$$h(m-k) v(k) = \begin{cases} 0 & \text{If } k \text{ is not integer multiple of } I \\ \text{Nonzero} & \text{If } k \text{ is integer multiple of } I. \end{cases}$$

Hence we can replace ' $k$ ' by  $KI$  in the above equation of convolution i.e.

$$y(m) = \sum_{k=-\infty}^{\infty} h(m - KI) v(KI)$$

But  $v(KI) = x(K)$ , hence above, equation becomes,

$$y(m) = \sum_{K=-\infty}^{\infty} h(m - KI) x(K) \quad \text{--- (1)}$$

This is an equation for output signal.

### Relationship between the spectrums of $x(n)$ and $y(n)$

$V(m)$  can be expressed in terms of  $x(n)$  as.

$$V(m) = \begin{cases} x\left(\frac{m}{I}\right), & m = 0, \pm I, \pm 2I, \dots \\ 0 & \text{Otherwise.} \end{cases}$$

This means  $V(m)$  is non zero only at integer multiples of  $I$ .

Taking  $Z$  transform of  $V(m)$  we get.

$$V(z) = \sum_{m=-\infty}^{\infty} V(m) z^{-m}$$

Note that  $V(n) = 0$ , except when 'm' is integer multiple of 1

Hence above equation can be written as

$$V(z) = \sum_{m=-\infty}^{\infty} V(mI) z^{-mI}$$

But  $V(mI) = x(m)$  Hence the above equation becomes.

$$V(z) = \sum_{m=-\infty}^{\infty} x(m) z^{-mI}$$

$$= \sum_{m=-\infty}^{\infty} x(m) (z^I)^{-m}$$

$$V(z) = X(z^I)$$

— (2)

If  $V(z)$  is evaluated on unit circle, then its spectrum

can be obtained.

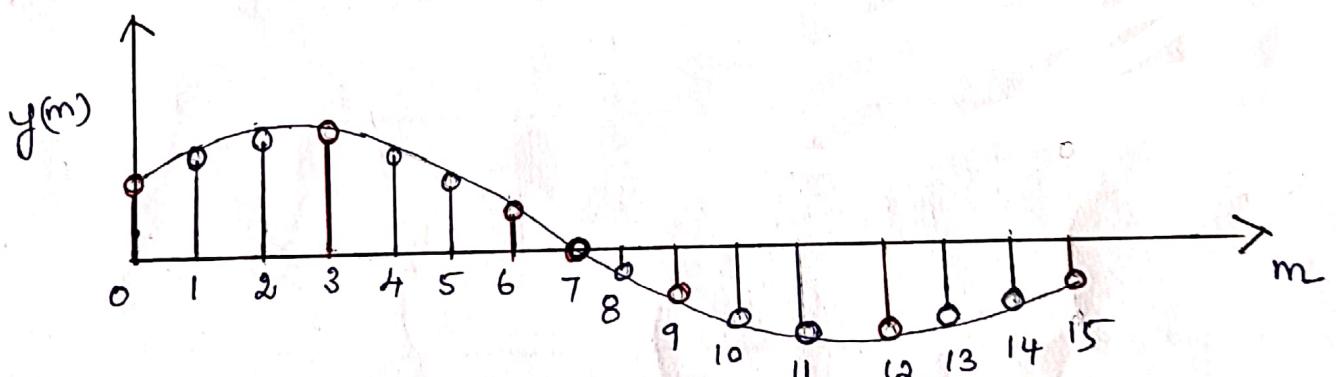
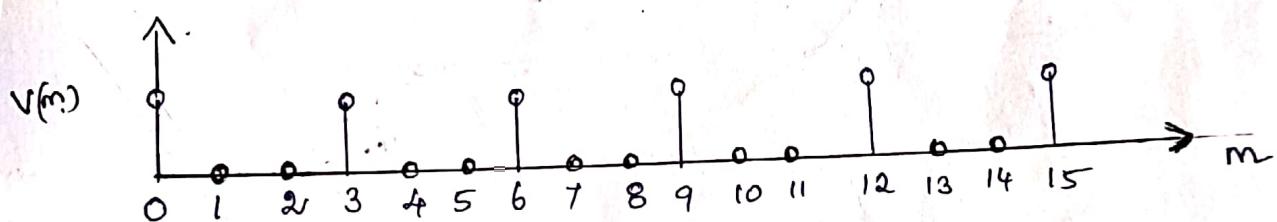
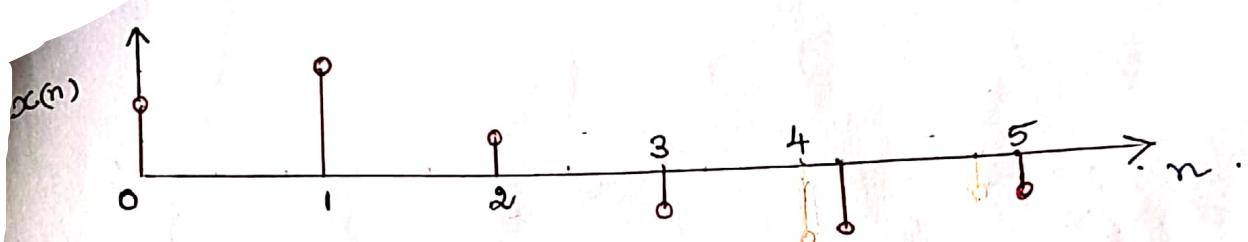
putting  $z = e^{j\omega y}$  in above equation.

