

VLSI Signal Processing



Multiple Inputs Multiple outputs
(MIMO)



Single Input Single Output
(SISD)

↳ Using:

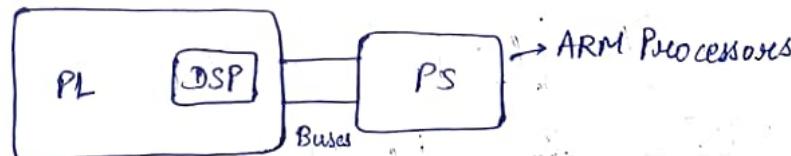
- ① OBC
- ② DSP - Processors, GPGPU - CUDA
- ③ Custom Hardware - FPGA (semi-custom)
- ④ ASIC

↳ Full custom

→ FPGA + dedicated DSP Unit

FPGA + Processor (PS)

↳ PS + PL



↳ Software-Hardware Co-processing

→ Soft core

→ Hard core

• RISC V - for space application

Buses → AXI → used nowadays

PLB

} AMBA (Advanced Microcontroller Bus Architecture)

Major Functionalities of DSP

① Convolution

$$y[n] = \sum_{k=0}^{N-1} x[k] h[n-k]$$

↳ Reflection
 ↳ Dot product / Multiplier
 ↳ Summer

→ DSP: non-terminating units

↳ I/p coming continuously

↳ Not applicable for non-linear system

e.g. Multi-rate: multiple frequencies / sampling rate

↳ Systems with multiple rates

↳ Connected using upsampling / downsampling

↳ Linear time variant / invariant

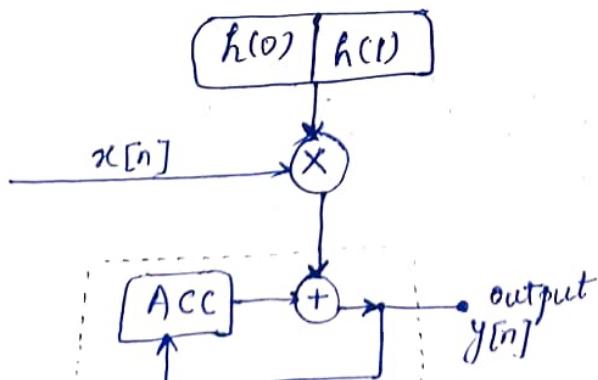
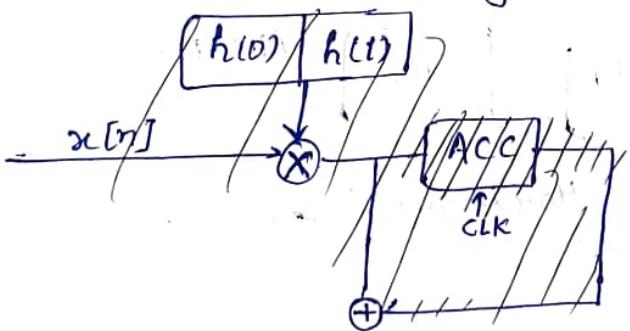
→ Convolution is valid for LTI system.

$$n=0, y[0] = x[0] h[0] \rightarrow \text{one multiplier (zero addition)}$$

$$n=1, y[1] = x[0] h[1] + x[1] h[0] \rightarrow \text{2 multipliers for parallel processing}$$

↳ Addition of previous result

↳ Memory → Accumulator



MAC
Multiply-Accumulate unit

→ Sequential: only one adder & multiplier

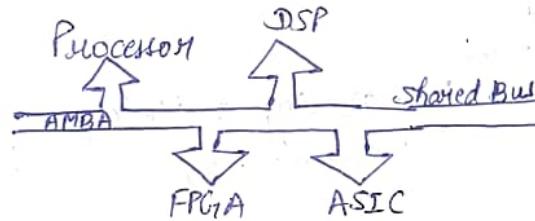
Accumulation → Memory Addition

2-D Convolution:

$$y[m, n] = \sum_{\substack{k=0 \\ l=0}}^{N-1} x[k, l] h[m-k, n-l]$$

↳ Linear independence property

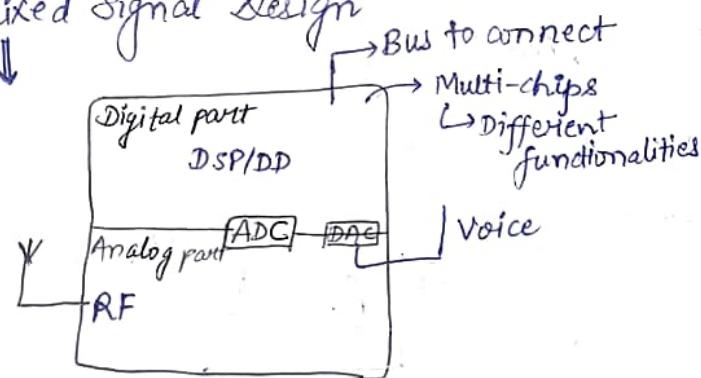
↳ Apply in row, Keeping column fixed.



17-01-2024

Mobile platform

Mixed Signal Design



SOC → System on chip

NoC → Network on chip

- ↳ Has SOC, more no. of processing elements
- ↳ Each connected with dedicated router
- ↳ Each can communicate faster
- ↳ High throughput
- ↳ More area

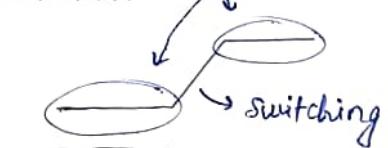
Performance Metrics :

- Power consumption → Provisions to reduce power through algorithm
 - Speed
 - Area
- ↳ To go to sleep mode
↳ Power Retention logic
↳ Power ↓, Area ↑ (complexity↑)

$$\text{Power, } P = P_{\text{sp}} + P_{\text{dp}}$$

(Static) (Dynamic)

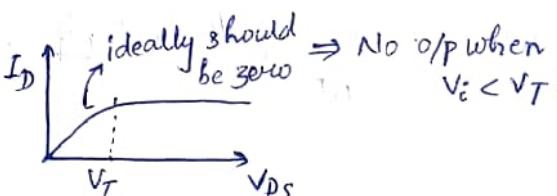
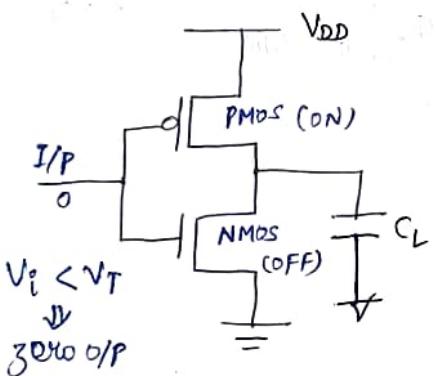
Device in steady-state condition



↳ During switching

↳ cannot be zero, can be reduced.

↳ Zero (ideally), but leakage.



$$P_{\text{dp}} = C_L V_{\text{DD}}^2 f, \quad f: \text{frequency of switching}$$

$$\text{or } C_L V_{\text{DD}} \cdot V_f \\ = C_L V_{\text{DD}} \cdot V_{\text{DD}} f$$



$$\text{Derivation: } i_L = -C \frac{dV}{dt}$$

$$\text{Worst case, } i_L = -C \frac{V_{\text{DD}}}{t_p}$$

$$= C L V_{\text{DD}} f$$

$$\therefore P_{\text{dp}} = V \times i_L = V_{\text{DD}} \times C_L \times V_{\text{DD}} \times f$$

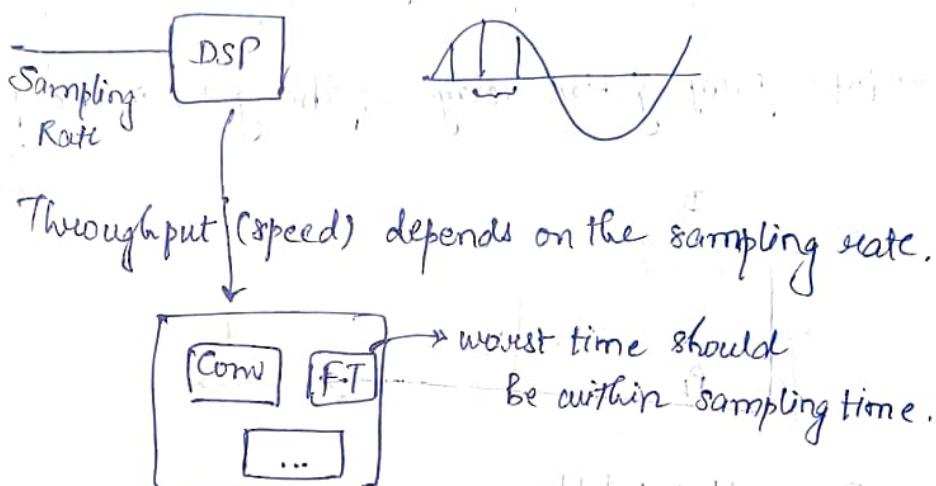
$P_{\text{dp}} \downarrow \Rightarrow f \downarrow$ [But we want high speed]

$\Rightarrow V_{\text{DD}} \downarrow$ (But only upto a level; risk of noise)

Scaling
Profiling

DSP Representation

① Realtime application



② Data driven applications

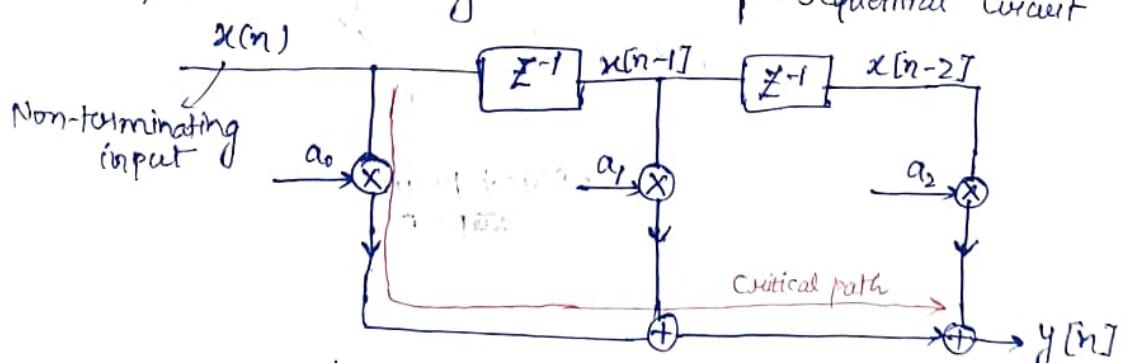
↳ Each node should be functional / excited only when all the data is received.

Eg. FIR filter:

$$y[n] = a_0 x[n] + a_1 x[n-1] + a_2 x[n-2]$$

↳ 3-tap filter $\rightarrow a_0, a_1, a_2$

Representation using structures:



↳ Direct Form I Structure

↳ Hardware Required: 2 adders

3 multipliers

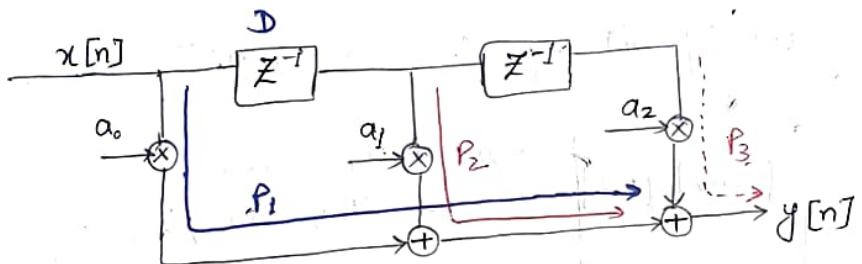
2 delay blocks

Latency → Delay to get the o/p after the application of i/p.

Critical path [in sequential circuit]: Path b/w 2 storage elements
OR, path ~~b/w~~ (longest possible) b/w input & output
[in both Sequential or Combinational circuit]

- ↳ Consists of 2 adders + 1 multiplier in the e.g.
- ↳ Reduce the critical path without reducing functionalities.
- ↳ Total propagation delays of the gates in the path

21-01-2025



$$P_3 = 1 \text{ M} + 1 \text{ Adder}$$

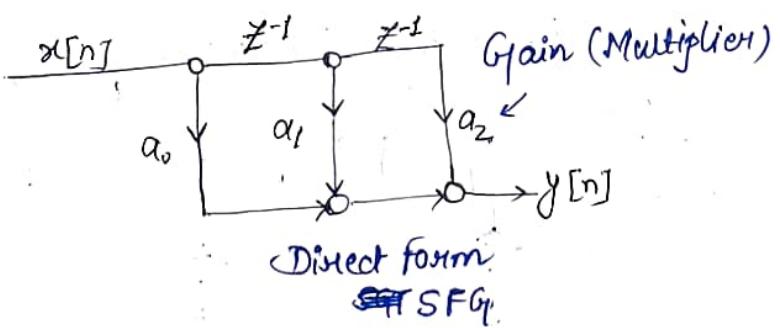
$$P_1 = 1 \text{ M} + 2 \text{ Adder} \quad \checkmark$$

$$P_2 = 1 \text{ M} + 2 \text{ Adder} \quad \checkmark \quad \} \text{ worst possible}$$

→ Block Diagram Representation

Signal Flow Graph

↳ vertex, node



Transposition Structure

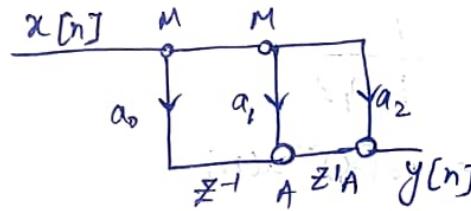
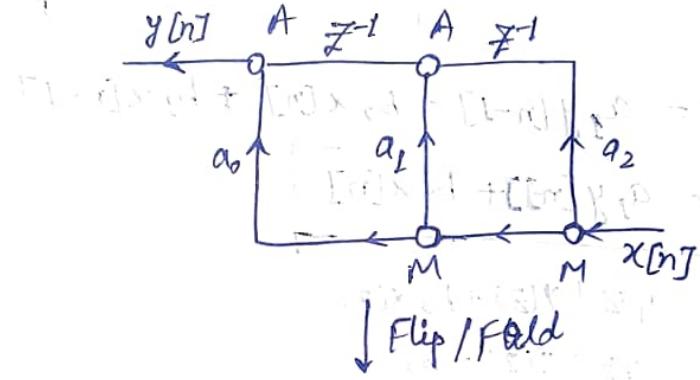
(Direct Form II)

Input → Output

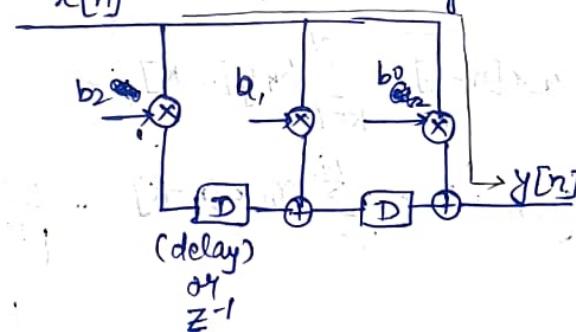
Incoming → Outgoing

Accumulation

→ ⊕ ←



↓ Block diagram



→ Data Broadcast
Structure

$$y[n] = a_0 z^{-1} + a_1 z^{-2} + a_2 z^{-3}$$

Data Flow Graph (DFG)

↳ DSP runs on Data Driven Property

○ Node/Actor

Multiplexer

Execution / fire

$$\text{IIR Filter: } y[n] = \sum_{k=0}^M a_k y[n-k] + \sum_{k=0}^{N-1} b_k x[n-k] \quad (M=N \text{ or } M < N)$$

1st order IIR filter: $y[n] = a_1 y[n-1] + b_0 x[n] + b_1 x[n-1]$

$$y[n] = a_1 y[n-1] + b_0 x[n]$$

$$\Rightarrow Y(z) = a_1 z^{-1} Y(z) + b_0 X(z)$$

$$\Rightarrow H(z) = a_1 z^{-1} H(z) + b_0$$

$$\Rightarrow H(z) = \frac{b_0}{1 - a_1 z^{-1}}$$

$$\Rightarrow h(n) = b_0 (a_1)^n u(n).$$

23-01-2024

IIR Filter:

$$y[n] = \sum_{k=0}^M a_k x[n-k] - \sum_{k=1}^N b_k y[n-k]$$

$$\sum_{k=0}^N b_k y[n-k] = \sum_{k=0}^M a_k x[n-k]$$

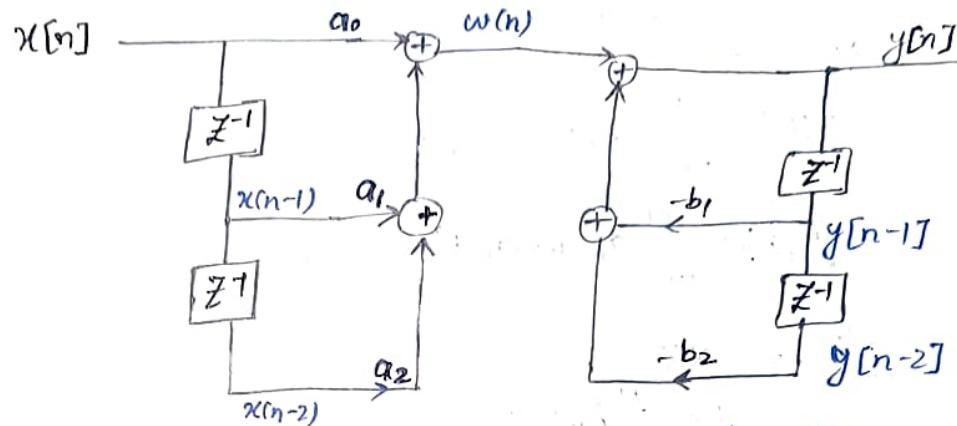
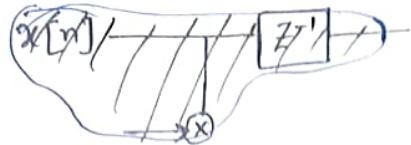
... LCCDE
(Linear Constant Coefficient Difference Equations)

- $M=N$
- $M>N$
- $M<N$

Simple form: $y[n] = b_1 y[n-1] + a_1 x[n]$

for $M=2, N=2,$

$$\begin{aligned} y[n] &= (a_0 x[n] + a_1 x[n-1] + a_2 x[n-2]) \\ &\quad - b_1 y[n-1] - b_2 y[n-2] \\ \Rightarrow & \boxed{b_0 = 1} \end{aligned}$$



↳ Direct Form I
↳ Non-canonical

Canonical form: No. of delay elements is equal to the order of the system.

