SAURABH KUMAR
SC22B146

# AVM867 - VLSI Signal Processing
## Assignment - I

① $-2^{N-1} \le x \le 2^{+(N-1)} - 1$ : N-bit fixed point integer $x$

   (i) $N = 12$

    $\Rightarrow -2^{11} \le x \le 2^{+11} - 1$

    $\Rightarrow -2^{11} \le x \le 2^{+11} - 2^{-0}$

      $Q_{11.0}$ can be represented in the range $[-2^{11}, 2^{11} - 2^{-0}]$

      ∴ Most negative no. $= -2^{+11} = -2048$

   (ii) $N = 12$

    Most positive no. $= 2^{+11} - 1$

                $= +2047$

② $N = 24$

   $-2^{N-1} \le x \le 2^{(N-1)} - 1$

   $\Rightarrow -2^{23} \le x \le 2^{23} - 1$

   Dynamic range $= \dfrac{\text{Max value}}{\text{Min value}}$

             $= \dfrac{2^{N-1} - 1}{1} = \dfrac{2^{23} - 1}{1} \approx 2^{23}$

   Dynamic range in dB $= 20 \log_{10}(\text{Dynamic Range})$

             $= 20 \log(2^{23})$

             $= 20 \times 23 \times 0.301$

             $= 138.46 \text{ dB}$

③ N-bit signed fractional $\rightarrow -1 \le x < 1$

   (i) $N = 24$

    $DR = \dfrac{\text{Max value}}{\text{Min value}} = \dfrac{(1 - 2^{-(N-1)})}{(2^{-(N-1)})}$

    Largest possible positive fraction $= 1 - $ resolution

    $[2^{-(N-1)} : \text{Resolution}]$

    $(2^{-(N-1)}) \rightarrow$ smallest possible +ve fraction = resolution

        $\approx \dfrac{1}{2^{-(N-1)}} = 2^{N-1}$

             $= 2^{23}$

   $DR_{dB} = 20 \log_{10}(2^{23}) = 20 \times 23 \times 0.301 = 138.46 \text{ dB}$

(ii) Precision = resolution    23 bits

$$= 2^{-(N-1)}$$
$$= \frac{1}{2^{23}} \approx 1.19 \times 10^{-7}$$

(iii)  N = 48,

$$DR = 2^{48-1} = 2^{47}$$

↓

$$DR_{dB} = 20 \times 47 \times 0.301$$
$$= 282.47\, dB \longrightarrow 24 \times 20 \times 0.301 = 6.02\,N\, dB\ more$$

$$Precision = 2^{-(48-1)} = 2^{-47} \qquad 48\ 47\ bit \longrightarrow 24\ bits\ more$$
$$\approx 7.1 \times 10^{-15} \longrightarrow \frac{1}{2^N}\ times$$

(iv)  N' : no. of bits (new)

$$\frac{2^{-(N'-1)}}{2^{-(N-1)}} = \frac{1}{8}$$
$$\Rightarrow 2^{-N'+N} = 2^{-3}$$
$$\Rightarrow -N' + N = -3$$
$$\Rightarrow N' = N+3 \longrightarrow 3\ more\ bits$$

④  In $Q_{1.4}$ format,

$$(0.4375)_{10} = 0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4}$$

↓

$$(0.0111)_2$$

ⓐ Signed Representation : (In $Q_{1.4}$ format)

+ 0.4375 → 0  0  0111

− 0.4375 → 1  0  0111

↓

signed bit

ⓑ One's complement

+ 0.4375 → same → 0  0  0111

− 0.4375 → flip all bits → 1  1  1000
            of 000111

ⓒ Two's complement :

+ 0.4375 → same → 0  0  0111

− 0.4375 → +1 to → 1  01  01001
            1's comp.

⑤

$$0.25_{10} \rightarrow 0.0100_2 \rightarrow 0010_2$$
$$-0.625_{10} \rightarrow 0.1010_2 \rightarrow 0101_2$$

$\underbrace{\phantom{0101}}_{Q_{1.3}\ format}$

$\downarrow$ 1's complement

$$1010$$
$$+1$$
$$\overline{1011} \rightarrow \text{2's complement}$$
$$(-0.625)_{10}$$

∴ $0.25 - 0.625$

$$= \quad 0010$$
$$+ \ 1011$$

$\left.\rule{0pt}{20pt}\right\}$ Subtraction

$$\overline{1101} \rightarrow \text{Negative no.}$$

$\downarrow$ 1's comp.

$$0010$$
$$+1$$
$$\overline{0011} \rightarrow \text{2's complement}$$
$$\downarrow$$
$$0.375$$