$$V(e^{j\omega y}) = X(e^{j\omega y I})$$

$V(e^{j\omega y})$  can be written as  $V(\omega y)$  and  $X(e^{j\omega y})$  as  $X(\omega y)$ .

$$V(\omega y) = X(\omega y I). \quad — (3)$$

Waveform of up sampling



But  $F_y = \frac{1}{I} f_s$ .

$$\omega_y = 2\pi \frac{f}{F_y} = 2\pi \frac{f}{\frac{1}{I} f_s} = \frac{\omega_s}{I}$$

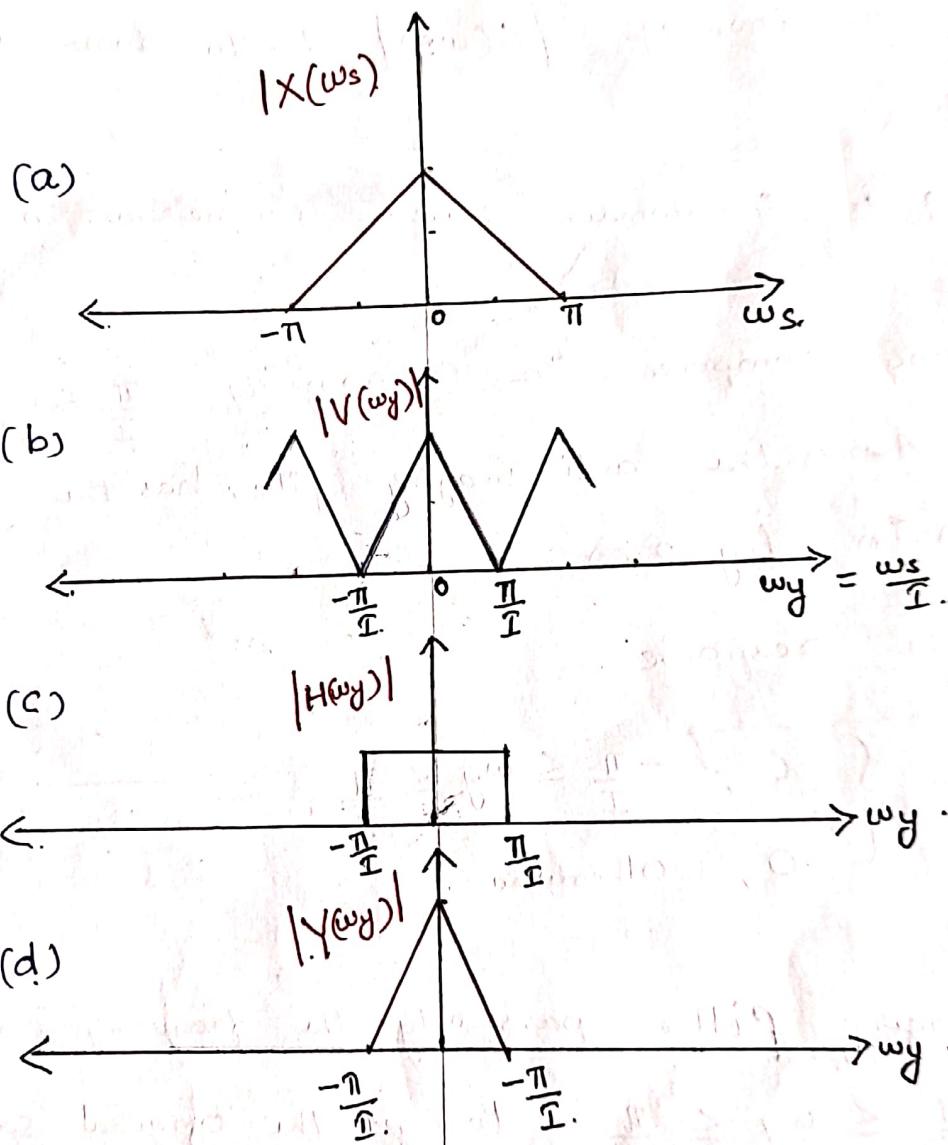


Fig 2.

Spectrum of Signals in upsampling operation.

The fig 2(a) shows an arbitrary spectrum of  $x(n)$ .

The spectrum of  $x(n)$  has the frequencies from  $-\pi$  to  $\pi$ .

$$\text{i.e. } -\pi \leq \omega_s \leq \pi.$$

The fig 2(b) shows the spectrum of  $V(n)$ . It is the  $I$ -fold periodic repetition of  $|X(\omega_s)|$ , due to upsampling operation.

Upsampling in time domain causes compression in frequency domain.

The frequency components in the range of  $-\frac{\pi}{I} \leq \omega_y \leq \frac{\pi}{I}$  are unique. Hence the anti-Imaging filter has the response as shown in fig 2(c).