↳ Hardware Required:

- 4 adders
- 4 delay elements (DFFs)
- 5 multipliers

↳ Complex structure

↳ Reduce complexity: Combine delay elements for $w(n)$ and $y(n)$.

$$H(z) = \frac{Y(z)}{X(z)} = \frac{Y(z)}{W(z)} \cdot \frac{W(z)}{X(z)} \\ = H_1(z) \cdot H_2(z)$$

$$H_1(z) = \frac{Y(z)}{W(z)}$$

$$\Rightarrow Y(z) = W(z) \cdot H_1(z)$$

$$H_2(z) = \frac{W(z)}{X(z)}$$

$$\Rightarrow W(z) = H_2(z) \cdot X(z)$$

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

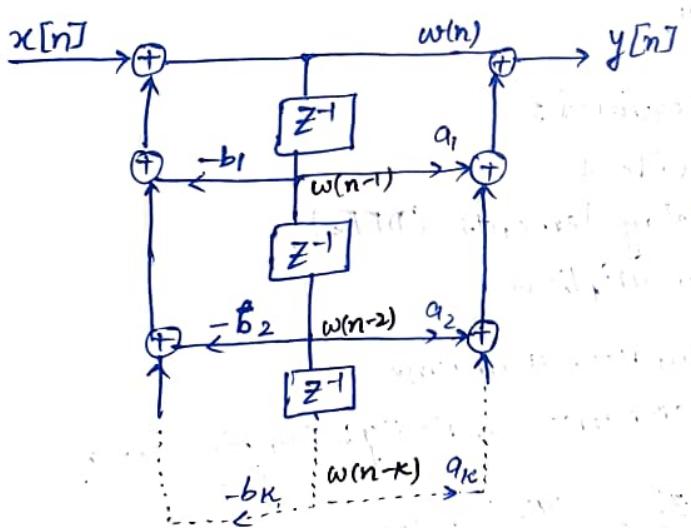
$$\frac{Y(z)}{W(z)} = \sum_{k=0}^M a_k z^{-k}$$

$$\Rightarrow Y(z) = \sum_{k=0}^M a_k z^{-k} \cdot W(z)$$

$$y[n] = \sum_{k=0}^M a_k w[n-k] \quad \dots \textcircled{2}$$

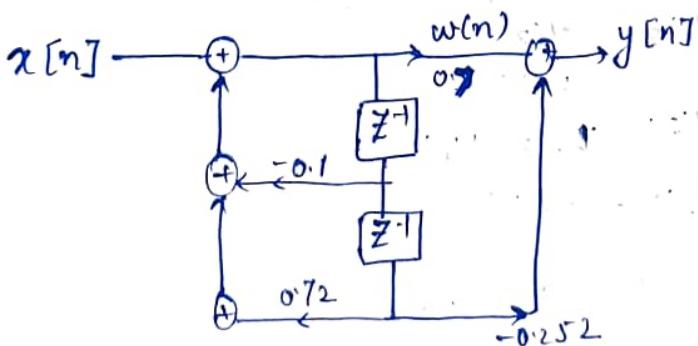
$$\underline{W(z)} = \frac{X(z)}{1 + \sum_{k=1}^N b_k z^{-k}}$$

$$\Rightarrow w[n] = x[n] - \sum_{k=1}^N b_k w[n-k] \quad \dots \textcircled{1}$$



- Half of form I)
- Less delay elements
- Canonical
- Direct form II

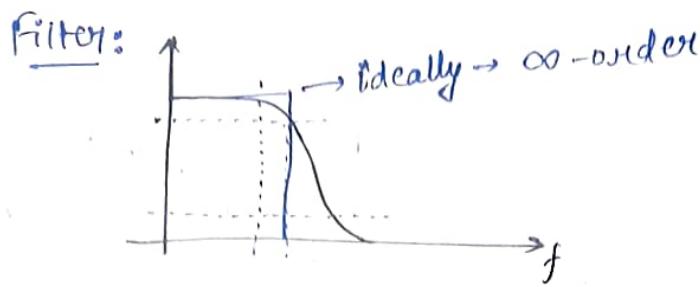
Eg. $y[n] = (-0.1 y[n-1] + 0.72 y[n-2]) + (0.7 x[n] - 0.252 x[n-2]).$



- 2 delays
- Filter

$$w(n) = x(n) - 0.1 w(n-1) + 0.72 w(n-2)$$

$$y(n) = 0.7 w(n) - 0.252 w(n-2).$$

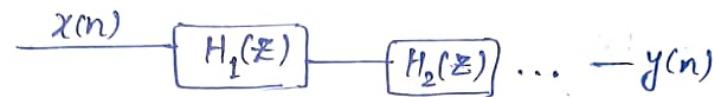


Cascade form structure:

$$H(z) = H_1(z) \cdot H_2(z) \cdots H_N(z)$$

→ To reduce quantization error

(28-01-2025)

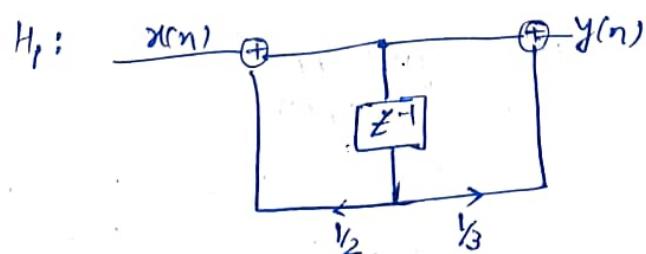
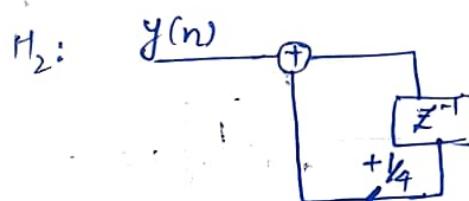


$$y[n] = \frac{3}{4}y[n-1] - \frac{1}{3}y[n-2] + x[n] + \frac{1}{3}x[n-1]$$

$$\begin{aligned} H(z) &= \frac{1 + \frac{1}{3}z^{-1}}{1 - \frac{3}{4}z^{-1} + \frac{1}{3}z^{-2}} = \frac{1 + \frac{1}{3}z^{-1}}{(1 - \frac{1}{2}z^{-1})(1 - \frac{1}{4}z^{-1})} \\ &= H_1(z) \cdot H_2(z) \end{aligned}$$

$$H_1(z) = \frac{(1 + \frac{1}{3}z^{-1})}{(1 - \frac{1}{2}z^{-1})}$$

$$H_2(z) = \frac{1}{(1 - \frac{1}{4}z^{-1})}$$



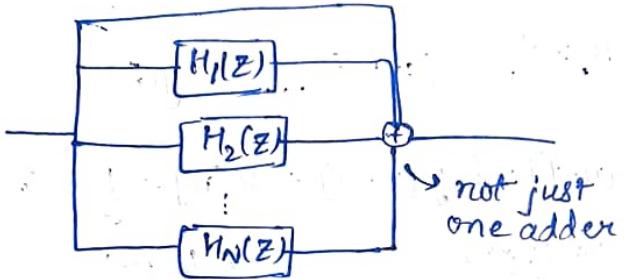
Parallel form:

$$H(z) = C + \sum_{k=1}^N \frac{c_k}{1-p_k z^{-1}} \quad (\text{Partial Expansion})$$

→ for IIR as well as FIR

$$H(z) = C + \frac{c_1}{1-p_1 z^{-1}} + \frac{c_2}{1-p_2 z^{-1}} + \dots + \frac{c_N}{1-p_N z^{-1}}$$

$$H(z) = C + H_1(z) + H_2(z) + \dots + H_N(z)$$



Linear Phase System

Type-I: $h(n) = h(n-N)$ [symmetric]
 $h(N-n)$ (for causal), N : Even or odd

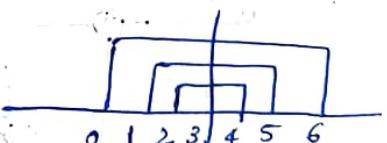
$$H(z) = \sum_{k=0}^{N-1} h(k) z^{-k}$$
 ↳ Delay elements can be shared (only 4 required)

$$h[0] = h[+6]$$

$$h[1] = h[+5]$$

$$h[3] = h[3]$$

$$h[2] = h[4]$$



$$\# SQNR = \frac{\text{Input Signal (RMS)}}{\text{Noise (RMS)}}$$

(Signal to
Quantization
Noise Ratio)

$= 6.02 + 1.76 \rightarrow$ For every one bit improvement in the original, SQNR improves by 6 dB.

Types:

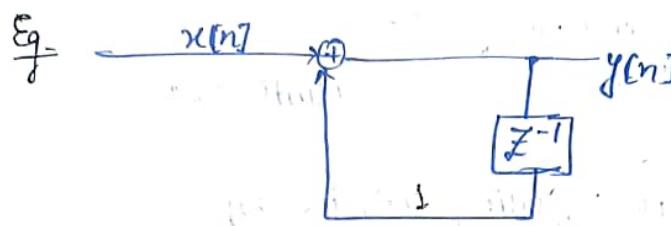
→ Analog to digital (Input) Q. error

→ Multiply & accumulation (Product) error

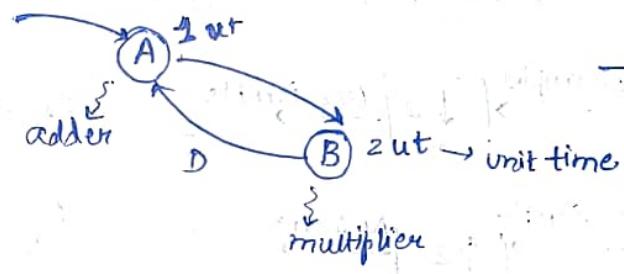
- Truncation introduces more error.
- Rounding less error.

Data Flow Graph

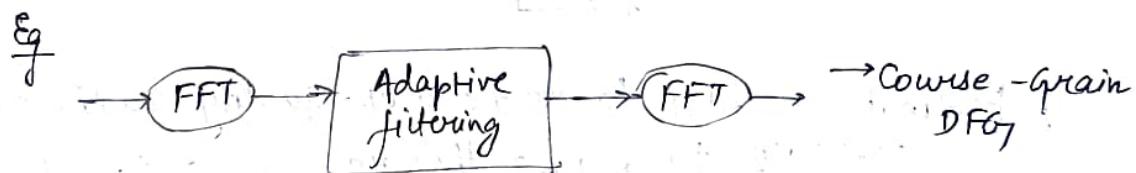
- Nodes → Represent computations (multiply & adder).
- Directed edges → Represent data paths.
- Edge → Delay.



$$y[n] = y[n-1] + x[n]$$



→ fine DFG
 ↳ Only adder and multiplier computations

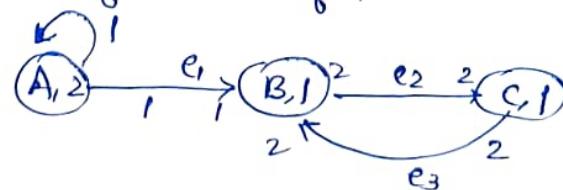


30-01-2025

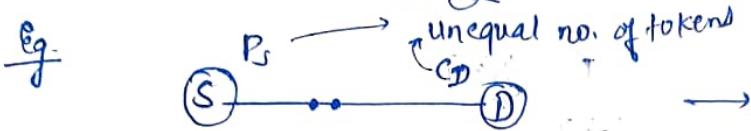
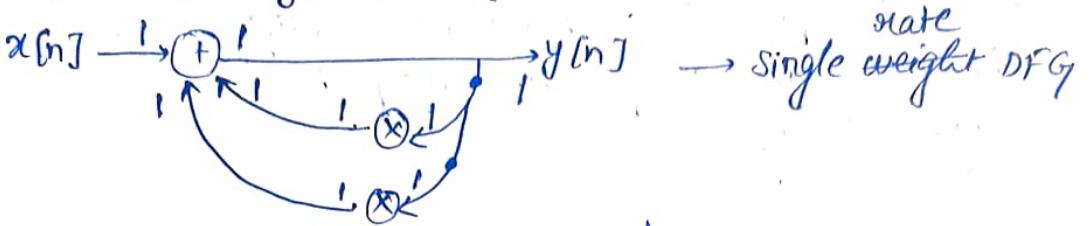
→ Exposes hidden concurrency among different parts.

Synchronous Dataflow Graph:

- ↳ No. of data samples produced or consumed by each node in each execution is specified.
- ↳ Each tasks operate on pre-defined no. of samples and produce a fixed no. of o/p's.

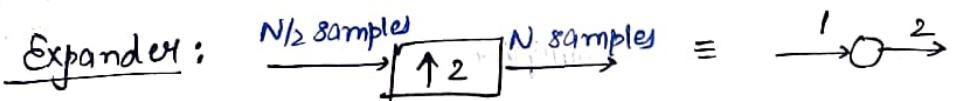
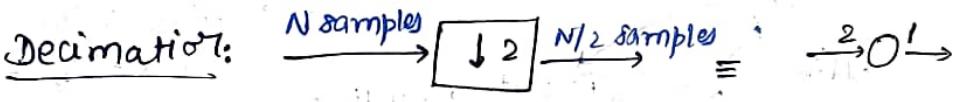


$$\text{Eq. } y[n] = x[n] + a_1 y[n-1] + a_2 y[n-2]$$



$$f_s P_S = f_D C_D : \text{fixing relationship}$$

Single-Rate and Multi-rate SDFGs:



Eg Format converter to store a CD-quality audio sampled at 44.1 kHz to a DAT that records a signal at 48 kHz.

$$\frac{f_{\text{out}}}{f_{\text{CD}}} = \frac{480}{441} = \frac{160}{147} = \frac{2}{1} \times \frac{4}{3} \times \frac{5}{7} \times \frac{9}{7}$$

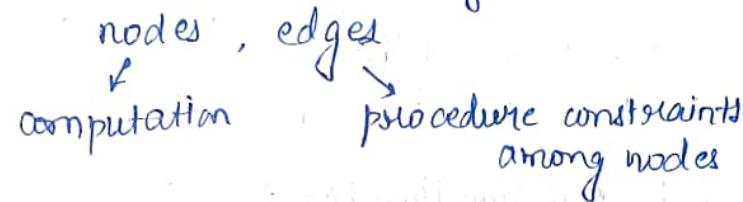


Homogeneous SDFG:

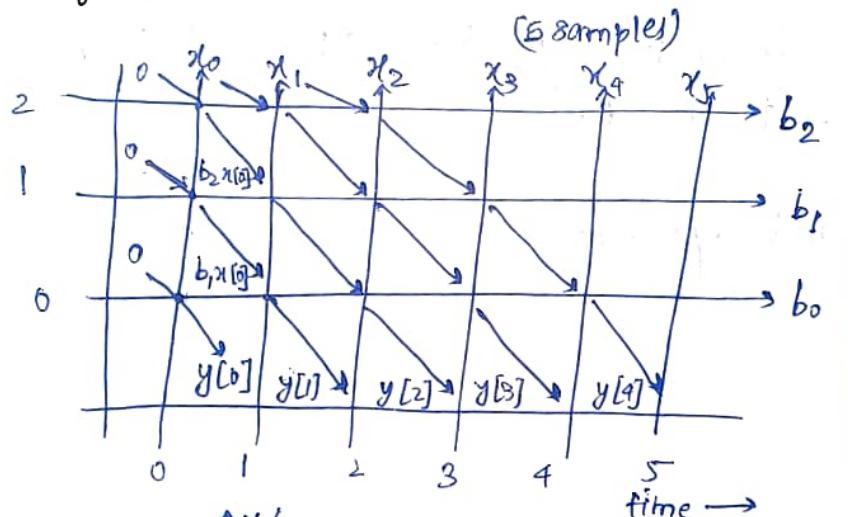
↳ Single-rate graph where each node produces one token

Dependency Graph:

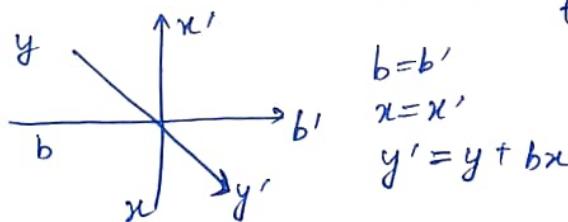
↳ Directed graph which shows the dependence of the computation in an algorithm.



Eq. $y[n] = b_0 x[n] + b_1 x[n-1] + b_2 x[n-2]$



$$\begin{aligned}
 y[0] &= b_0 x[0] + b_1 x[-1] \\
 &\quad + b_2 x[-2] \\
 y[1] &= b_0 x[1] + b_1 x[0] \\
 &\quad + b_2 x[-1] \\
 y[2] &= b_0 x[2] + b_1 x[1] \\
 &\quad + b_2 x[0]
 \end{aligned}$$



Implement 4×4 Matrix multiplication using ~~systolic~~ table in Verilog.

Transformations on a Dataflow Graph:

- ↳ To minimize critical path delay or no. of registers.
- Retiming
- Folding
- Unfolding
- Pipelining
- Parallel processing

Performance Measures for Design Tradeoffs

- ↳ Iteration period: Time taken by system to compute all the operations in one iteration of an algorithm.
 - ↳ Measured in terms of the no. of clock cycles required to compute one o/p sample.
- ↳ sampling-period: Average time b/w two successive data samples.
- ↳ Latency: Time delay for algorithm to produce output in response to an input.
- ↳ Power
 - ↳ static and dynamic.