∴ Result $= -0.375$

⑥  $-0.72 \rightarrow Q_{0.7}$ Representation

$$0.72_{10} \rightarrow (0.1011100)_2$$

In $Q_{0.7}$ Representation,

$$\underset{\underset{sign\,bit}{\downarrow}}{1}\ \underset{\underset{fraction\ part}{\downarrow}}{1011100}$$

⑦  $0xABCD \rightarrow (10 \times 16^3) + (11 \times 16^2) + (12 \times 16^1) + (13 \times 16^0)$

$$= (43981)_{10}$$
$$= (1010 \bullet 1011\ 1100\ 1101)_2$$

Q.15 format $\rightarrow$ 0.1.010 1011 1100 1101

Q8.7 format $\rightarrow$ 0 1010 1011 . 11 00 110

⑧

$x = -0.5 \xrightarrow{Q0.3} 1\ 100$

$y = 0.875 \xrightarrow{Q0.3} 0\ 111$

$X = -0.5 \times 2^3 = -4$

$Y = 0.875 \times 2^3 = 7$

$Z = XY = -4 \times 7 = -28$

$\Rightarrow z = -\dfrac{28}{64} = -0.4375 \qquad$ (for Q.6 format).

and

$xy = -0.5 \times 0.875 = -0.4375$

As $z = xy$, there is <u>no</u> quantization error.

⑨ <u>Double precision</u> floating pt.

1 sign bit — 11 exp. bit — 52 frac. bit

<u>Min$^m$</u>:
(Subnormal)

$exp \to 0\ 0\ 0\ \text{---}$  (11 zeros)

$frac \to 0\ 0\ 0\ \text{---}$  (52 zeros)

$\hookrightarrow 2^{-1022} \times 2^{-52} = \boxed{2^{-1074}}$

<u>Min$^m$</u>:
(Normal)

$exp \to 0\ 0\ \text{--}\ 1$

$frac \to 0\ 0\ 0\ \text{---}$

$\hookrightarrow 2^{1-1023} \times (1.00)$

$= \boxed{2^{-1022}}$

<u>Max</u>:

$exp \to 1\ 1\ 1\ 1\ \text{----}$

$frac \to 0\ 0\ 0\ \text{---}$

$\hookrightarrow 2^{2047 - 1023} = (2^{1024} \times 1.00\ldots)$

$= \boxed{2^{1024}}$

⑩ 32-bit

ⓐ Fixed-point number :

- **Accuracy:** Fixed accuracy ; loses accuracy for bigger no.s.
- **Precision :** $2^{-(N-1)} = 2^{-31}$ 31 bits
- **Dynamic Range :** $20 \log_{10} \left( \dfrac{\text{Max value}}{\text{Min value}} \right)$

$$= 20 \log_{10} \left( \dfrac{2^{31}-1}{2^{-31}} \right) \approx 20 \log_{10} (2^{62})$$

$$= 20 \times 62 \times 0.301$$

$$= 373.24 \text{ dB.}$$

ⓑ Floating-point no.:

- **Accuracy:** Accurate for bigger no.s.
- **Precision :** 23 bits
- **Dynamic Range :** $20 \log \left( \dfrac{2^{128}}{2^{-149}} \right)$

$$= 20 \log (2^{277}) \approx 1667.54 \text{ dB.}$$

⑪ $(0.752456)_{10} \longrightarrow$ in Q0.4 format.

$\downarrow$

$(0.110000)_2$

In Q0.4 format, $\underset{\uparrow}{\underline{0}} \; \underline{1100}$

sign bit

$(0.1100)_2 \longrightarrow (0.75)_{10}$
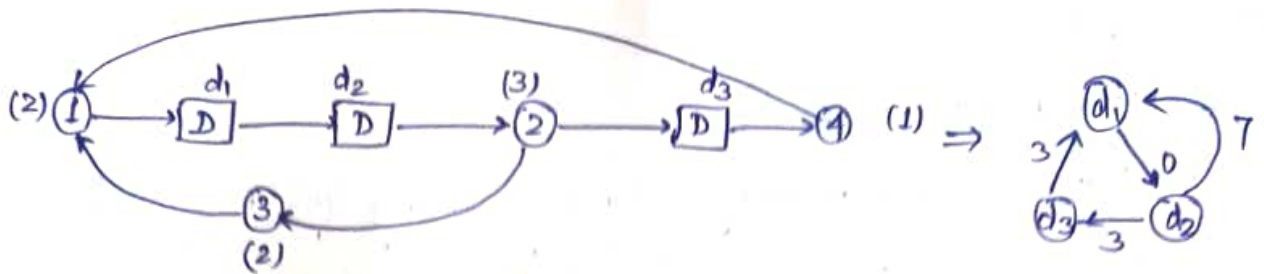
$\therefore$ Error $= 0.752456 - 0.75$