It has the response.

$$H(\omega_y) = \begin{cases} C, & -\frac{\pi}{I} \leq \omega_y \leq \frac{\pi}{I} \\ 0, & \text{Otherwise} \end{cases} \quad \text{--- (4)}$$

The antiImaging filter passes only the frequency components.

from  $-\frac{\pi}{I} \leq \omega_y \leq \frac{\pi}{I}$  to get the original spectrum back

The Signal  $V(n)$  is passed through the filter. Hence from equation (3) and equation (4) we can write the spectrum of output signal as.

$$Y(\omega_y) = \begin{cases} C \times (C/I), & -\frac{\pi}{I} \leq \omega_y \leq \frac{\pi}{I} \\ 0, & \text{Otherwise.} \end{cases} \quad \text{--- (5)}$$

This equation gives the relationship between the Spectrums of output and input signals.

### Value of Scale factor C .

The scale factor is selected such that

$$Y(m) = \infty \left( \frac{m}{\pi} \right) \quad \text{for } m=0, \pm 1, \pm 2, \dots$$

By inverse fourier transform  $y(m)$  can be obtained as .

$$y(m) = \frac{1}{2\pi} \int_{-\pi}^{\pi} Y(\omega_y) e^{j\omega_m \omega_y} d\omega_y .$$

For  $m=0$

$$y(0) = \frac{1}{2\pi} \int_{-\pi}^{\pi} Y(\omega_y) e^{j0} \cdot d\omega_y .$$

$$\boxed{\text{But } Y(\omega_y) = C \times (\omega_y)^{-1}}$$

$$y(0) = \frac{1}{2\pi} \int_{-\frac{\pi}{\omega_s}}^{\frac{\pi}{\omega_s}} C \times (\omega_y)^{-1} d\omega_y .$$

$$\text{But } \omega_y = \frac{\omega_s}{I} \quad \text{or} \quad \omega_s = \omega_y I$$

Hence the integration limits will be  $-\pi$  to  $\pi$  i.e .

$$y(0) = \frac{1}{2\pi} \int_{-\pi}^{\pi} C \times \left( \frac{\omega_s}{I} \right)^{-1} d\omega_s \frac{1}{I}$$

$$y(0) = \frac{C}{I} \left[ \int_{-\pi}^{\pi} x(\omega_s) d\omega_s \right]$$

$$y(0) = \frac{C}{I} x(0)$$

Thus  $y(0) = x(0)$ , if  $C = I$ .

This is the value of scaling factor.

## Sampling Rate conversion by a Rational factor $I/D$

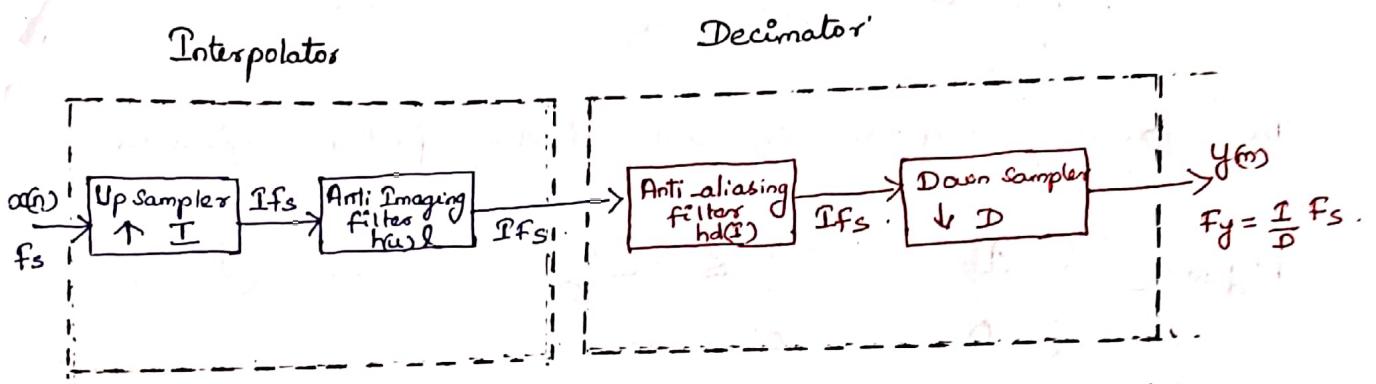


Fig: Cascade of Interpolator and Decimator to obtain Sampling rate conversion by  $I/D$ .

- \* The interpolation by a factor  $I$  is obtained first to increase the sampling rate to  $Ifs$ .
- \* The output of interpolator is then decimated by a factor  $D$ , so that the final output rate is  $f_y = \frac{I}{D} fs$ .
- \* Such operation can be implemented by cascade connection of interpolator and decimator.
- \* In the above figure there are two low pass filters.
- \* The overall cutoff frequency will be minimum of the two cutoff frequencies.

The frequency response of anti Imaging filter is given as.

$$H_u(\omega) = \begin{cases} C & -\frac{\pi}{I} \leq \omega \leq \frac{\pi}{I} \\ 0 & \text{Otherwise.} \end{cases} \quad \text{--- (1)}$$

$$H_u(\omega) = \begin{cases} 1 & -\frac{\pi}{I} \leq \omega \leq \frac{\pi}{I}, \text{ since } C=I \text{ (Scaling)} \\ 0 & \text{otherwise} \end{cases}$$

The frequency response of anti-aliasing filter is given as,

$$H_d(\omega) = \begin{cases} 1 & -\frac{\pi}{D} \leq \omega \leq \frac{\pi}{D} \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

- \* The overall cascading effect of lowpass filters will have a cutoff frequency which is minimum of  $\frac{\pi}{I}$  and  $\frac{\pi}{D}$ .