Number Representation

- Fixed point
- Floating point

Fixed Point Representation

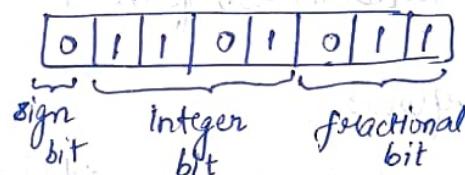
- ↳ Computer assumes any bits/no. as integer.
- ↳ An integer X is used to represent a real no. x , with

$$N = m + n + 1 \text{ bits}$$

N : wordlength

m : no. of integer bits (excluding sign bit)

n : no. of fractional bits



Q-format: X sometimes called $Q_{m,n}$ or Q_m number.

Conversion to and from fixed point:

Real to fixed-point: $x := \text{Round}(x \cdot 2^n)$

fixed-point to Real: $x := x \cdot 2^{-n}$

e.g. $x = 13.4 \rightarrow$ Represent using $Q_{4,3}$ format.

$$x = \text{Round}(13.4 \cdot 2^7) = 107 (= 01101011_2)$$

| Sign | 2^3 | 2^2 | 2^1 | 2^0 | 2^{-1} | 2^{-2} | 2^{-3} |
|------|-------|-------|-------|-------|----------|----------|----------|
| 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 |

$8 \quad 4 \quad 2 \quad 1 \quad 0.5 \quad 0.25 \quad 0.125$

$$x = \frac{107}{8} = 13.375 \neq 13.4$$

$$\rightarrow \text{Dynamic Range} = \frac{\text{Max value}}{\text{Min value}} = \frac{2^{N-1}-1}{1} \quad (\text{signed integer format})$$

Range v8. Resolution for fixed-point no.s:

A Q_{m.n} fixed point no. can represent no. in the range

$$[-2^m, 2^m - 2^{-n}]$$

Fractional part resolution = $2^{-(n-1)}$

Eg. To represent x in signed 8-bit variable.

$$-28.3 < x < 17.5$$

$m=5$ to represent the integer part:

$$(2^4 = 16 < 28.3 < 32 = 2^5)$$

$n = 8-1-m = 2$ bits left for the fractional part.

x should be stored in Q_{5.2} format.

Eg. Addition with overflow:

Two nos in Q_{4.3} format are added:

$$\begin{array}{r}
 x = 12.25 \Rightarrow x = 98 \quad 01100010 \\
 y = 14.75 \Rightarrow y = 118 \quad + 01110110 \\
 \hline
 z = x + y = 216 \quad \quad \quad 11011000
 \end{array}$$

This no. is out of range and will be represented as

$$216 - 256 = -40 \Rightarrow z = -5.0$$

$$(2^8)$$

06-02-2025

Fixed-Point Multiplication and Division

If operands and results are in same Q-format:

$$z = xy \Leftrightarrow z = (x \cdot y) / 2^n$$

$$z = x/y \Leftrightarrow z = (x \cdot 2^n) / y$$

→ Double wordlength is needed for the intermediate result.

→ Division (multip.) by 2^n is implemented by a right (left)-shift by n -bits.

→ Lowest bits in the result are truncated (round-off noise).

→ Risk of overflow.

Eg. Two nos in Q5.2 format:

$$x = 6.25 \Rightarrow x = 25$$

$$y = 4.75 \Rightarrow y = 19$$

$$x \cdot y = 47.5$$

$$\Rightarrow z = 47.5 / 2^2 = 11.8 \Rightarrow z = 29.5$$

In Q5.3 format:

$$x = 6.25 \Rightarrow x = 50$$

$$y = 4.75 \Rightarrow y = 38$$

$$\Rightarrow z = \frac{1900}{2^3} = 237.5 \Rightarrow z = 29.625$$

With different Q-format:

In general, multiplication of 2 fixed-point nos. Q_{m₁,n₁} and Q_{m₂,n₂} gives an intermediate result in the format:

$$Q_{m_1 + m_2, n_1 + n_2}$$

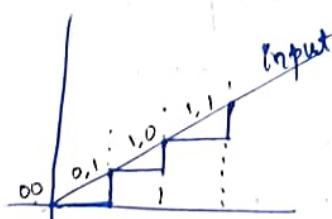
which may be then right-shifted n₁+n₂-n₃ steps and stored in the format

$$Q_{m_3, n_3}$$

Quantization error in representing H(z):

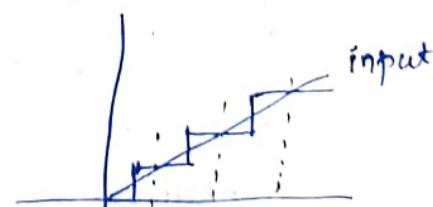
- Input Q.E.
- Product Q.E.
- Output Q.E.

Truncation:



$$\begin{aligned}\hookrightarrow \text{worst-case error} &= 2^{-(B-1)} \\ &= 2(\text{Rounding error})\end{aligned}$$

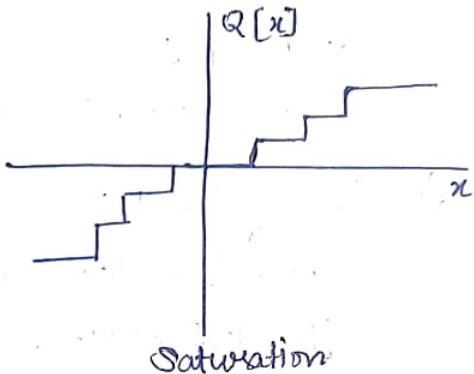
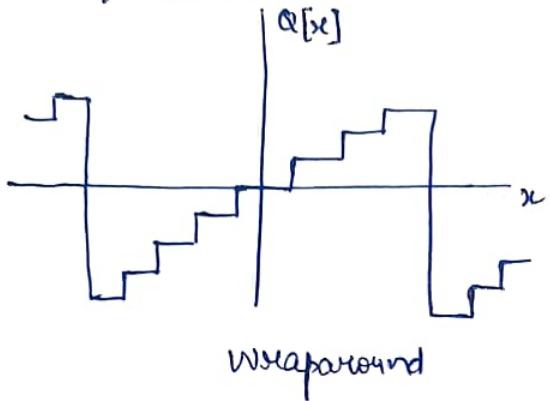
Rounding:



SIMULINK → HDL TOOLBOX
MATLAB → fixed-point toolbox

Overflow

- Wraparound
- Saturation



Schwartz's inequality:

$$|y[n]| \leq \sqrt{\sum_{n=-\infty}^{\infty} h^2[n] \sum_{n=-\infty}^{\infty} x^2[n]} \rightarrow \text{To find upper bound on the o/p values of a LTI system.}$$

Fixed Point Quantization

- ↳ The fractional B-bit 2's complement no. representation evenly distributes 2^B quantization levels b/w -1 and $1 - 2^{-(B-1)}$.
- ↳ Spacing b/w quantization levels : $\frac{2}{2^B} = 2^{-(B-1)} = \Delta_B$.
- ↳ Any signal value falling b/w two levels is assigned to one of the two levels.

$$e = \underbrace{Q[x] - x}_{\text{quantization notation: } X_Q} : \text{quantization error.}$$

Floating Point Number

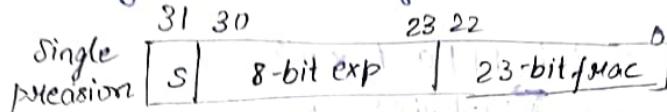
- ↳ $F \times r^E$, F: fraction
E: exponent
r: radix ($= 10$ for decimal nos.)

In 32-bit single precision floating-point representations:

- ↳ Most significant bit : sign bit (s) $\rightarrow 0$ for +ve
 1 for -ve
- ↳ Following 8 bits : biased exponent (E)
- ↳ Remaining 23 bits : fraction (F).

Single precision Floating Point No.

- ↳ Easy (and lazy) way of dealing with scaling problem.
- ↳ 32-bit single precision floating point



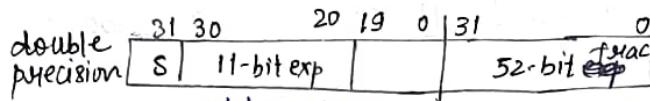
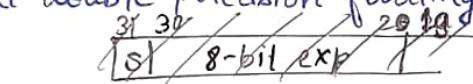
$$x = -1^s \times 2^{\text{exp}-127} \times 1.\text{frac}$$

$$1.175 \times 10^{-38} < |x| < 1.7 \times 10^{38}$$

- ↳ MSB is sign bit (same as fixed point).
- ↳ 8-bit exponent in bias-127 integer format (i.e., add 127 to it).
- ↳ 23-bit to represent only the fractional part of the mantissa.

The MSB of the mantissa is ALWAYS '1', therefore it is not stored.

- ↳ 64-bit double precision floating point:



odd register (e.g., A5) even register (e.g., A4)

$$x = -1^s \times 2^{\text{exp}-1023} \times 1.\text{frac}$$

$$2.2 \times 10^{-308} < |x| < 1.7 \times 10^{308}$$

- ↳ MSB is sign-bit (same as fixed point).

- ↳ 11-bit exponent in bias-1023 integer format (i.e., add 1023 to it).

- ↳ 52-bit to represent only the fractional part of the mantissa.

The MSB of the mantissa is ALWAYS '1', therefore it is not stored.

Short and Long IEEE 754 Formats: Features

| Feature | Single/Short | Double/Long | Coded no. $x = (-1)^s \cdot 2^e \cdot f$ |
|---------------------------|--|---|---|
| Word width in bits | 32 | 64 | |
| Significand in bits | 23 + 1 hidden | 52 + 1 hidden | |
| Significand range | $[1, 2 - 2^{-23}]$ | $[1, 2 - 2^{-52}]$ | |
| Exponent bits | 8 | 11 | |
| Exponent bias | 127 | 1023 | |
| Zero (± 0) | $e + \text{bias} = 0, f = 0$ | $e + \text{bias} = 0, f = 0$ | |
| Subnormal | $e + \text{bias} = 0, f \neq 0$ | $e + \text{bias} = 0, f \neq 0$ | |
| Infinity ($\pm \infty$) | Represents $\pm 0 \cdot f \times 2^{+126}$ | Represents $\pm 0 \cdot f \times 2^{-1022}$ | |
| Not-a-no. (NaN) | $e + \text{bias} = 255, f = 0$ | $e + \text{bias} = 2047, f = 0$ | |
| min | $2^{-126} \approx 1.2 \times 10^{-38}$ | $2^{-1022} \approx 2.2 \times 10^{-308}$ | |
| max | $2^{128} \approx 3.4 \times 10^{38}$ | $2^{104} \approx 1.8 \times 10^{308}$ | |

Eg. Represent -13.8 using IEEE 754 32-bit single precision floating point representation.

$$(13.8) \rightarrow (1101.11001)$$

$$1101.11001 = 1.10111001 \times 2^3$$

$\xrightarrow{\text{shifted right}}$

Actual exponent = 3

Biased exponent = $3 + 127 = 130 = (100000010)_2$

sign of fraction / Mantissa (8) = -ve = 1

| S | Biased exponent | Mantissa/fraction |
|---|-----------------|--------------------------------|
| 1 | 100000010 | 101110010000000000000000000000 |

Eg. 32-bit pattern:

1 1000 0001 011000 00000 00000 00000 00000

s=1 (-ve)

Biased exponent = 1000 0001 = 129

Actual exponent = $129 - 127 = 2$

$$1.01100000000000000000000000 = 1 + 1 \times 2^{-2} + 1 \times 2^{-3}$$

$$1.01100000000000000000000000 = 1.375$$

$$1.01100000000000000000000000 = 1 - 1.375 \times 2^2 = -5.5$$

single/short

double/long

Ordinary no.

et+bias $\in [1, 254]$

et+bias $\in [1, 2046]$

ee $\in [-126, 127]$

ee $\in [-1022, 1023]$

represents $s.f \times 2^e$

represents $s.f \times 2^e$

Example:

$$\begin{cases} 0 & 0.(8).0 \\ 1 & 0.(8).0 \end{cases} 0...23.0 = 0$$

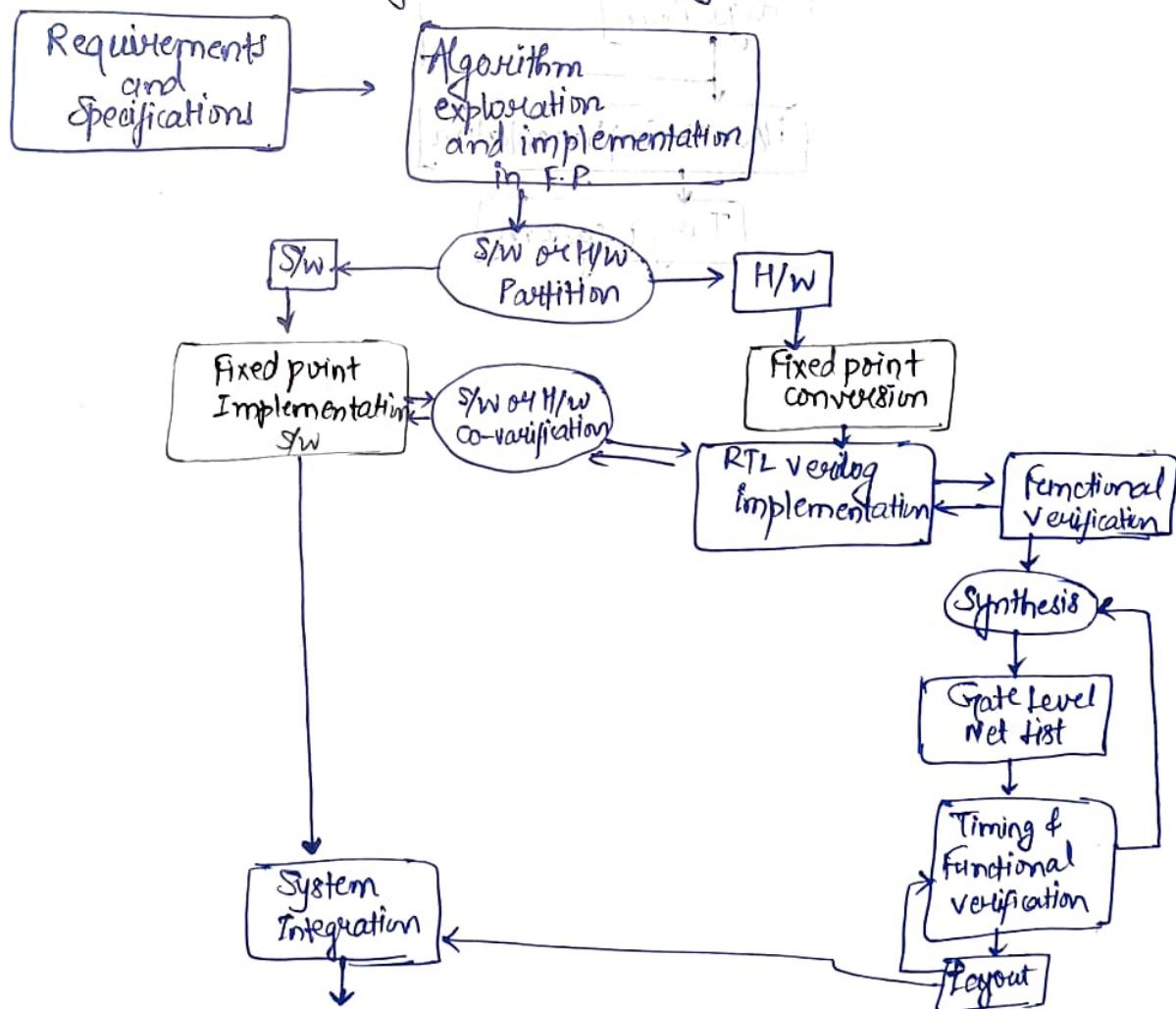
$$\begin{cases} 0 & 1.(8).1 \\ 1 & 1.(8).1 \end{cases} 0...23.0 = \text{infinity}$$

$$\begin{cases} 0 & 1.(8).1 \\ 1 & 1.(8).1 \end{cases} 0...23.0 = -\text{infinity}$$

$$\begin{cases} 0 & 1.(8).1 \\ 1 & 1.(8).1 \end{cases} 0000010.(17).0 = \text{NaN}$$

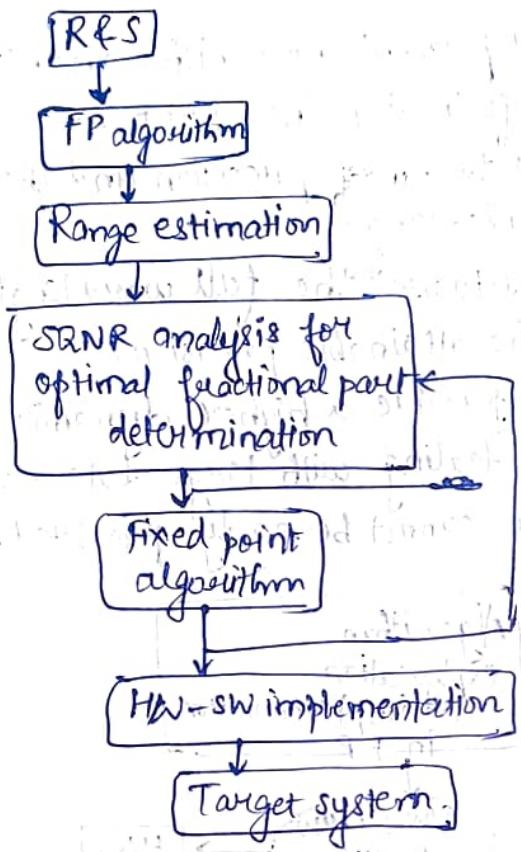
$$\begin{cases} 1 & 1.(8).1 \end{cases} 001000100010010101010 = \text{NaN}$$

- Floating-point DSPs represent nos with a mantissa and an exponent, nowadays following the IEEE 754.
- The mantissa dictates the no. of precision and the exponent controls its dynamic range.
- Nos are scaled so as to use the full word length available, hence maximizing the attainable precision.
- Floating-point nos provide a higher dynamic range, which can be essential when dealing with large data sets and with data sets whose range cannot be easily predicted.