$$= 0.002456$$

**Range:**

Most positive : $1 - 2^{-4} = 0.9375$

Most negative : $0 - 2^{-4} = -0.0625$

Series of matrices : $L^{(m)}, m = 1, 2, ..., d ;$   $d = no. of delays$
   $m = delays$

$$\ell_{ij}^{(m+1)} = \max_{\substack{k \in K \\ \downarrow \\ [1,d]}} \left(-1, \ell_{i,k}^{(1)} + \ell_{k,j}^{(m)}\right) : longest\ computation\ time\ from\ d_i\ to\ d_j$$

<u>L(1):</u>

$\ell_{11} = \ell_{22} = \ell_{33} = -1$  (no path without delay)
$\ell_{12} = 0$ (no computational delay)
$\ell_{13} = -1$
$\ell_{14} = 3 + 2 + 2 = 7$
$\ell_{23} = 3$
$\ell_{31} = 1 + 2 = 3$
$\ell_{32} = -1$

$$\therefore L^{(1)} = \begin{pmatrix} -1 & 0 & -1 \\ 7 & -1 & 3 \\ 3 & -1 & -1 \end{pmatrix}$$
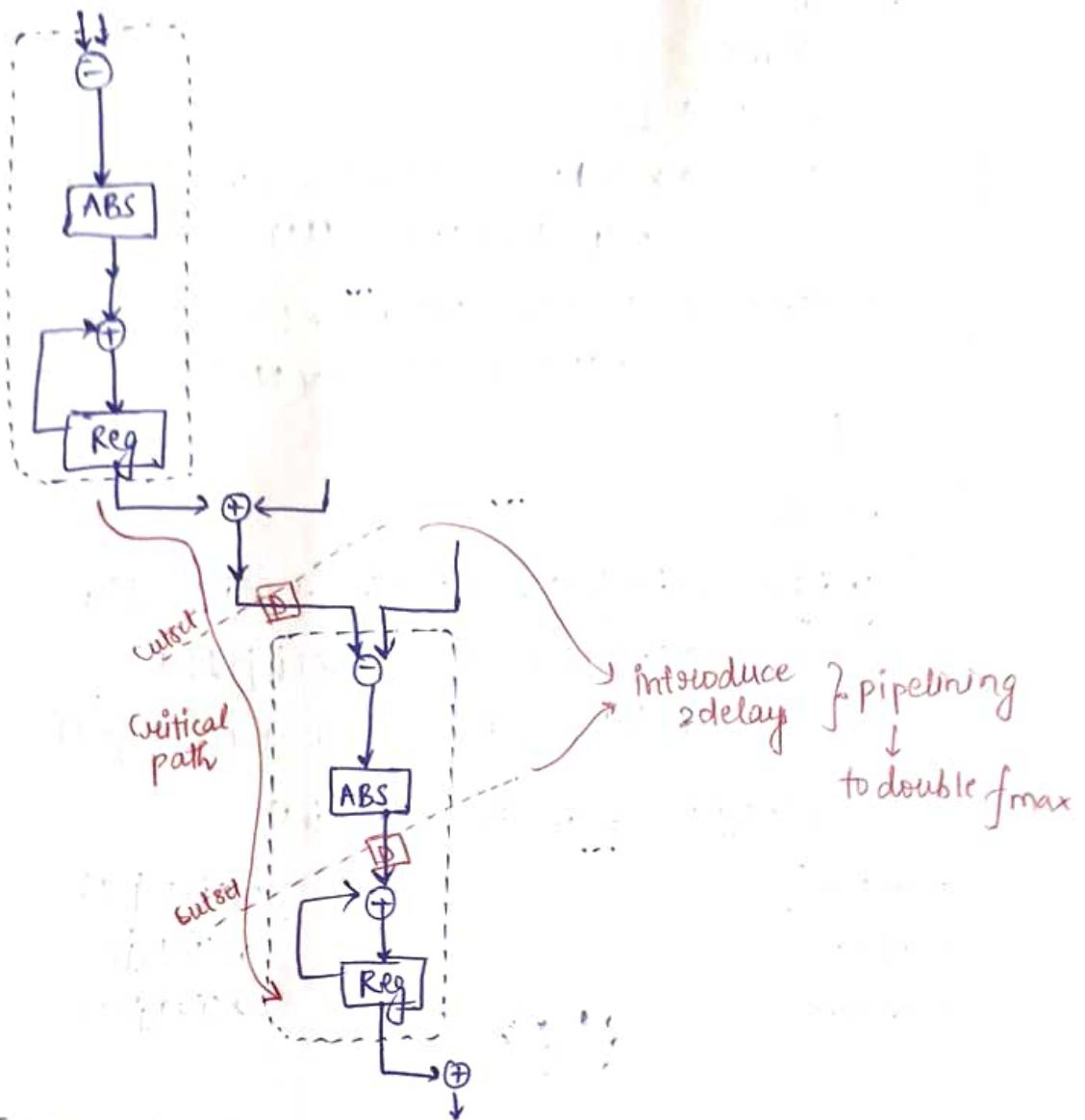
<u>L(2):</u>