Hence we can write frequency response of combined filter as.

$$H(\omega) = \begin{cases} 1 & |\omega| \leq \min\left(\frac{\pi}{D}, \frac{\pi}{I}\right) \\ 0 & \text{Otherwise} \end{cases} \quad (3)$$

- \* Thus, Single filter can be used, as defined by above equation. The block diagram can be modified as follows.

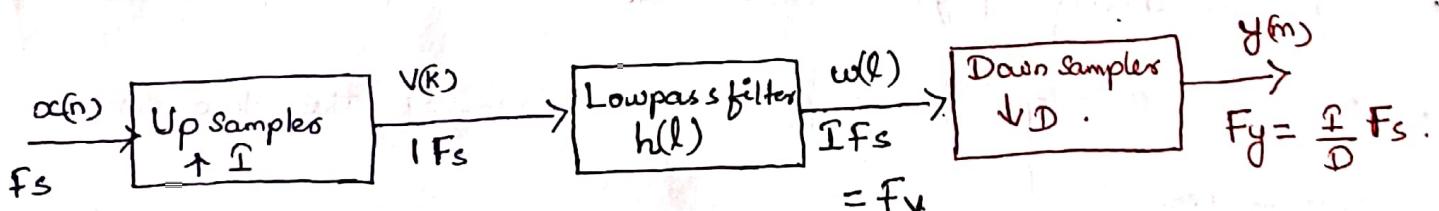


Fig:- Sampling rate conversion by a factor  $\frac{1}{D}$ .

## Derivation for output $y(m)$

The output of lowpass filter is given as

$$\begin{aligned} w(l) &= \sum_{k=-\infty}^{\infty} h(l-k) v(k) \\ &= \sum_{k=-\infty}^{\infty} h(l-kT) x(k) \end{aligned} \quad \text{--- (4)}$$

\* The output of down sampler is given as.

$$y(m) = w(mD)$$

$$y(m) = \sum_{k=-\infty}^{\infty} h(mD-kT) x(k) \quad \text{--- (5)}$$

This is the equation for output sequence.

## Relationship between the spectrums of $x(n)$ and $y(m)$

\* The output signal of upsampler is  $v(k)$  with frequency  $\omega_v$ .

for this signal, using  $V(\omega_y) = X(\omega_y T)$

we can write.

$$V(\omega_v) = X(\omega_v T) \quad \text{--- (6)}$$

\* The signal  $v(k)$  is passed through a low pass filter and the output signal is  $w(l)$

Hence we can write.

$$W(\omega_v) = H(\omega_v) \cdot V(\omega_v)$$

putting the values from equation 3 and equation 6 in above equation we get

$$W(\omega_v) = \begin{cases} \frac{1}{D} \times (\omega_v I), & |\omega_v| \leq \min\left(\frac{\pi}{D}, \frac{\pi}{I}\right) \\ 0 & \text{Otherwise} \end{cases} - (7)$$

We know that the spectrum of downsampled signal is given as.

$$Y(\omega_y) = \frac{1}{D} \times \left( \frac{\omega_y}{D} \right)$$

But the input to the down sampler is  $\omega_v$ .

Hence above equation can be written as.

$$Y(\omega_y) = \frac{1}{D} W\left(\frac{\omega_y}{D}\right) \quad (8)$$

$$\text{But } \omega_v = \frac{2\pi f_v}{f_v} = \frac{2\pi f_v}{f_y} = 2\pi \frac{f_v}{f_y} \cdot \frac{1}{D} = \frac{\omega_y}{D}$$

$$\text{so } \omega_v = \frac{\omega_y}{D}$$

Hence the right hand side of above equation (7) can be written as.

$$Y(\omega_y) = \frac{1}{D} W(\omega_v)$$

putting  $\omega_v$  from equation (6)

$$Y(\omega_y) = \begin{cases} \frac{1}{D} \times (\omega_v I) & |\omega_v| \leq \min\left(\frac{\pi}{D}, \frac{\pi}{I}\right) \\ 0 & \text{Otherwise} \end{cases}$$

Since  $\omega v = \frac{\omega y}{D}$  above equation can be modified as

$$Y(\omega y) = \begin{cases} \frac{\pi}{D} \times \left( \frac{\pi}{D} \omega y \right), & \left| \frac{\omega y}{D} \right| \leq \min \left( \frac{\pi}{D}, \frac{\pi}{I} \right) \\ 0 & \text{Otherwise} \end{cases}$$

In the above equation  $\left| \frac{\omega y}{D} \right| \leq \min \left( \frac{\pi}{D}, \frac{\pi}{I} \right)$  can be written as

$$|\omega y| \leq \min \left( \pi, \frac{D}{I} \pi \right)$$

Then we have :

$$Y(\omega y) = \begin{cases} \frac{\pi}{D} \times \left( \frac{\pi}{D} \omega y \right) & |\omega y| \leq \min \left( \pi, \frac{D}{I} \pi \right) \\ 0 & \text{Otherwise} \end{cases}$$

This is the required equation that relates the spectrums of input and output.