FP to fixed point conversion

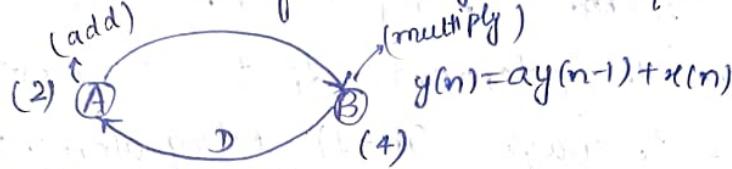
- Start with the algorithm (FP).
- Estimate the range.
- Determine n.
- The fractional part m requires detailed analysis of signal quantization noise.



Loop Bound and Iteration Bound

Iteration Bound (of DSP program) is the lower bound on the iteration of DSP program regardless of the amount of computing resources available as fundamental limit on the minimum sampling period.

Loop bound: Time for a loop to complete one iteration.



Precedence constraint:

inter-iteration : delay : $B_n \rightarrow A_{n+1}$

intra-iteration : no-delay : $A_n \rightarrow B_n$

iteration : $A_n \rightarrow B_n \Rightarrow A_{n+1} \rightarrow B_{n+1}$

Loop: same begin & end node
 Iteration of node
 ↳ execution of node exactly once
 Iteration of DFG
 ↳ execution of each node exactly once

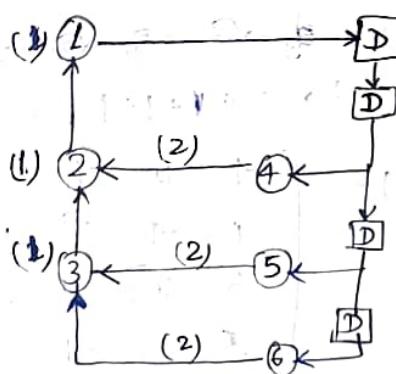
Loop bound of i^{th} loop : $\frac{t_i}{w_i}$, t_i : loop computation time
 w_i : no. of delay in the loop.



$$\text{Loop bound} = \frac{G^k}{2} = 3.$$

↳ 1 iteration of loop completed in 3 u.t.
 [Iteration bound = 3 u.t.]

Eg:



$$1 \rightarrow 4 \rightarrow 2 \rightarrow 1 : \text{l.p.} = \frac{4}{2} \text{ u.t.}$$

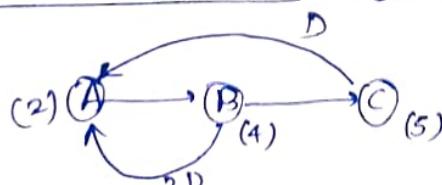
$$1 \rightarrow 5 \rightarrow 3 \rightarrow 2 \rightarrow 1 : \text{l.p.} = \frac{5}{3} \text{ u.t.}$$

$$1 \rightarrow 6 \rightarrow 3 \rightarrow 2 \rightarrow 1 : \text{l.p.} = \frac{5}{4} \text{ u.t.}$$

Critical loop: Loop which has max^m loop bound.

Iteration bound : Max^m loop bound:

$$T_{\infty} = \max \left\{ \frac{t_i}{w_i} \right\}$$



$$\text{l.b.} = 11 \text{ u.t. } (A \rightarrow B \rightarrow C \rightarrow A)$$

Longest path Matrix Algorithm "Iteration Bound" (LPM Method)

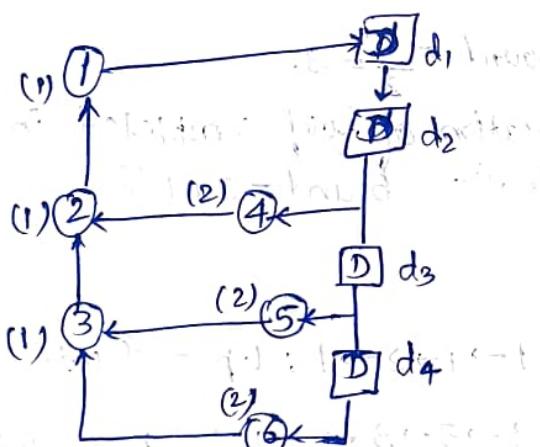
$L^{(m)}$, $m=1, 2, \dots, d$, d : no. of delays in the DFG
 m : delays.
 ↳ series of matrix.

→ value of $l_{i,j}^{(m)}$: Longest computation time of all paths from delay element d_i to delay element d_j that passes through $(m-1)$ delay elements; if no such path, it is set to -1.

higher-order matrices: $l_{i,j}^{(m+1)} = \max_{k \in K} (-1, l_{i,k}^{(m)} + l_{k,j}^{(m)})$

\downarrow neither should be -1 ↳ -1, if so.

$$T_{\infty} = \max_{i, m \in \{1, 2, \dots, d\}} \left\{ \frac{l_{i,i}^{(m)}}{m} \right\}$$



- $m-1=0$ & consider zero delays
- $l_{11}^{(1)}$: no path without delay → -1
- $l_{12}^{(1)}$: no computational node → 0
- $l_{13}^{(1)}$: no path without delay → -1
- ⋮
- $l_{31}^{(1)}$: $d_3 \rightarrow 5 \rightarrow 3 \rightarrow 2 \rightarrow 1 \rightarrow d_1$
 $\hookrightarrow 2+1+1+1=5$

$$\therefore L^{(1)} = \begin{bmatrix} -1 & 0 & -1 & -1 \\ 4 & -1 & 0 & 1 \\ 5 & 0 & -1 & 0 \\ 5 & -1 & -1 & -1 \end{bmatrix}$$

$$\underline{L^2}: l_{2,3}^{(2)} = \max_{K \in \mathbb{O}} (-1, l_{2,K}^{(1)} + l_{K,3}^{(1)}) \rightarrow \text{no valid case} \rightarrow -1$$

$$l_{3,2}^{(2)} = \max_{K \in \mathbb{I}} (-1, l_{3,K}^{(1)} + l_{K,2}^{(1)}) \rightarrow \text{valid for } K=1 \rightarrow 5$$

$$L(2) = \begin{bmatrix} 4 & -1 & 0 & -1 \\ 5 & 4 & -1 & 0 \\ 5 & 5 & -1 & -1 \\ -1 & 5 & -1 & -1 \end{bmatrix} \quad l_{i,j}^{(2)} = \max_{K \in \mathbb{K}} (-1, l_{i,K}^{(1)} + l_{K,j}^{(1)}), \\ l_{i,K}, l_{K,j}^{(1)} \neq -1.$$

$$l_{1,1}^{(2)} = \max_{K \in \mathbb{O}} (-1, l_{1,K}^{(1)} + l_{K,1}^{(1)}) \rightarrow K=1: \text{not valid} \\ K=2: \text{valid} \rightarrow 4$$

L₃:

$$l_{i,j}^{(3)} = \max_{K \in \mathbb{K}} (-1, l_{i,K}^{(1)} + l_{K,j}^{(2)}) \rightarrow \text{from L}(1) \text{ and L}(2).$$

$$L(3) = \begin{bmatrix} 5 & 4 & -1 & 0 \\ 8 & 5 & 4 & -1 \\ 9 & 5 & 5 & -1 \\ 9 & -1 & 5 & -1 \end{bmatrix}$$

L₄:

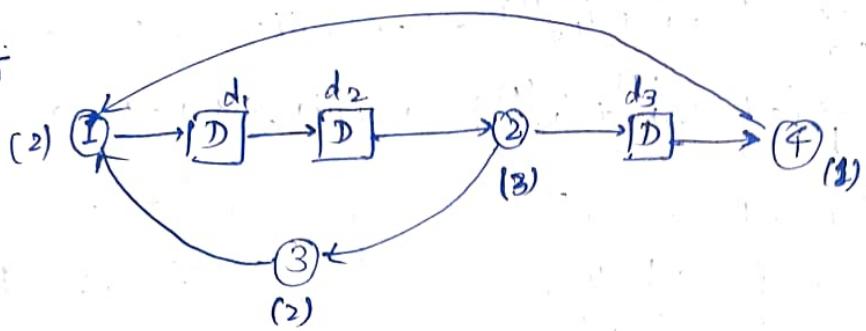
$$L(4) = \begin{bmatrix} 8 & 5 & -4 & -1 \\ 9 & 8 & 5 & 4 \\ 10 & 9 & 5 & 5 \\ 10 & 9 & -1 & 5 \end{bmatrix}$$

$$T_\infty = \max \left\{ \frac{4}{2}, \frac{4}{2}, \frac{5}{3}, \frac{5}{3}, \frac{5}{3}, \frac{8}{4}, \frac{8}{4}, \frac{5}{4}, \frac{5}{4} \right\} = 2$$

diagonal values (if not)
delay [L(m)]

Cycle Algorithm

Eg

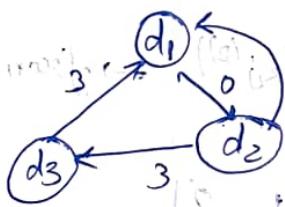


$$d_1 \rightarrow d_2 = 0$$

$$d_2 \rightarrow d_1 : d_2 \rightarrow 2 \rightarrow 3 \rightarrow 1 \rightarrow d_1 = 7$$

$$d_2 \rightarrow d_3 : d_2 \rightarrow 2 \rightarrow 3 \rightarrow d_3 = 3$$

$$d_3 \rightarrow d_1 : d_3 \rightarrow 4 \rightarrow 1 \rightarrow d_1 = 3$$



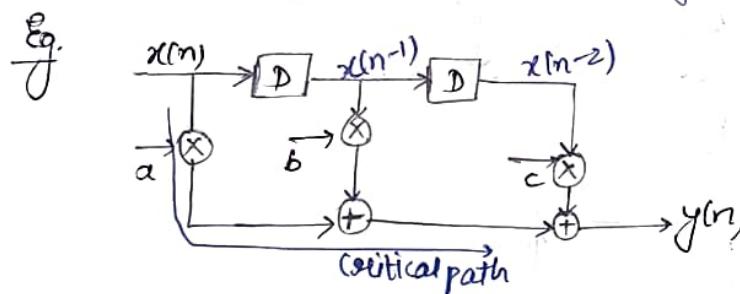
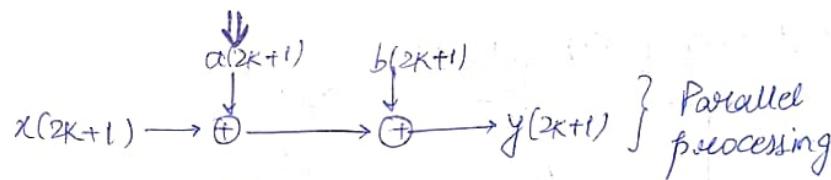
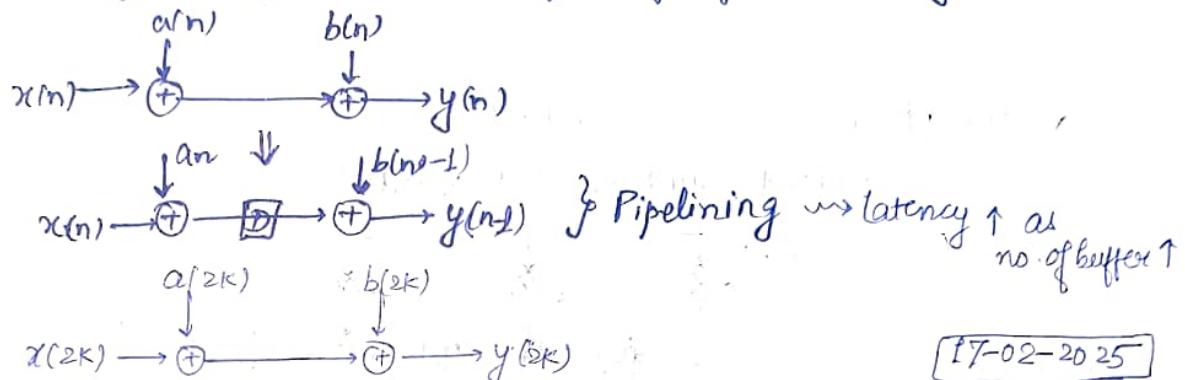
$$L^{(1)} = \begin{bmatrix} -1 & 0 & -1 \\ 7 & -1 & 3 \\ 3 & -1 & -1 \end{bmatrix}$$

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

$$T_{\max} = \max \left\{ \frac{7}{2}, \frac{7}{2}, \frac{6}{3}, \frac{6}{3}, \frac{6}{3} \right\} = 3.5$$

Pipelining

- ↪ Reduce the critical path.
- ↪ Either increase the clock speed or decrease the power consumption.
- ↪ Inserting pi buffer without affecting functionality.



T_M : multiplication-time
 T_A : addition-time

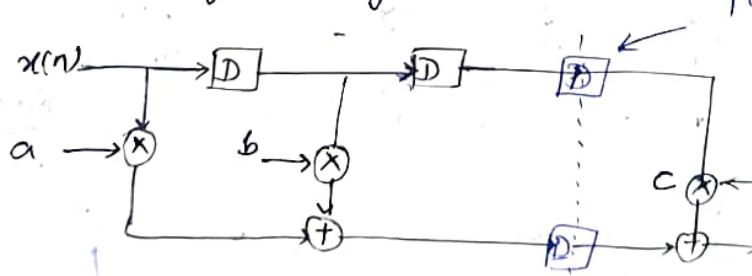
$$T_{\text{sample}} \geq T_M + 2T_A$$

$$f_{\text{sample}} \leq \frac{1}{T_M + 2T_A}$$

→ For faster input rate/sample rate, direct form structure cannot be used.

Pipelining of FIR Digital Filters

Place latches across ~~the~~ feed-forward cutset of the graph.



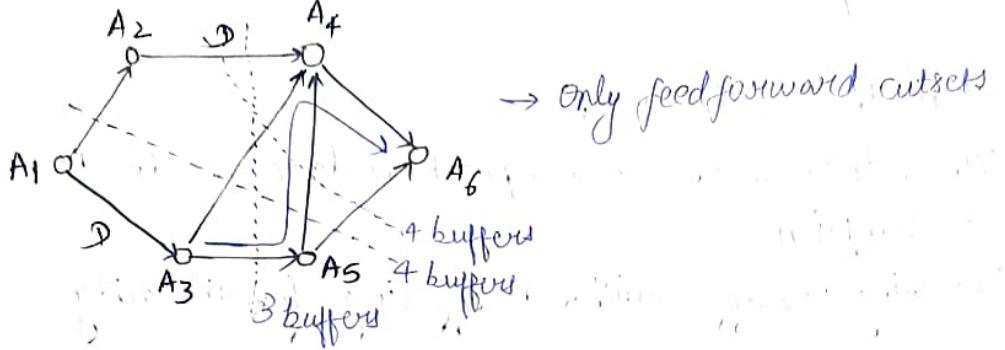
→ Performance \uparrow
 hardware complexity \uparrow

→ Latency \uparrow

→ Speed of DSP architecture (or clock period) is limited by the longest path b/w any 2 latches, or b/w an i/p and a latch or, b/w a latch and an o/p; or b/w i/p & o/p.

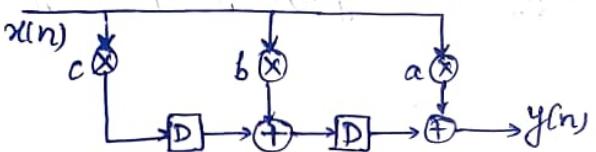
→ No. of latches \uparrow

Eg.

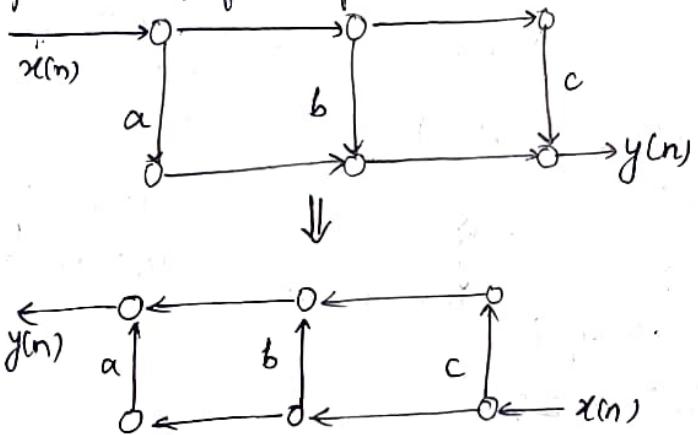


Data Broadcast Structure

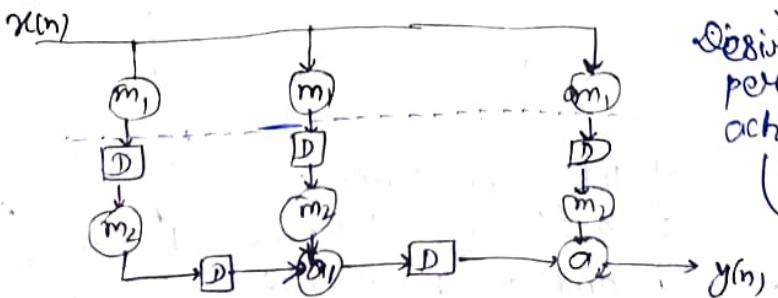
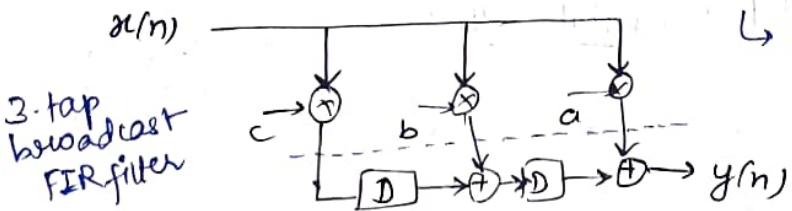
Data Broadcast structure of FIR filters:



Transposed SFG of FIR filter:



Fine Grain Pipelining



Let $T_M = 10 \text{ units}$, $T_A = 2 \text{ units}$

↳ Multiplier broken into 2 smaller units with processing times 6 units & 4 units.