$\ell_{11} = 3 + 2 + 2 = 7, \quad \ell_{12} = \cdots$

$$L^{(2)} = \begin{bmatrix} 7 & -1 & 3 \\ 6 & 7 & -1 \\ -1 & 3 & -1 \end{bmatrix}, \quad L^{(3)} = \begin{bmatrix} 6 & 7 & -1 \\ 14 & 6 & 10 \\ 10 & -1 & 6 \end{bmatrix}$$

Iteration bound, $T_\infty = \max \left\{ \frac{7}{2}, \frac{7}{2}, \frac{6}{3}, \frac{6}{3}, \frac{6}{3} \right\}$

$$= 3.5\ ut.$$

⑬ ⓐ $T_S = 5ns$, $T_{AB} = 7ns$, $T_A = 6ns$.



$T_{critical} = T_S + T_B + 2T_A$

$= 24\ ns = T_{clk}$.

Max. frequency: $f_{max} = \dfrac{1}{T_{critical}} = \dfrac{1}{24ns} = 41.67\ MHz$

ⓑ To double the working frequency,

$T_{crit}' = \dfrac{T_{crit}}{2} = 12\ ns$.

↓

Introduce 2 pipelining delays in every critical path.

∴ $T_{critical}' = T_S + T_{ABS} = 12\ ns$.

↪ $f_{max}' = \dfrac{1}{T_{crit}'} = \dfrac{1}{12ns} = 83.33\ ns$

⑭ $y(i,j) = \sum\limits_{m=-1}^{1} \sum\limits_{n=-1}^{1} x(i+m, j+n)$

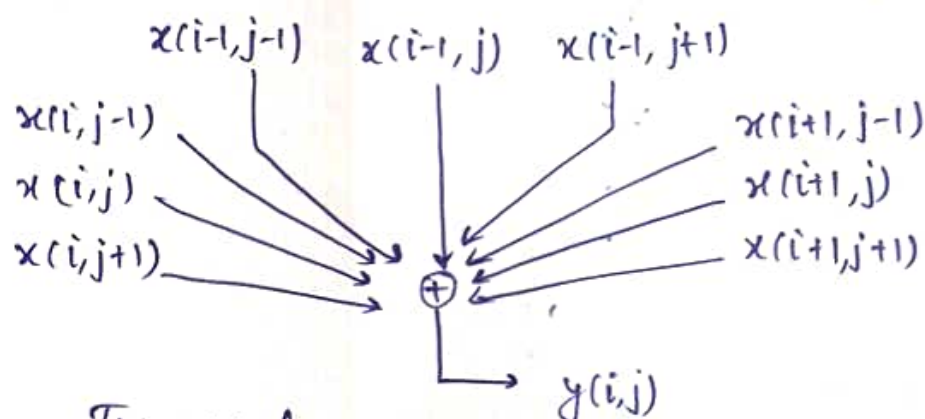$$i, i+m \in [0, W]$$
$$j, j+n \in [0, H]$$

$y(0,0) = x(-1,-1) + x(-1,0) + x(-1,+1) + x(0,-1) + x(0,0) + x(0,1) +$
$\qquad x(1,-1) + x(1,0) + x(1,1)$

$y(0,1) = x(-1,0) + x(-1,1) + x(-1,2) + x(0,0) + x(0,1) + x(0,2) +$
$\qquad x(1,0) + x(1,1) + x(1,2)$

$\vdots$

In general,

$y(i,j) = x(i-1, j-1) + x(i-1, j) + x(i-1, j+1) +$
$\qquad x(i, j-1) + x(i, j) + x(i, j+1) +$
$\qquad x(i+1, j-1) + x(i+1, j) + x(i+1, j+1).$



$x(i-1,j-1) \quad x(i-1, j) \quad x(i-1, j+1)$

$x(i, j-1)$
$x(i, j)$
$x(i, j+1)$

$x(i+1, j-1)$
$x(i+1, j)$
$x(i+1, j+1)$

$\oplus \longrightarrow y(i,j)$

This will happen for all $i, j$, and the actual
dependence graph will be 3D.