## 8.11 Polyphase Structure of Decimator

In section 6.9.4 we studied about the polyphase realization of FIR filters, where the transfer function  $H(z)$  is decomposed into  $M$  branches given by

$$H(z) = \sum_{m=0}^{M-1} z^{-m} P_m(z^M) \quad (8.22)$$

where

$$P_m(z) = \sum_{n=0}^{[(N+1)/M]} h(Mn + m)z^{-n} \quad (8.23)$$

The  $z$ -transform of an infinite sequence is given by

$$H(z) = \sum_{n=-\infty}^{\infty} h(n)z^{-n} \quad (8.24)$$

In this case  $H(z)$  can be decomposed into  $M$ -branches as

$$\begin{aligned} H(z) &= \sum_{m=0}^{M-1} z^{-m} P_m(z^M) \\ \text{where } P_m(z) &= \sum_{r=-\infty}^{\infty} h(rM + m)z^{-r} \end{aligned}$$

$$\begin{aligned} x_m(r) &= x(rM - m) \\ x_m(-r) &= x(-rM - m) \\ x_m(n - r) &= x[n - (rM + m)] \end{aligned}$$

$$\begin{aligned} H(z) &= \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} z^{-m} h(rM + m)z^{-rM} \\ &= \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} h(rM + m)z^{-(rM+m)} \end{aligned} \quad (8.25)$$

Let  $h(rM + m) = p_m(r)$

$$\begin{aligned} \implies H(z) &= \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} p_m(r) z^{-(rM+m)} \\ Y(z) &= \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} p_m(r) X(z) z^{-(rM+m)} \\ y(n) &= \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} p_m(r) x[n - (rM + m)] \end{aligned} \quad (8.26)$$

Let  $x_m(r) = x(rM - m)$  then

$$y(n) = \sum_{m=0}^{M-1} \sum_{r=-\infty}^{\infty} p_m(r) x_m(n - r) \quad (8.27)$$

$$\begin{aligned} &= \sum_{m=0}^{M-1} p_m(n) \star x_m(n) \\ &= \sum_{m=0}^{M-1} y_m(n) \end{aligned} \quad (8.28)$$

where  $y_m(n) = p_m(n) \star x_m(n)$

The operation  $p_m(n) \star x_m(n)$  is known as polyphase convolution, and the overall process is polyphase filtering.

If  $M = 3$  then

$$\begin{aligned} y(n) &= \sum_{m=0}^2 y_m(n) \\ &= y_0(n) + y_1(n) + y_2(n) \\ &= p_0(n) \star x_0(n) + p_1(n) \star x_1(n) + p_2(n) \star x_2(n) \end{aligned} \quad (8.29)$$

We know that  $x_m(n)$  can be obtained first by delaying  $x(n)$  by  $m$  units, and then downsampling by a factor  $M$ . Next  $y_m(n)$  can be obtained by convolving  $x_m(n)$  with  $p_m(n)$ . The structure of a polyphase decimator with 3 branches and a sampling rate reduction by a factor three is shown in Fig. 8.51.

For a general case of  $M$  branches and a sampling rate reduction by a factor  $M$ , the structure of polyphase decimator is shown in Fig. 8.52. The splitting of  $x(n)$  into the low rate sub-sequence  $x_0(n), x_1(n) \dots x_{M-1}(n)$  is often represented by a commutator.

In the configuration shown in Fig. 8.52 the input values  $x(n)$  enter the delay chain at high rate. Then the  $M$  down sampler sends the group of  $M$  input values