↳ Desired clock period can be achieved in $\left(\frac{T_M + T_A}{2}\right)$.

Parallel Processing (Block processing)

- ↳ Increases the sampling rate by replicating the hardware.
 - ↳ Multiple outputs are computed in parallel in a clock period.
 - ↳ Can be used to reduce power consumption.

Eg SISO FIR filter

$$y(n) = ax(n) + bx(n-1) + cx(n-2)$$

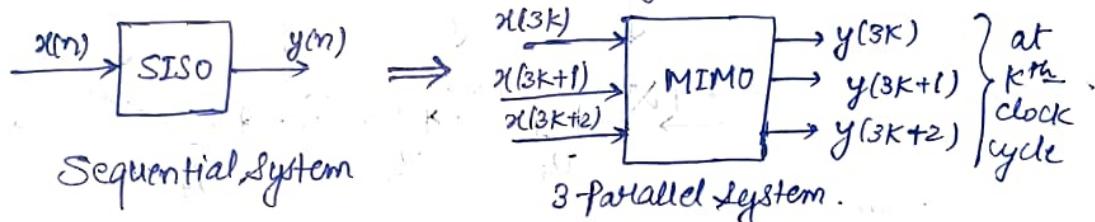
Convert to MIMO (with 3 i/p's per clock cycle to get parallel system):

$$y(3k) = ax(3k) + b\chi(3k-1) + c\chi(3k-2)$$

$$y(3k+1) = a x(3k+1) + b x(3k) + c x(3k-1)$$

$$g(3k+2) = a\chi(3k+2) + b\chi(3k+1) + c\chi(3k).$$

→ No. of i/p.s processed in a clock cycle : Block size.



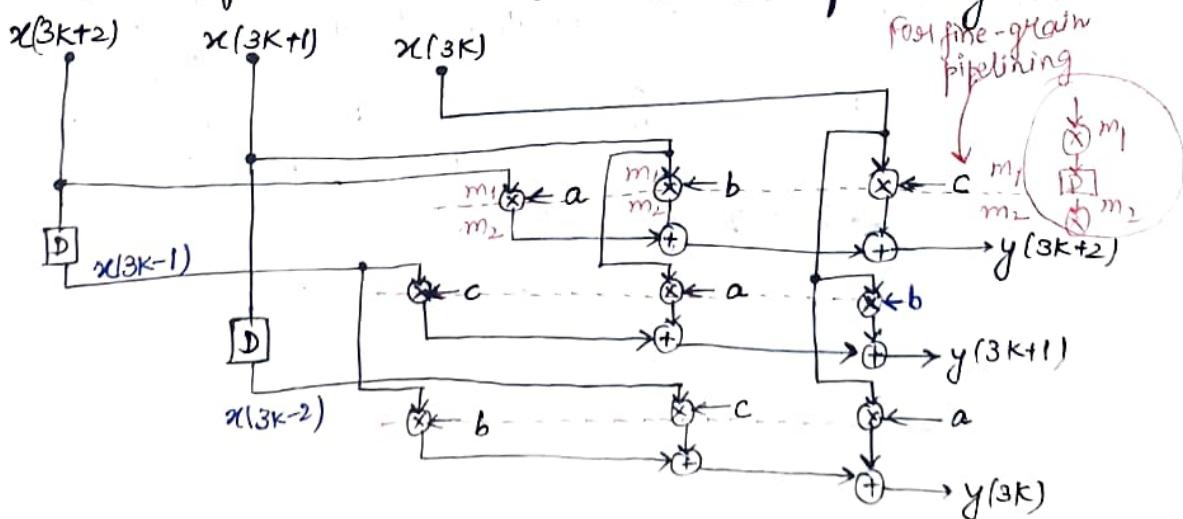
→ In MIMO structure, placing a latch at any line produces an effective delay of L (=block size) clock cycles at the sample state. So, each delay element known as block delay or L -slow.

↳ Power ↓

Performance maintained (critical path remains same).

Hardware Resources ↑

Eg. 3-tap FIR filter with block size 3 : Parallel processing arch.



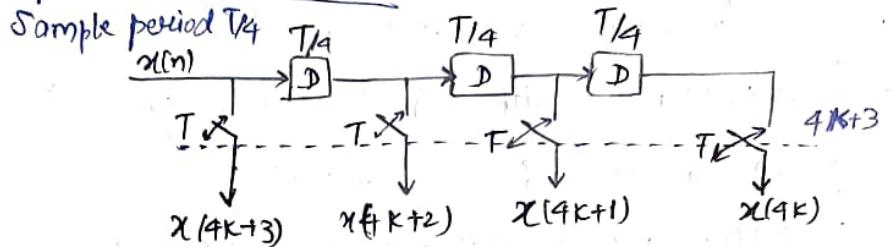
$$T_{clock} \geq T_M + 2T_A$$

$$T_{iteration} = T_{sample} = \frac{T_{clock}}{L} \geq \frac{T_M + 2T_A}{3}$$

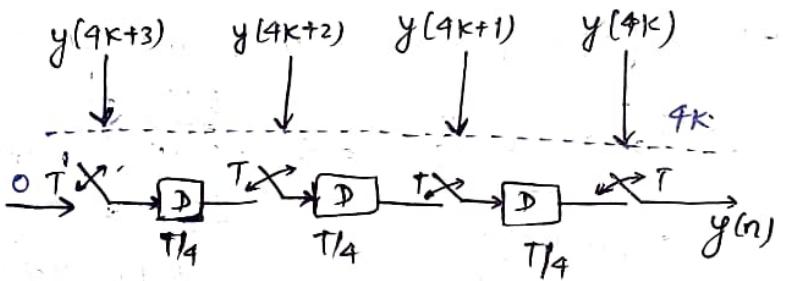
In parallel system: $T_{sample} \neq T_{clock}$

In pipelined system: $T_{sample} = T_{clock}$

Serial-to-parallel converter:



Parallel-to-serial converter:



why

When to use parallel processing:

↳ When the critical path is less than the I/O bound (system is communication bounded), then only pipelining can be used. Once this limit is reached, pipelining can no longer increase the speed.

↳ Combine pipelining with parallel processing to further increase the speed of DSP system.

Higher speed
↓
Lower power consumption

$$\Rightarrow T_{iteration} = T_{sample} = \frac{T_{clock}}{L \cdot M}$$

(for L block size parallel processing and)
M stage pipelining

→ Parallel processing can also be used for reduction of power consumption while using skew clocks.

For Low Power

Propagation delay (T_{pd}) of CMOS circuit:

$$T_{pd} = \frac{C_{charge} \cdot V_o}{K(V_o - V_t)^2}, \quad C_{charge}: \text{capacitance to be charged or discharged in a single clock cycle.}$$

Power consumption in CMOS circuit:

$$P_{CMOS} = C_{total} \cdot V_o^2 \cdot f, \quad C_{total}: \text{total capacitance of the CMOS ckt.}$$

In original sequential FIR filter:

$$P_{seq.} = C_{total} \cdot V_o^2 \cdot f, \quad f = \frac{1}{T_{seq}}$$

① For M-level pipelined system

↳ Critical path \downarrow to $1/M$ times.

↳ charge/disch. capacitance \downarrow to $1/M$ times.
in a single clock cycle

↳ Same clock speed (f) is maintained.

Supply voltage \downarrow to βV_o ($0 < \beta < 1$).

Pipelined filter:

$$P_{pip} = C_{total} \cdot \beta^2 \cdot V_o^2 \cdot f = \beta^2 \cdot P_{seq.} \downarrow \text{by } \beta^2$$

$$Pd: \quad T_{seq.} = \frac{C_{charge} \cdot V_o}{K(V_o - V_t)^2}, \quad T_{pip} = \frac{(C_{charge}/M) \cdot \beta V_o}{K(\beta V_o - V_t)^2}$$

For same clock speed,

$$\Rightarrow M(\beta V_o - V_t)^2 = \beta(V_o - V_t)^2 \rightarrow \text{solve to get } \beta.$$

② For L-parallel system

↳ Charging capac. does not change.

↳ Total capac. \uparrow by L times.

↳ To maintain same sampling rate

↳ Clock speed \uparrow to $L T_{seq}$.

↳ charging capac. is charged/disch. L times longer

Supply voltage \downarrow to βV_o .

$$T_{parallel} = \frac{(C_{charge}/L) \cdot \beta V_o}{K(\beta V_o - V_t)^2}$$

$$\Rightarrow L(\beta V_o - V_t)^2 = \beta(V_o - V_t)^2 \rightarrow \text{solve for } \beta.$$

$$\therefore P_{para} = (LC_{total}) (\beta V_o)^2 \frac{f}{L} = \beta^2 \cdot P_{seq}$$

③ Combining Pipelining and Parallel Processing for lower Power:

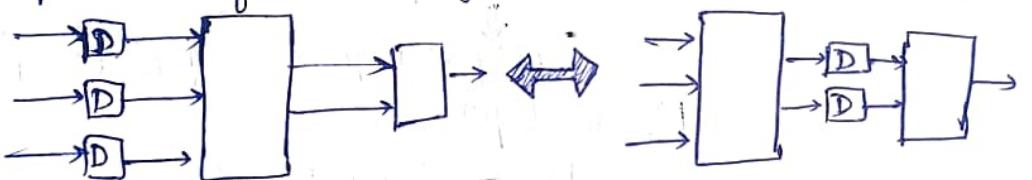
↪ Pipelining reduces the capacitance to be charged/discharged in 1 clock period, while parallel processing increases the clock period for charging/discharging the original capacitance.

$$p.d. : LT_{pd} = \frac{(C_{charge}/M) \beta V_o}{K(V_o - V_t)^2} = \frac{L \cdot C_{charge} \cdot V_o}{K(V_o - V_t)^2}$$

$$\Rightarrow LM \cdot f \beta (V_o - V_t)^2 = \beta (V_o - V_t)^2 \rightarrow \text{solve for } \beta.$$

Retiming

↪ A graph-transformation by which D-elements are relocated.



↪ Reduce critical path

↪ Improvement in speed (clock period \downarrow)

↪ Power consumption \downarrow

↪ No. of register \downarrow ↪ Used in logic synthesis.

↪ Latency \rightarrow does not alter

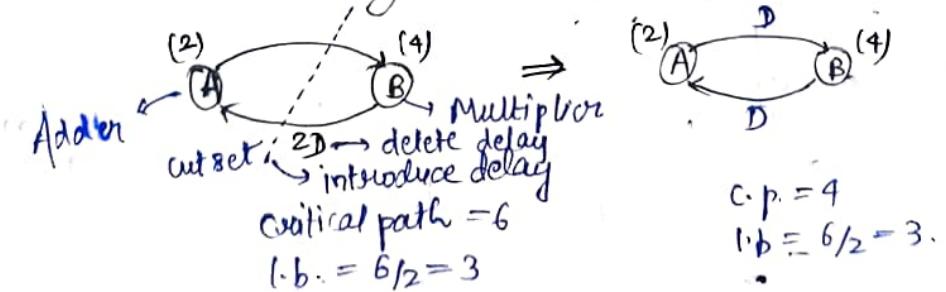
↪ Does not change

↪ delay in loop D

↪ iteration bound

Node Retiming: Move a D-element from each i/p edge to each edge (or vice versa).

Cutset Retiming: One type of pipelining.



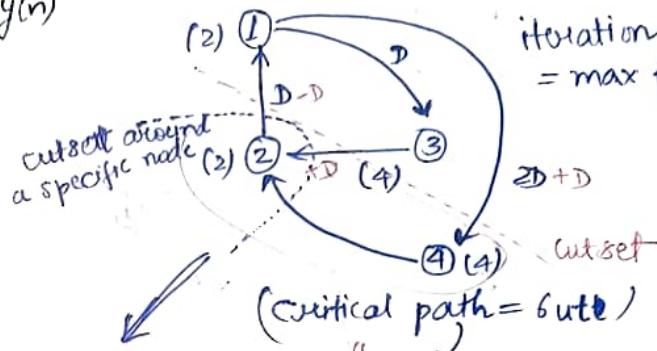
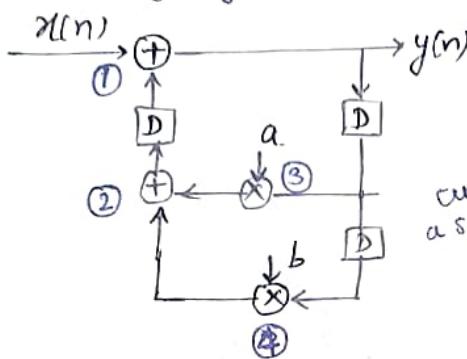
$$\begin{aligned} c.p. &= 4 \\ l.b. &= 6/2 = 3. \end{aligned}$$

$$y(n) = x(n) + y(n-2)$$

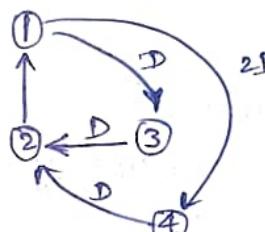
$$y(n) = x(n-1) + y(n-3)$$

Cutset Retiming:

→ Add delays to edges going one way & remove from ones going the other.

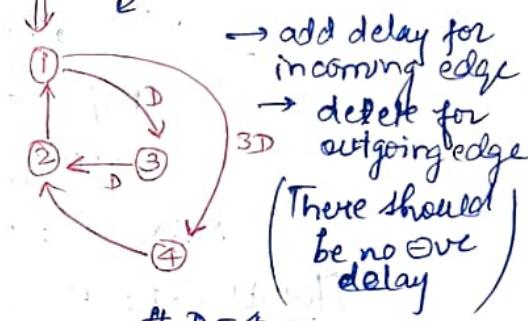


$$\text{iteration bound} = \max \left\{ \frac{8}{2}, \frac{8}{3} \right\} = \frac{8}{2}$$



Total critical = 4 u.t.

$$\# D = 5$$



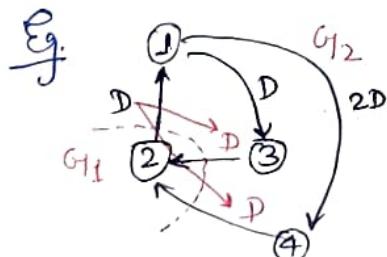
→ add delay for incoming edge
→ delete for outgoing edge
(There should be no over delay)

Retiming value, $r(v)$: No. of delays moved from each output edge of the node v to each of its i/p edges.

Feasibility constraints:

$$w_r(e) = w(e) + r(v) - r(u) \geq 0 : \text{no. of delay elements} \geq 0 \text{ for each edge } uv.$$

✓ ↓ ↓
after before (w: weight)
retiming retiming delay,



original weights

$$\begin{aligned} 1 \rightarrow 3 &= 1 \\ 1 \rightarrow 4 &= 2 \\ 2 \rightarrow 1 &= 1 \\ 3 \rightarrow 2 &= 0 \\ 4 \rightarrow 2 &= 0 \end{aligned}$$

$$w_r(e) = w(e) + r(N_{\text{receive}}) - r(N_{\text{send}})$$

$$\begin{aligned} 1 \rightarrow 3 &= 1 + 0 - 0 = 1 \\ 1 \rightarrow 4 &= 2 + 0 - 0 = 2 \\ 2 \rightarrow 1 &= 1 + 0 - 1 = 0 \\ 3 \rightarrow 2 &= 0 + 1 - 0 = 1 \\ 4 \rightarrow 2 &= 0 + 1 - 0 = 1 \end{aligned}$$

Choose retiming values:

$$G_2: r(1)=r(3)=r(4)=0$$

$$G_1: r(2)=1$$

Retimed weights

Pipelining: A special case of cutset retiming by placing delays at feedforward cutset.