### 8.30 Digital Signal Processing

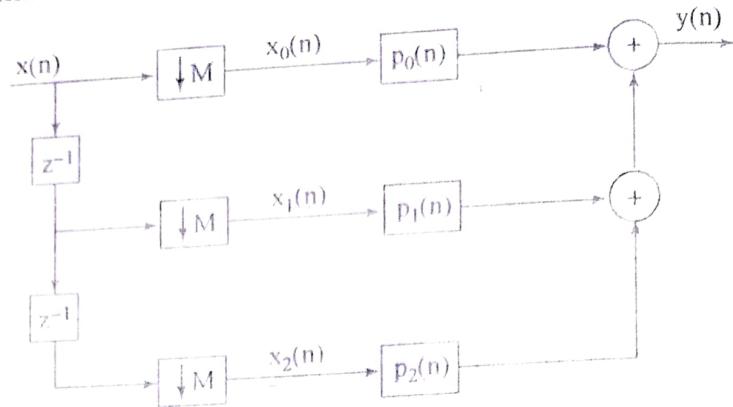


Fig. 8.51 Polyphase structure of a 3 branch decimator

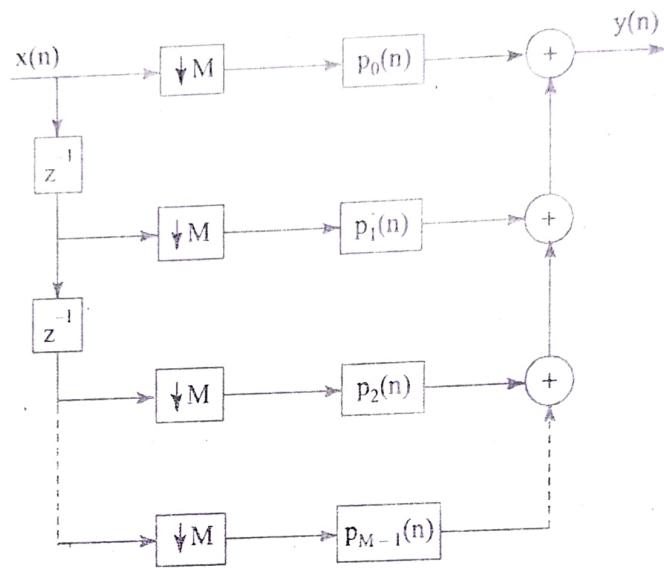


Fig. 8.52 Polyphase structure of a  $M$ -branch decimator

to  $M$  filters at time  $n = mM$ . For example at time  $n = 0$ , the value of  $x_0(0) = x(0), x_1(0) = x(-1), x_2(0) = x(-2), \dots, x_{M-1}(0) = x(-M+1)$  are sent.

In Fig. 8.53 to produce the output  $y(0)$ , the commutator must rotate in counter-clockwise direction starting from  $m = M - 1, \dots, m = 2, m = 1, m = 0$  and give the input values  $x(-M+1), \dots, x(-2), x(-1), x(0)$  to the filters  $p_{M-1}(n), \dots, p_1(n), p_0(n)$ .

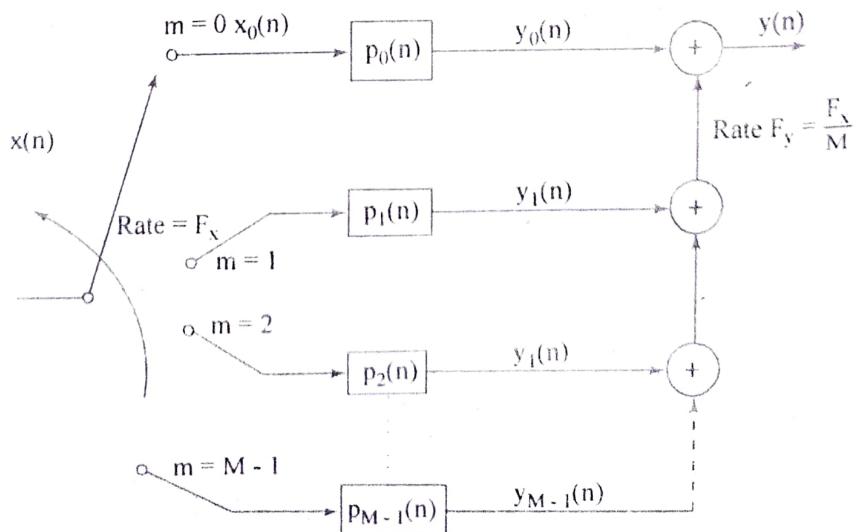


Fig. 8.53 Polyphase decimator with a commutator

## 8.12 Polyphase Decimation Using the $z$ -transform

The  $z$ -transform representation of a decimator is shown in Fig. 8.54.

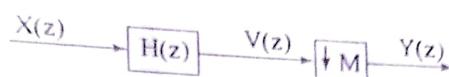


Fig. 8.54

The  $M$ -branch polyphase decomposition of  $H(z)$  is given by

$$H(z) = \sum_{m=0}^{M-1} z^{-m} P_m(z^M) \quad (8.30)$$

The subfilters  $P_0(z), P_1(z), \dots, P_{M-1}(z)$  are FIR filters and when combined in the right phase sequence produce the original filter  $H(z)$ .

Now replacing the transfer function  $H(z)$  in Fig. 8.54 by the polyphase structure results as in Fig. 8.55.

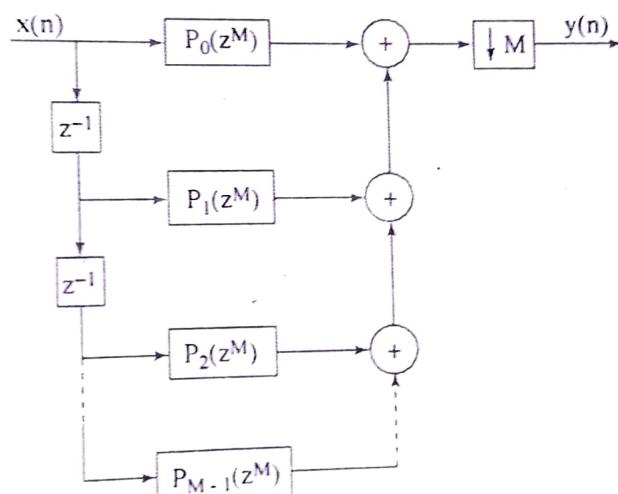


Fig. 8.55

Using first identity Fig. 8.55 can be changed as in Fig. 8.56.

The third identity in Fig. 8.25 can be used to derive the version in Fig. 8.57 in which both the number of filter operations and the amount of memory required are reduced by a factor of  $M$ .

## 8.13 Polyphase Structure of Interpolator

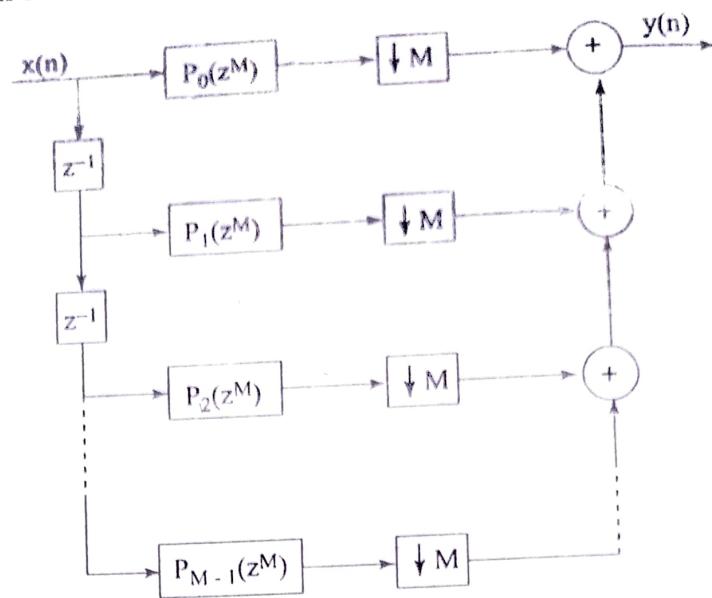
By transposing the decimator structure shown in Fig. 8.57, we can obtain polyphase structure for interpolator, which consists of a set of  $L$  sub-filters connected in parallel as shown in Fig. 8.58.

Here the polyphase components of impulse response are given by

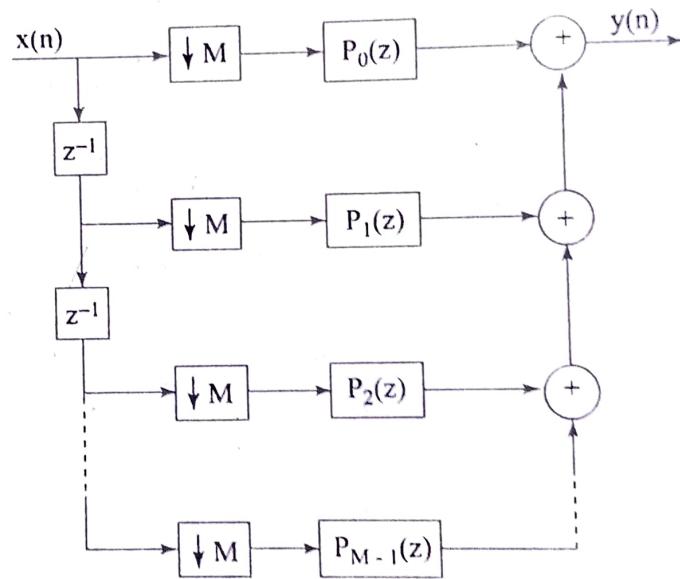
$$p_m(n) = h(nL + m) \quad m = 0, 1, 2, \dots, L - 1 \quad (8.31)$$

$$p_m(n) = h(nL + m) \quad m = 0, 1, 2, \dots, L - 1$$

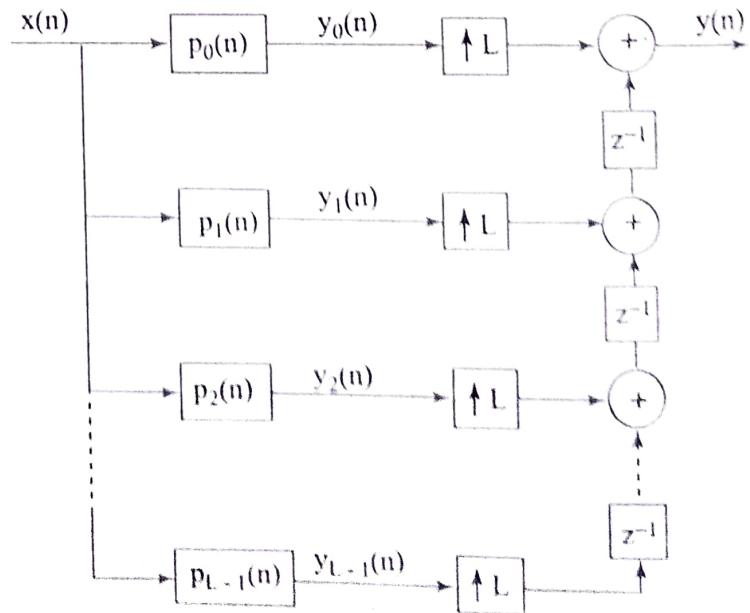
## 8.32 Digital Signal Processing



**Fig. 8.56**



**Fig. 8.57**



**Fig. 8.58**

where  $h(n)$  is the impulse response of anti-imaging filter. The output of  $L$  sub-filters can be represented as

$$y_m(n) = x(n)p_m(n) \quad m = 0, 1, 2, \dots, L-1 \quad (8.32)$$

By upsampling with a factor  $L$  and adding a delay  $z^{-m}$  the polyphase components are produced from  $y_m(n)$ . These polyphase components are all added together to produce the output signal  $y(n)$ .

The output  $y(n)$  also can be obtained by combining the signals  $x_m(n)$  using a commutator as shown in Fig. 8.59.