Constraints : [Bellman Ford Algorithm]

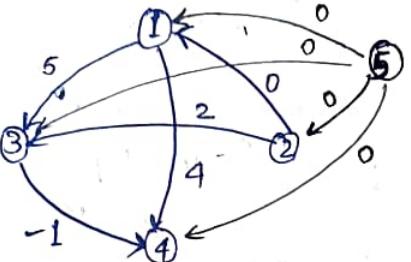
$$\gamma_1 - \gamma_2 \leq 0 \quad [\text{Edge } j \rightarrow i \text{ with weight } k]$$

$$\gamma_3 - \gamma_1 \leq 5$$

$$\gamma_4 - \gamma_1 \leq 4$$

$$\gamma_4 - \gamma_3 \leq -1$$

$$\gamma_3 - \gamma_2 \leq 2$$



w, r
weight \Rightarrow returned-vector

Reference node = 5th

$$w(5-1) = w(5-2) = w(5-3) = w(5-4) = 0$$

Ref. node : $r^0(N) = \{0, 0, 0, 0\}$

Iteration - 1 : 1st node: $w = r^0(1) + w(1 \rightarrow 1)$,
($N=1$)

$$r^0(1) + w(1 \rightarrow 1),$$

$$r^0(2) + w(2 \rightarrow 1),$$

$$r^0(3) + w(3 \rightarrow 1),$$

$$r^0(4) + w(4 \rightarrow 1)$$

$$= \{r^0(N) + w(N-1)\}$$

$$= \{0, \infty, 0, \infty, \infty\} \xrightarrow{\min} 0$$

$$w = r^0(1) + w(1 \rightarrow 2)$$

$$r^0(2) + w(2 \rightarrow 2)$$

$$r^0(3) + w(3 \rightarrow 2)$$

$$r^0(4) + w(4 \rightarrow 2)$$

$$= \{0, \infty, \infty, \infty, \infty\} \xrightarrow{\min} 0$$

$$w = \{5, 2, \infty, \infty\} \xrightarrow{\min} 0$$

$$w = \{0, 4, \infty, -1, \infty\} \xrightarrow{\min} -1$$

For no direct path existing
Take ∞ weight
For path going to itself

Iteration-2: $\gamma^1(v) = \{0, 0, 0, -1\}$

$$v=1: w = \gamma^0(1) + \overset{\infty}{w}(1-1)$$

$$\gamma^0(2) + \overset{\infty}{w}(2-1)$$

$$\gamma^0(3) + \overset{\infty}{w}(3-1)$$

$$\gamma^0(4) + \overset{\infty}{w}(4-1)$$

$$= \{\infty, 0, \infty, \infty\}$$

$$v=2: w = \{\infty, \infty, \infty, \infty\}$$

$$v=3: w = \{5, 2, \infty, \infty\}$$

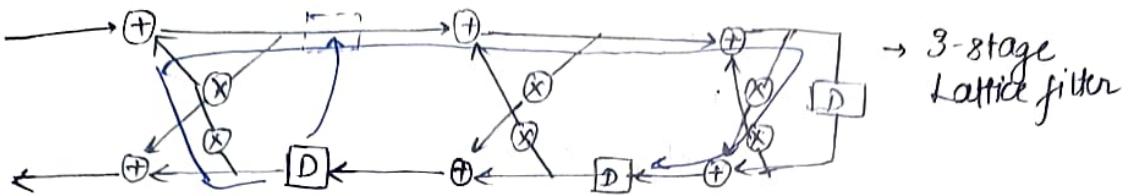
$$v=4: w = \{4, \infty, -1, \infty\}$$

$$\therefore \gamma^2(v) = \{0, 0, 0, -1\} \rightarrow \text{same as the prev.}$$

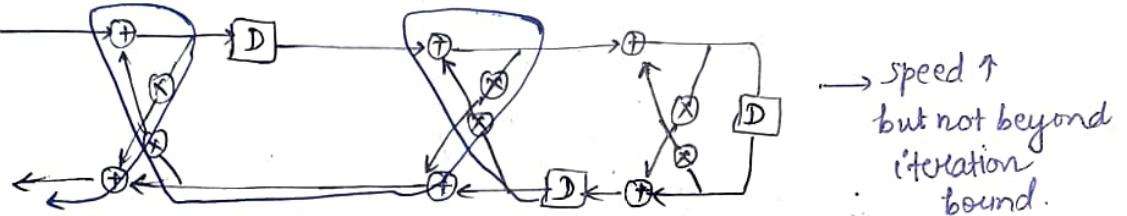
↓
stop.

$$\left. \begin{array}{l} \gamma(1)=0 \\ \gamma(2)=0 \\ \gamma(3)=0 \\ \gamma(4)=-1 \end{array} \right\} \text{Retimed values}$$

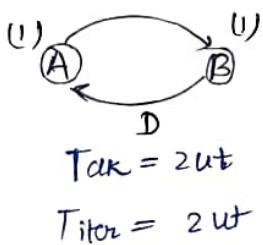
Pipelining



$$T_{\text{critical}} = 4T_A + 2T_M$$



K-Slow-down and Retiming

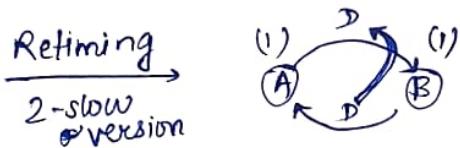


2 slow transformation → 

Replace each D by kD

→ K-1 null operation must be inserted

→ Hardware utilization is reduced (e.g., 50%).

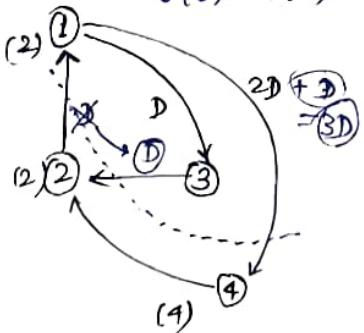


PCDT Constraints from DFG:

$$w_r(e) = w(r) + \gamma(V) - \gamma(U) \geq 0$$

$$\gamma(U) - \gamma(V) \leq w(e)$$

$$T_{\text{iter}} = T_{\infty} = \max \left\{ \frac{8}{2}, \frac{8}{3} \right\} = \frac{8}{2}$$



Nodes: cannot be reduced computation time → Total \rightarrow can be reduced ✓

$$\begin{aligned} \textcircled{1} &\rightarrow \textcircled{3} & 2+4 = 6 \text{ ut} &\rightarrow w_r(1-3) = w(1-3) + r(4) - r(1) \\ \textcircled{2} &\rightarrow \textcircled{1} & 4 \text{ ut } X &\rightarrow w_r(2-1) = 1 + r(1) - r(2) \\ \textcircled{1} &\rightarrow \textcircled{4} & 6 \text{ ut } \checkmark &\rightarrow w_r(1-4) = 2 + r(4) - r(1) \\ \textcircled{4} &\rightarrow \textcircled{2} & 6 \text{ ut } \checkmark &\rightarrow w_r(4-2) = 0 + r(2) - r(4) \\ \textcircled{3} &\rightarrow \textcircled{2} & 6 \text{ ut } \checkmark &\rightarrow w_r(3-2) = 0 + r(2) - r(3) \end{aligned}$$

$$w_r(1-3) \geq 1$$

$$w_r(2-1) \geq 0$$

$$w_r(1-4) \geq 1$$

$$w_r(4-2) \geq 1$$

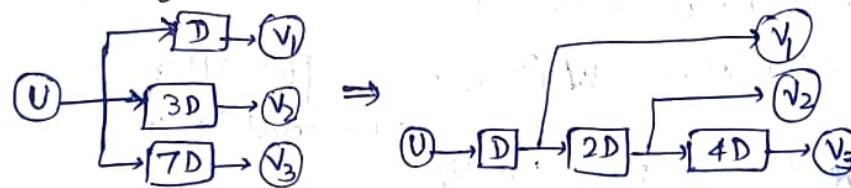
$$w_r(3-2) \geq 1$$

27-02-2025

Register Minimization

↳ Find a feasible soln with min. no. of registers in the design space constrained by feasibility and critical path constraint (the refined timeline space with min. clock period).

Eg. Multiple fan-out problem



↗ fan out ↑
 ↳ driving multiple o/p's
 ↳ more driving current ↑

Minimize Clock Period

Min. feasible clock period of a graph G_1 is

$$\phi(G_1) = \max \{ t(p) : w(p) = 0 \}.$$

$w(u,v)$: Min^m no. of registers on any path $u \rightarrow v$

$D(u,v)$: Max^m computation time among all paths from $u \rightarrow v$ with weight $w(u,v)$.

$\rightarrow M = t_{\max} \cdot n$, t_{\max} : $\max_{e \in G}$ computation time of any node in G .
 n : no. of nodes in G .

form a new graph G' which is same as G except the edge weights replaced by

$$w'(e) = M w(e) - t(v) \quad (e = u \rightarrow v)$$

solve for all-pairs shortest path on G' (S_{UV})

If $u \neq v$, then

$$W(u, v) = \lceil S_{UV}/M \rceil \xrightarrow{\text{next integer value}}$$

$$D(u, v) = M \times W(u, v) - S_{UV} + t(v)$$

If $u = v$, $w(u, v) = 0$,

$$D(u, v) = t(u, v)$$

Then, we use $w(u, v)$, $D(u, v)$ to find if there is a soln such that $\Phi(G) \leq c$.

construct the set of constraints:

- feasibility constraint:

$$\gamma(u) - \gamma(v) \leq w(e) \quad \text{for every edge } e \text{ in } G.$$

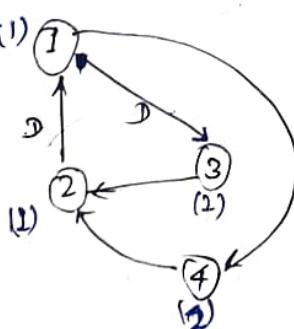
- critical path constraint:

$$\gamma(u) - \gamma(v) \leq W(u, v) - 1 \quad \text{for all nodes } u, v \text{ in } G \\ \text{such that } D(u, v) > c \quad (\text{cycle time}).$$

solve to get $\gamma(\cdot)$ (retiming values)

$$w_s(p) = w(p) + \gamma(b) - \gamma(a)$$

Eg.

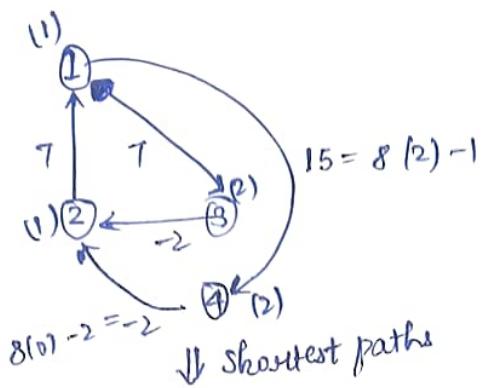


$$M = 2 \times 4 = 8$$

$$w'(e) = M \times w(e) - t(v)$$

$$y(n) = ay(n-2) + by(n-3) + x(n)$$

\Downarrow find G'

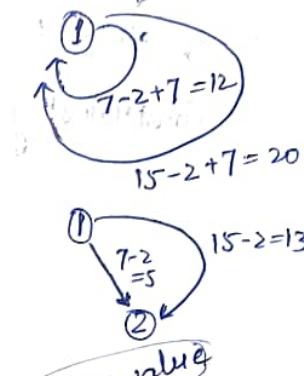


| S_{UV} | 1 | 2 | 3 | 4 |
|----------|----|----|----|----|
| 1 | 12 | 5 | 7 | 15 |
| 2 | 7 | 12 | 14 | 22 |
| 3 | 5 | -2 | 12 | 20 |
| 4 | 5 | -2 | 12 | 20 |

↓ Find $w(u,v)$, $D(v,v)$

if $u \neq v$, $w(u,v) = [S_{UV}/M]$ next integer value
 $v = u$, $w(u,v) = 0$, $D(v,v) = t(v) = t(u)$

| $w(u,v)$ | 1 | 2 | 3 | 4 |
|----------|---|---|---|---|
| 1 | 0 | 1 | 1 | 2 |
| 2 | 1 | 0 | 2 | 3 |
| 3 | 1 | 0 | 0 | 3 |
| 4 | 1 | 0 | 2 | 0 |



| $D(v,v)$ | 1 | 2 | 3 | 4 |
|----------|---|---|---|---|
| 1 | 1 | 4 | 3 | 3 |
| 2 | 2 | 1 | 4 | 4 |
| 3 | 4 | 3 | 2 | 6 |
| 4 | 4 | 3 | 6 | 2 |

Feasibility constraints:

$$\cancel{w(e)}: r(v) - r(u) \leq w(e)$$

$$\therefore r(2) - r(1) \leq 1$$

$$r(1) - r(3) \leq 1$$

$$r(1) - r(4) \leq 2$$

$$r(3) - r(2) \leq 0$$

$$r(4) - r(2) \leq 0$$

Critical path constraints:

$$r(u) - r(v) \leq w(u,v) - 1$$

for all vertices u, v in G such that $D(u,v) > 3$.

$$\therefore r(1) - r(2) \leq 0$$

$$r(2) - r(3) \leq 1$$

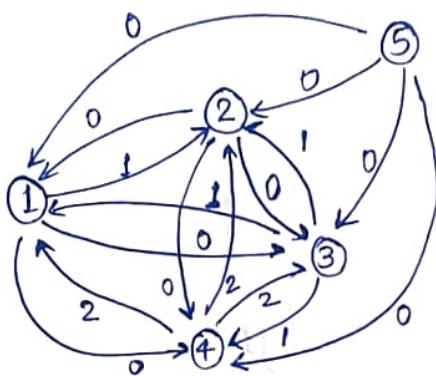
$$r(2) - r(4) \leq 2$$

$$r(3) - r(1) \leq 0$$

$$r(3) - r(4) \leq 2$$

$$r(4) - r(1) \leq 0$$

$$r(4) - r(3) \leq 1.$$

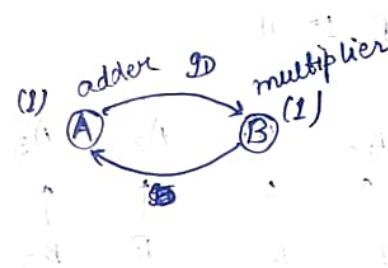
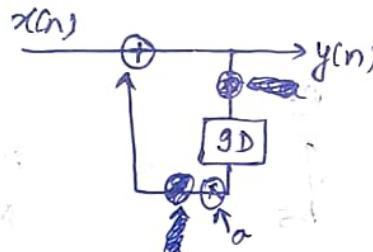


Constraint graph for the inequalities

Unfolding

- Transformation technique applied to a DSP program to create a new program describing more than one iteration of the original program.
- Unfolding factor: $J \rightarrow J$ consecutive iterations of the original program.
Also used to design bit parallel and word parallel architectures from bit serial and word serial arch.

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



- Unfolding factor, $J=2$.

$$y(2k) = ay(2k-1) + x(2k)$$

$$y(2k+1) = ay(2k) + x(2k+1)$$

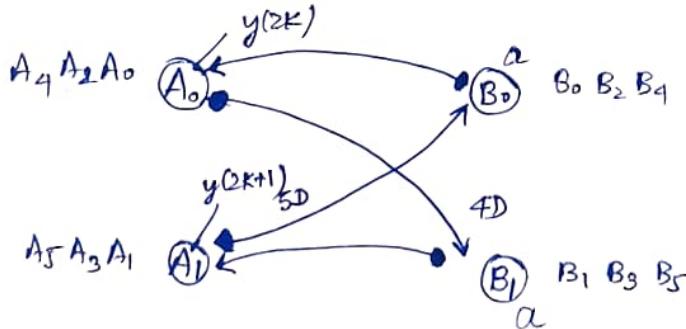
$$\Rightarrow y(2k) = ay(2k-1) + x(2k)$$

$$= ay(2(k-1)+1) + x(2k)$$

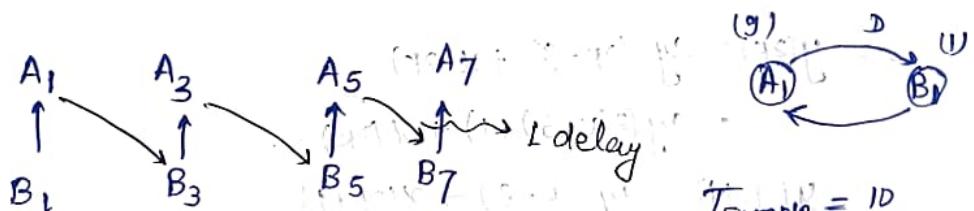
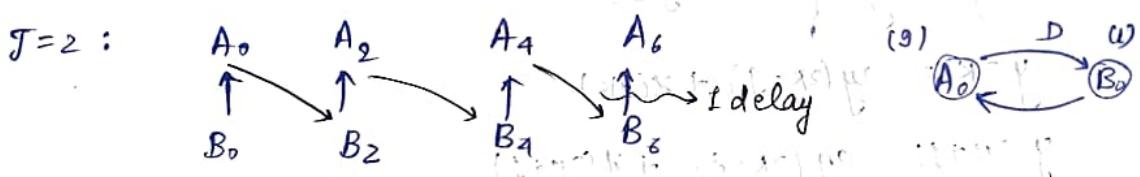
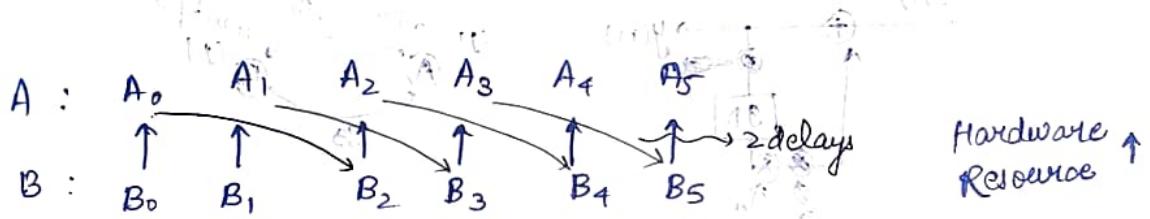
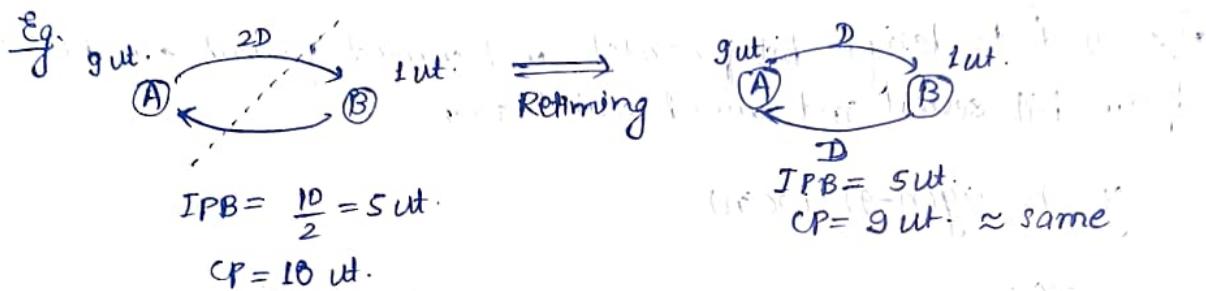
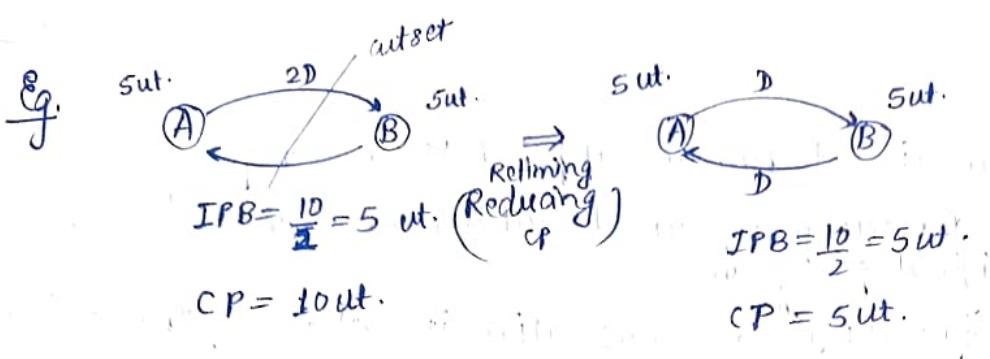
$$y(2k+1) = ay(2k) + x(2k+1)$$

$$= ay(2(k-1)+0) + x(2k+1)$$

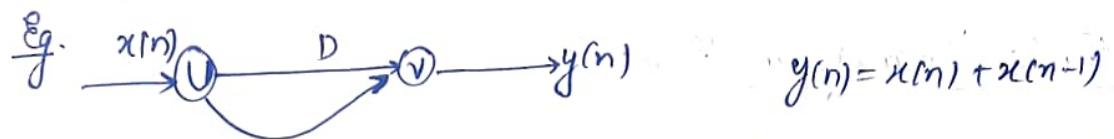
2-architectures: even & odd



A_0, B_0 : even
 A_1, B_1 : odd



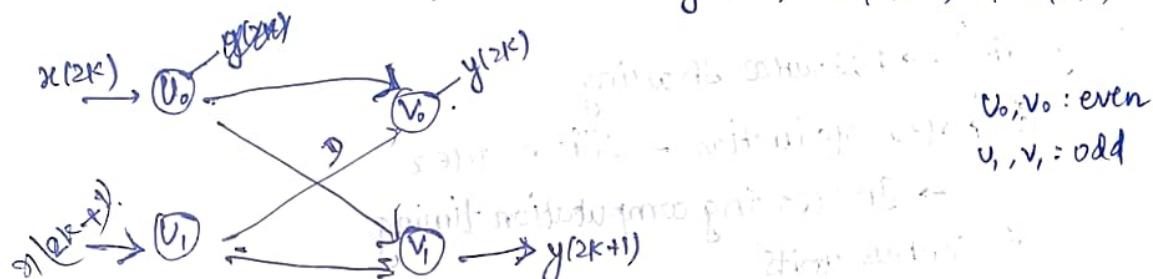
$$T_{\text{sample}} = \frac{10}{2} \text{ No. of unfoldings}$$



$J=2$ (unfolding factor)

$$y(2k) = x(2k) + x(2k-1),$$

$$y(2k+1) = x(2k+1) + x(2k)$$



v_0, v_o : even
 v_1, v_i : odd

$$v_1, v_2 = \text{odd}$$

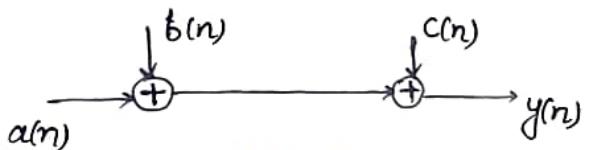
For Best wife

Algorithmic Transformations:

- ① Unfolding
- ② Retiming
- ③ Folding → Resource sharing
 - i) Area reduction - silicon area
→ Increasing computation timing
 - ii) Control units
 - iii) Adding of registers.