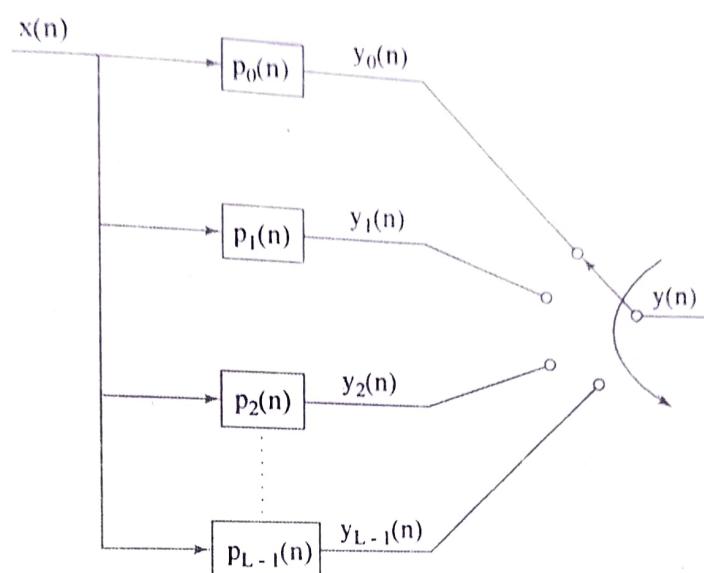


Fig. 8.59

## 8.14 Polyphase Interpolation Using the $z$ -transform

The block diagram of an interpolator is shown in Fig. 8.60.

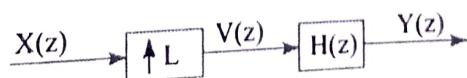


Fig. 8.60

The transfer function  $H(z)$  of the interpolator is given by

$$H(z) = \sum_{m=0}^{L-1} z^{-m} p_m(z^L) \quad (8.33)$$

$H(z)$  in Fig. 8.60 can be realized using polyphase structure as shown in Fig. 8.61.

Using Fourth and Sixth identity an efficient structure for interpolator can be obtained as shown in Fig. 8.62.

### 8.34 Digital Signal Processing

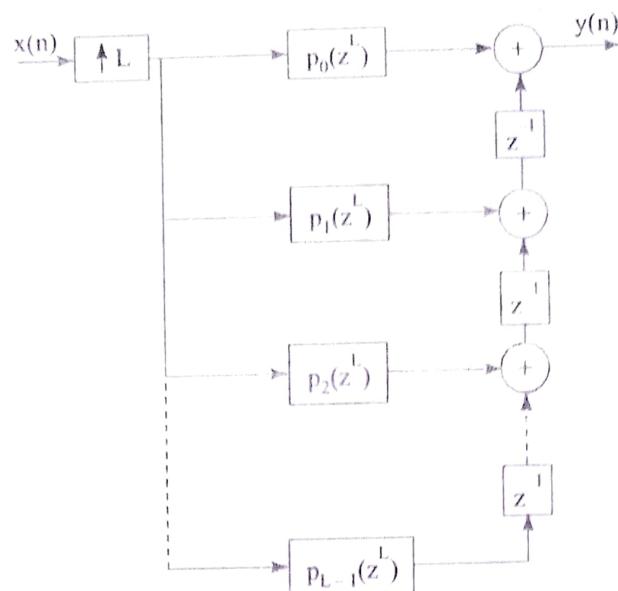


Fig. 8.61

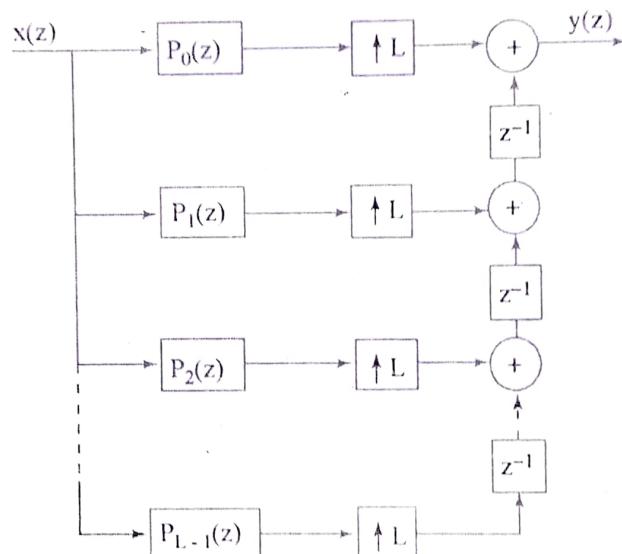


Fig. 8.62

### 8.15 Multistage Implementation of Sampling Rate Conversion

If the decimation factor  $M$  and/or interpolation factor  $L$  are much larger than unity, the implementation of sampling rate conversion in a single stage is computational inefficient. Therefore for performing sampling rate conversion for either  $M \gg 1$  and /or  $L \gg 1$  we go in for multistage implementation.

If the interpolation factor  $L \gg 1$ , then we express  $L$  into a product of positive integers as

$$L = \prod_{i=1}^N L_i \quad (8.34)$$

Then each interpolator  $L_i$  is implemented and cascaded to get  $N$  stages of interpolation and filtering as shown in Fig. 8.63.