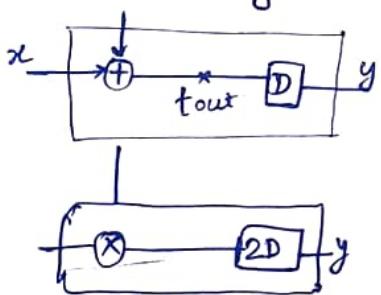
Folding

E.g.



$N \geq 2$: no. of clock cycles required for single iteration

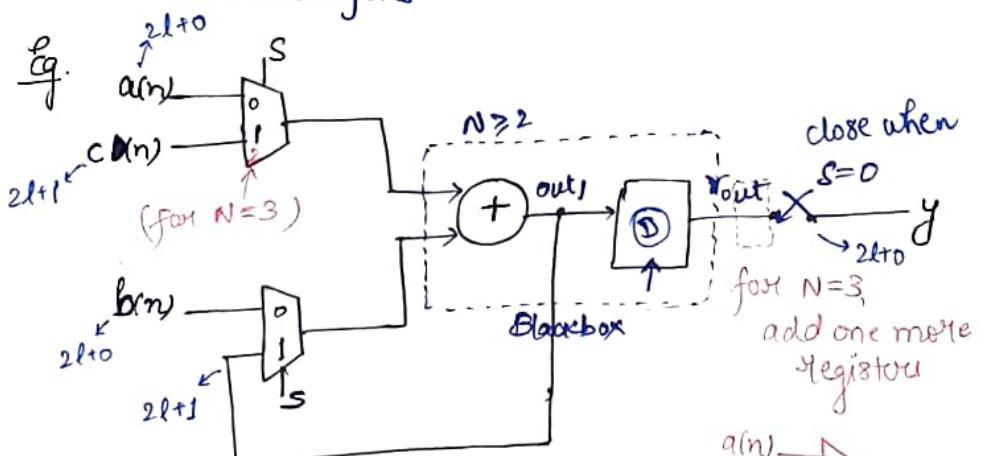
Pipeline stage:



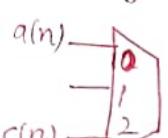
| t: | 0 | 1 | 2 |
|-------|-------------|-------------|-------------|
| in 1: | a_0 | a_1 | a_2 |
| in 2: | b_0 | b_1 | b_2 |
| tout: | $a_0 + b_0$ | $a_1 + b_1$ | $a_2 + b_2$ |
| - | $a_0 + b_0$ | $a_1 + b_1$ | |

[Adder takes 1 ut.]

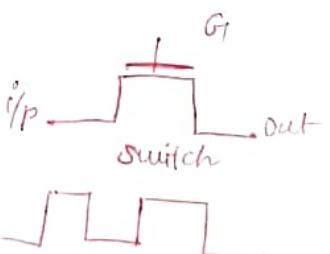
Latency:
2 clock cycles



for $N=3$,
add one more
register



control
signal

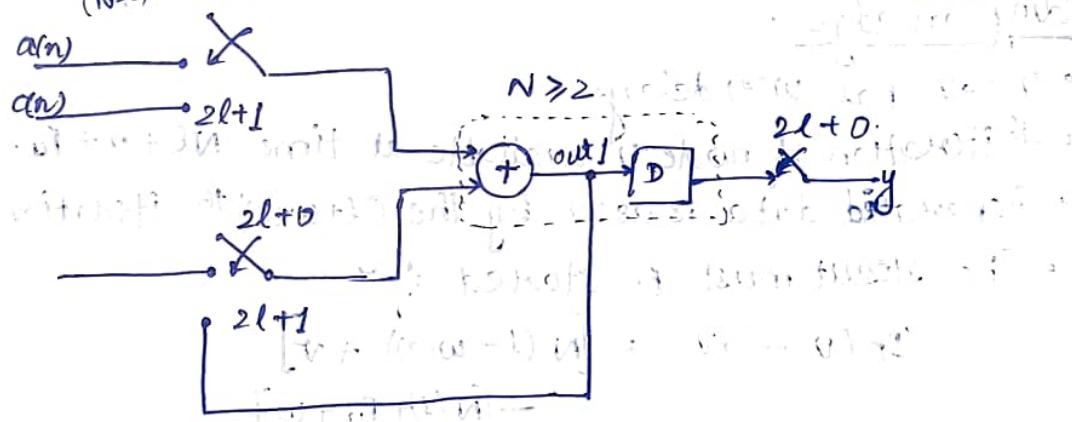


| | | | | | | |
|---------------------|-------------|-------------------|-------------|---|---|---|
| S: | 0 | 1 | 0 | 1 | 0 | 1 |
| timestamp: (CLK) | 0 | 1 | 2 | 3 | 4 | 5 |
| in1 | a_0 | $a_0 + c_0$ | a_1 | | | |
| in2 | b_0 | $b_0 + a_0$ | b_1 | | | |
| out1 | $a_0 + b_0$ | $a_0 + b_0 + c_0$ | $a_1 + b_1$ | | | |
| rout1 | $a_0 + b_0$ | $a_0 + b_0 + c_0$ | $a_1 + b_1$ | | | |
| y | - | $a_0 + b_0$ | - | | | |

Take O/P when $[\text{Tap the output}]$

Time to wait $\therefore N \geq 2$: (Minm clock cycles)

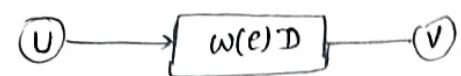
$N=2$: folding factor



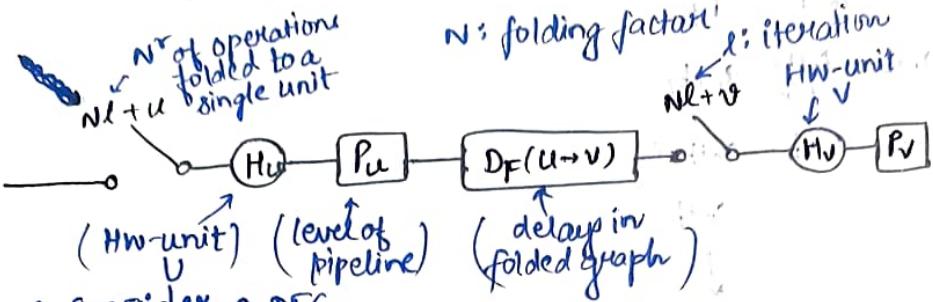
With N^2 clock cycles

Folding Transformation

$$w(e) \geq 0$$



$u & v$ are folding orders, i.e., scheduled time
 $0 \leq u, v \leq N-1$

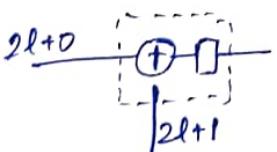
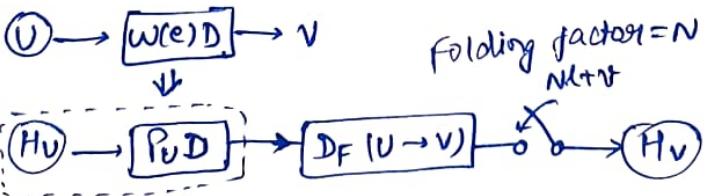


- Consider a DFG.
- An edge e connecting nodes U and V with $w(e)$ delays.
- Executions of the l th iterations of U & V at time units $NL + u$ & $NL + v$.
- u and v (folding orders) : time partitions at which the nodes are scheduled to execute & satisfy: $0 \leq u, v, N$.
- N : folding factor \rightarrow no. of operations folded to a single HW-unit.
- $HU, HV \rightarrow$ l units to execute nodes U and V (HV pipelines by PU stages).

Folding an edge:

- $U \xrightarrow{e} V$ has $w(e)$ delays.
- l th iteration of node V available at time $NL + u + P_u$.
- Generated data is used by the $(l + w(e))$ th iteration of V .
- The result must be stored for

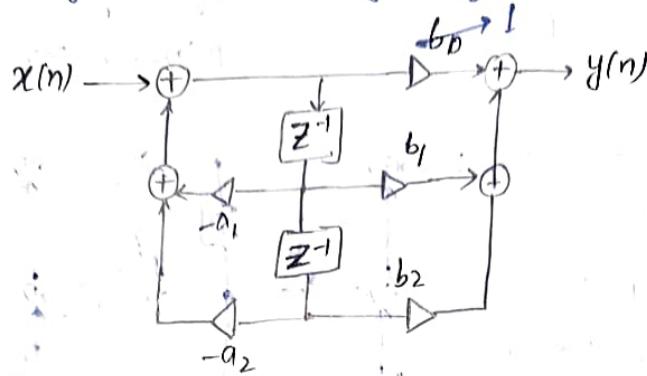
$$\begin{aligned} DF(U \xrightarrow{e} V) &= [N(l + w(e)) + v] \\ &\quad - [N(l + P_u + u)] \\ &= Nw(e) - P_u + v - u \end{aligned}$$



$$\begin{aligned} DF &= [2(0+0)+1] - [0+1+0] \\ &= 1 - 1 = 0 \end{aligned}$$

Biquad Filter

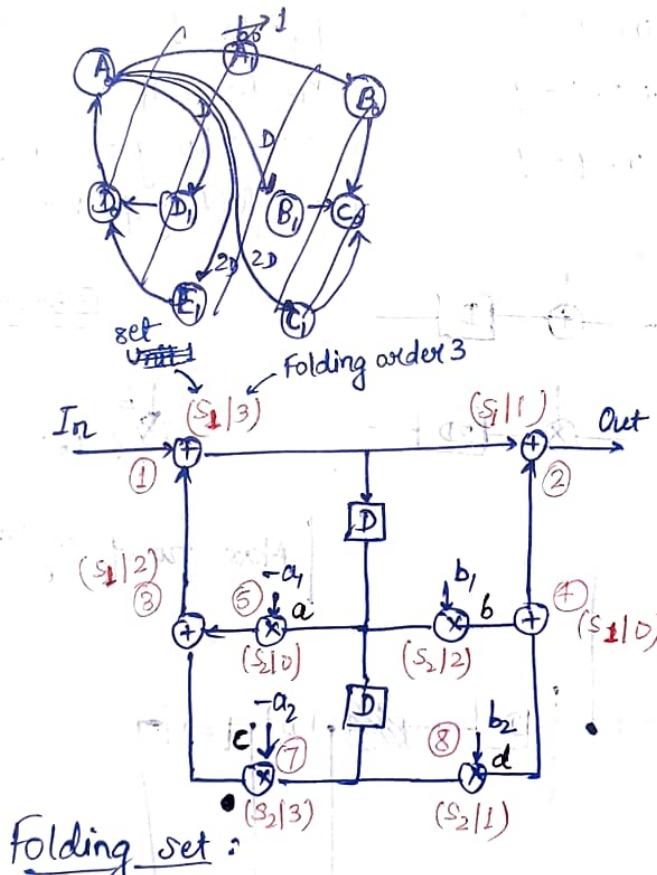
↳ 2 poles, 2 zeros digital filter



Police: improved response

Zeros : reduce \Rightarrow suppress performance ↓

$$H(z) = g \cdot \frac{1 + \beta_1 z^{-1} + \beta_2 z^{-2}}{1 + \alpha_1 z^{-1} + \alpha_2 z^{-2}}$$



An order of operations executed by same hardware.

$$\text{E.g. } S_1 = \{A_1, \phi, A_2\}$$

$$A_1: (S_1)D), A_2(S_1)2$$

Biquad filter :

Folding equation:

$$DF(U \rightarrow V) = Nw(e) - Pu + v - u^*$$

receive send

$$DF(1 \rightarrow 2) = 4(1) - 1 + 1 - 3 = 1$$

$$DF(1 \rightarrow 5) = 4(1) - 1 + 0 - 3 = 0$$

$$DF(1 \rightarrow 6) = 4(1) - 1 + 2 - 3 = 2$$

$$DF(1 \rightarrow 7) = 4(1) - 1 + 3 - 3 = 3$$

$$DF(1 \rightarrow 8) = 4(2) - 1 + 1 - 3 = 5$$

$$DF(3 \rightarrow 1) = 4(0) - 1 + 3 - 2 = 0$$

$$DF(4 \rightarrow 2) = 4(0) - 1 + 1 - 0 = 0$$

$$DF(5 \rightarrow 3) = 4(0) - 1 + 1 - 0 = 0$$

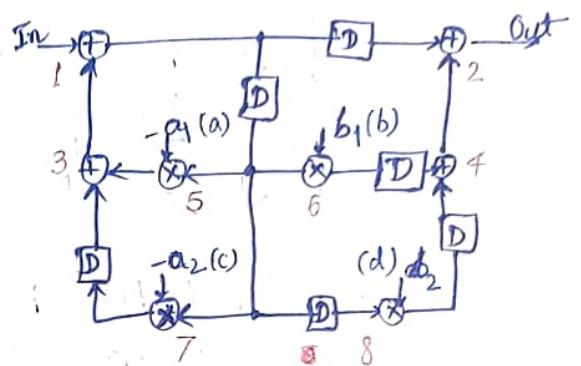
$$DF(6 \rightarrow 4) = 4(1) - 2 + 0 - 2 = 0$$

$$DF(7 \rightarrow 3) = 4(1) - 2 + 2 - 3 = 1$$

$$DF(8 \rightarrow 4) = 4(1) - 2 + 0 - 1 = 1$$

Valid folding:

$$D(U \rightarrow V) > 0$$



$P_u = P_A = 1$, if U is adder

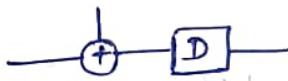
$P_u = P_M = 2$, if U is multiplier

$w(e)$: delay b/w U & V

N : folding factor = 4.

25-03-2025

$$S_1 = \{ \overset{0}{4}, \overset{1}{2}, \overset{2}{3}, \overset{3}{1} \}$$

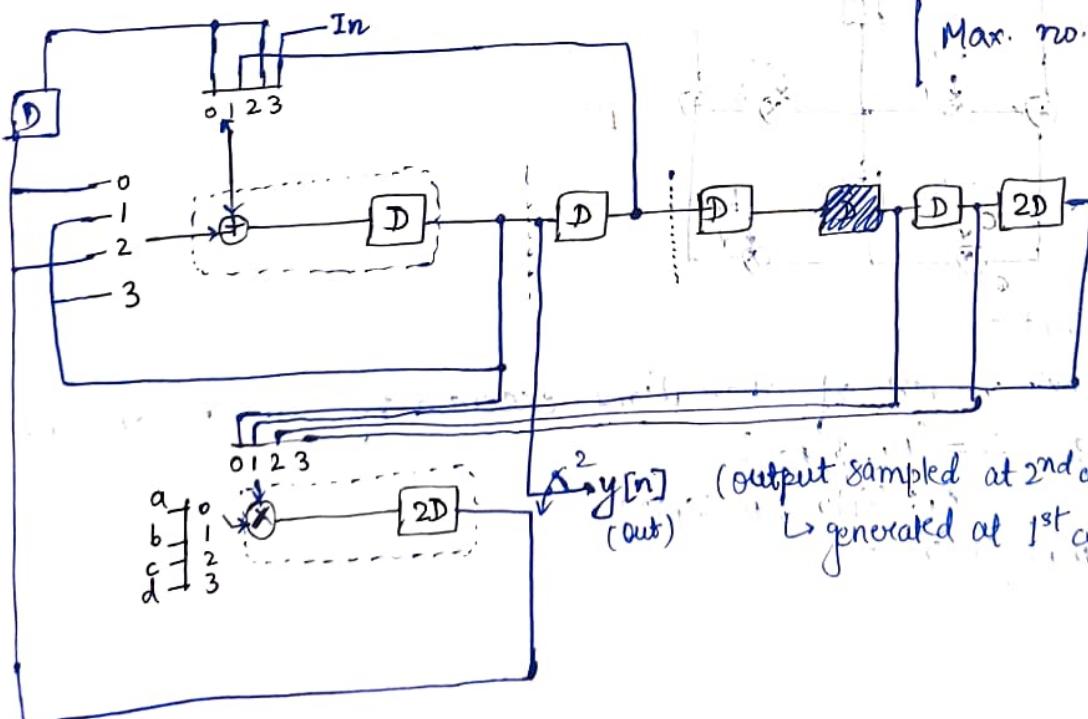


$$S_2 = \{ \overset{0}{5}, \overset{1}{8}, \overset{2}{6}, \overset{3}{7} \}$$

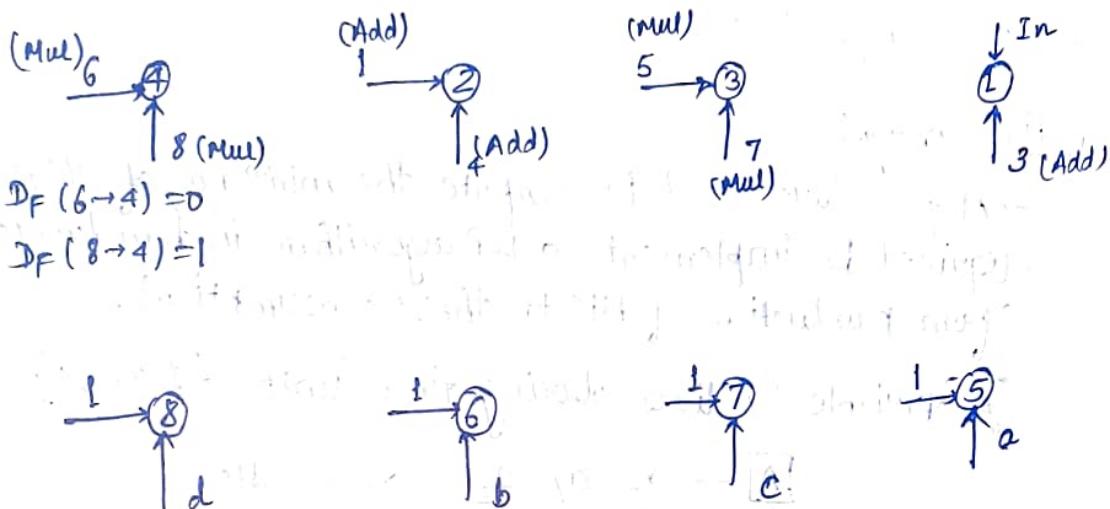


$$\overset{b}{X} (S_2 | 0)$$

Max. no. of Registers = 5



$y[n]$ (output sampled at 2nd cycle)
(out) ↴ generated at 1st cycle ($S_1 | \downarrow$)



27-03-2025

Write verilog code to implement the biquad filter.

↳ Control unit

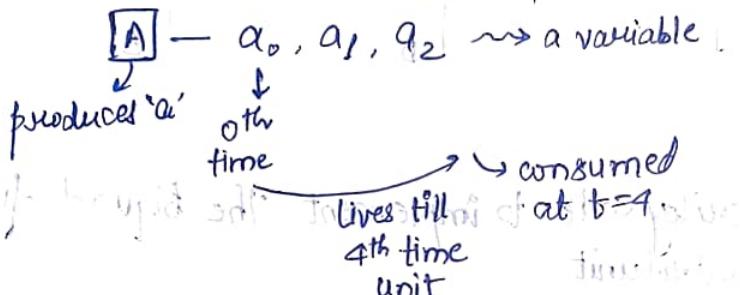
↳ ALU / Data path unit

Register Minimization

Lifetime Analysis

→ A procedure used to compute the min^m no. of registers required to implement a DSP algorithm in hardware (from production of bits to their consumption).

Eg. variable 'a' lives during time unit {1, 2, 3, 4}



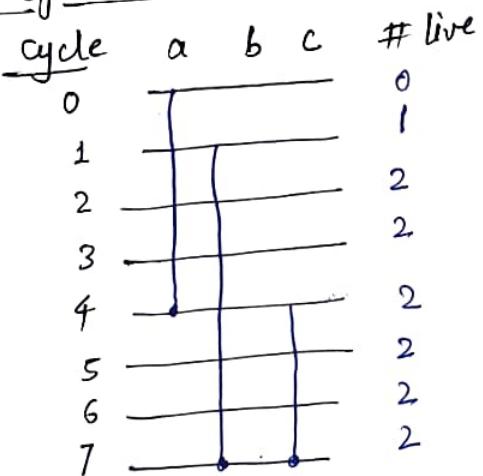
b lives during {2, 3, 4, 5, 6, 7}
c lives during {5, 6, 7}

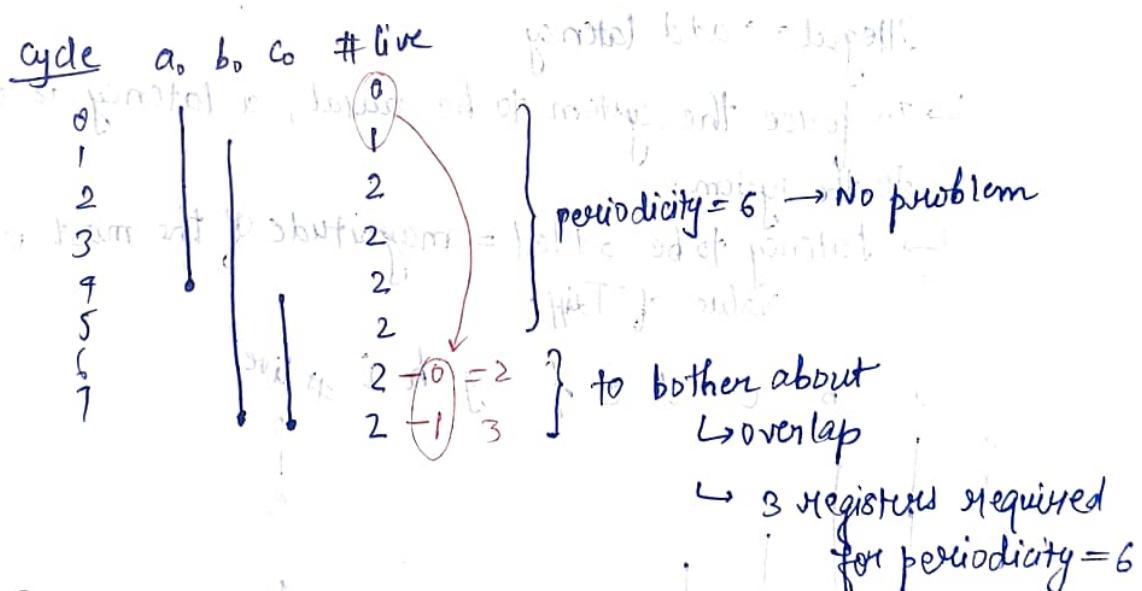
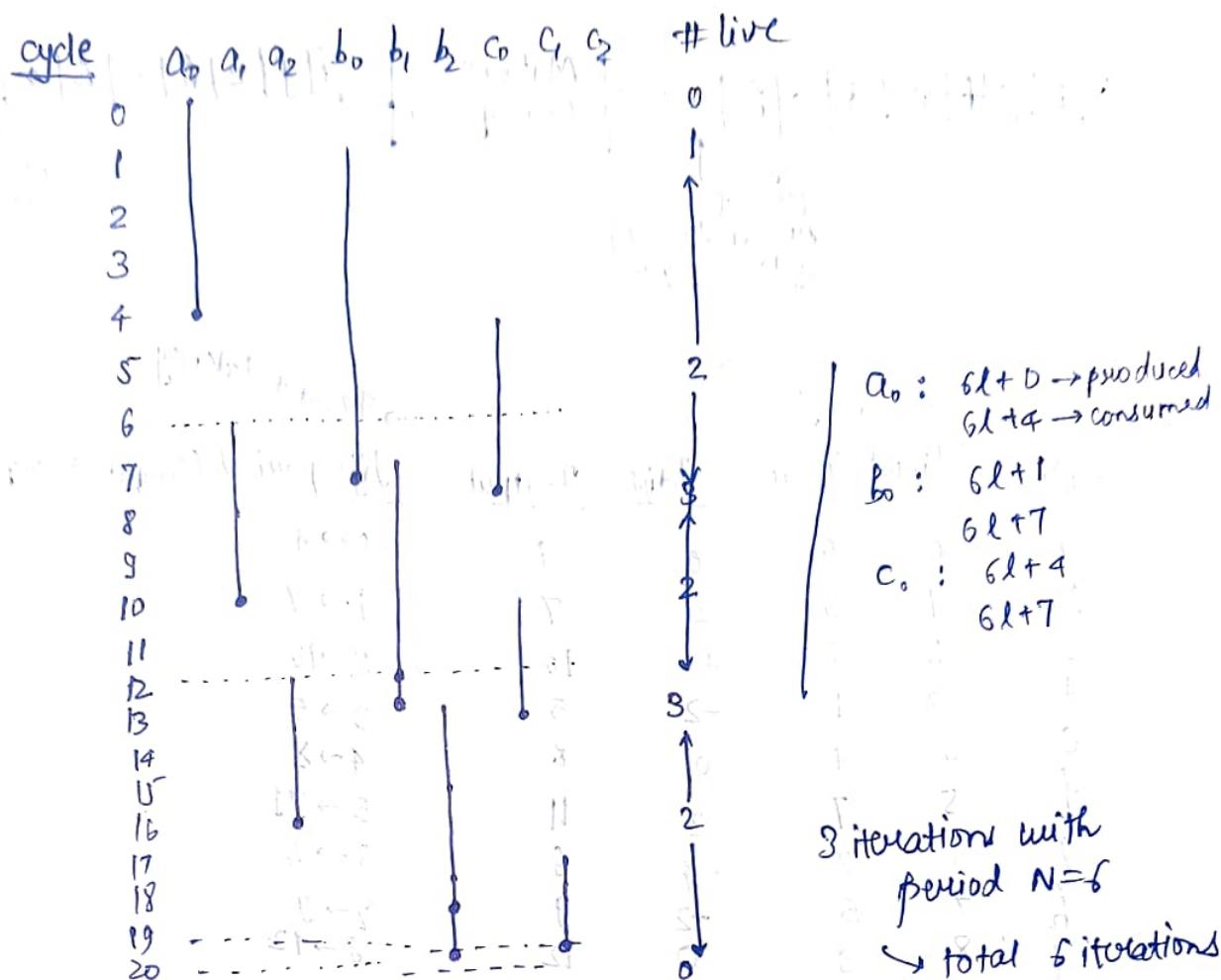
time stamp : {1, 2, 3, 4, 5, 6, 7}

a a.b a.b a.b b,c b,c b,c → variables living
1 2 2 2 2 2 2 → registers required.

Min^m Regs required = max{1, 2, 2, 2, 2, 2, 2} = 2.

Linear lifetime chart

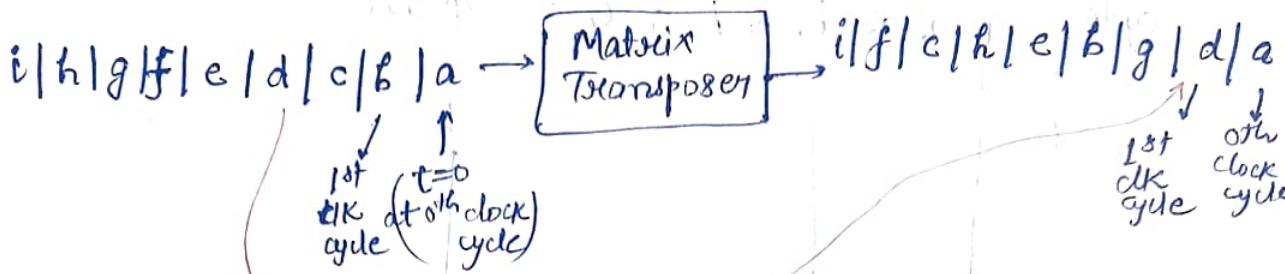




3x3 Matrix Transpose

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \Rightarrow \begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix}$$

1 iteration \rightarrow 9 clock cycles



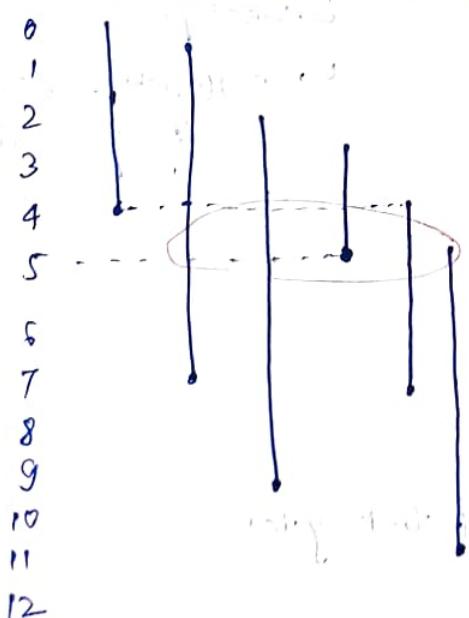
| Sample | <u>T_{input}</u> | <u>T_{zout}</u> | <u>T_{diff}</u> | <u>T_{output}</u> | <u>T_{zout} + latency</u> | <u>Life period (T_{input} → T_{output})</u> |
|--------|--------------------------|-------------------------|-------------------------|---------------------------|-----------------------------------|---|
| | | | | | | |
| a | 0 | 0 | 0 | 4 | 0 → 4 | |
| b | 1 | 3 | 2 | 7 | 1 → 7 | |
| c | 2 | 6 | 4 | 10 | 2 → 10 | |
| d | 3 | 1 | (-2) | 5 | 3 → 5 | |
| e | 4 | 4 | 0 | 8 | 4 → 8 | |
| f | 5 | 7 | 2 | 11 | 5 → 11 | |
| g | 6 | 2 | (-4) | 6 | 6 → 6 | |
| h | 7 | 5 | (-2) | 9 | 7 → 9 | |
| i | 8 | 8 | 0 | 12 | 8 → 12 | |

Illegal → add latency

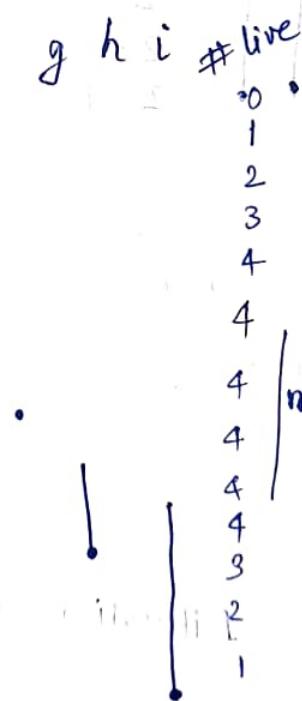
↳ To force the system to be causal, a latency is added to the system.

↳ Latency to be added = magnitude of the most negative value of T_{diff}:

a b c d e f g h i # live



max # live = 4 registers



Forward-Backward Register Allocation Technique

Steps:

- ① Determine the minm no. of registers using lifetime analysis.
- ② Input each variable at the time step corresponding to the beginning of its lifetime. If multiple variables are input in a given cycle, these are allocated to multiple registers with preference given to the variable with the longest lifetime.
- ③ Each variable is allocated in a forward manner until it is dead or it reaches the last register. In forward allocation, if the Register 'i' holds the variable in the current cycle, then register (i+1) holds the same variable in the next cycle. If (i+1)th register is not free then use the first available forward register.
- ④ Being periodic, the allocation repeats in each iteration. So hash out the register R_j for the cycle $i+N$ if it holds a variable during cycle i .
- ⑤ For variables that reach the last register and are still alive, they are allocated in a backward manner on a first come first serve basis.
- ⑥ Repeat steps 4 and 5 until, the allocation is complete.

E.g. Register minimization:

$$a = \{1, 2, 3, 4\}$$

$$b = \{2, 3, 4, 5, 6, 7\}$$

$$c = \{5, 6, 7\}$$

[$N=6$]

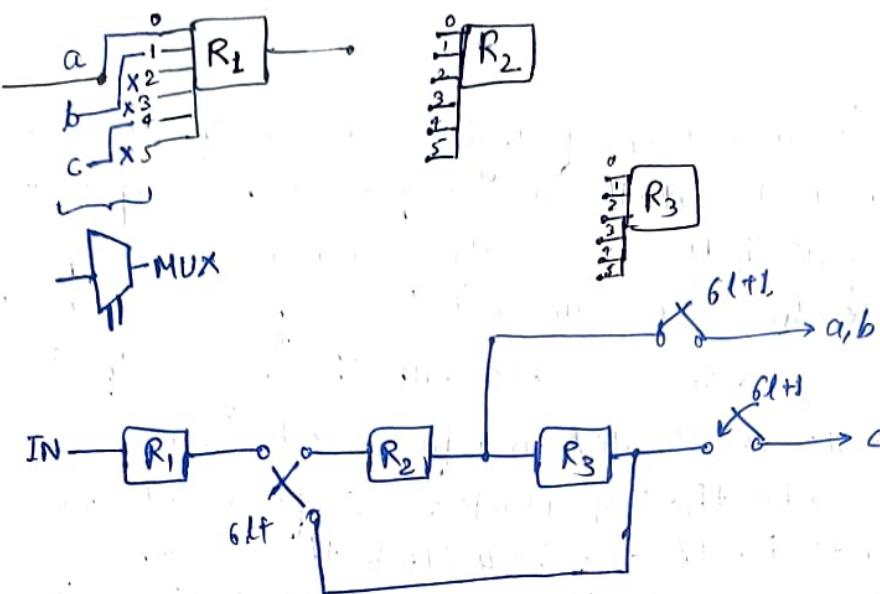
Min^m no. of Registers = 3. [Single iteration]

Register Allocation:

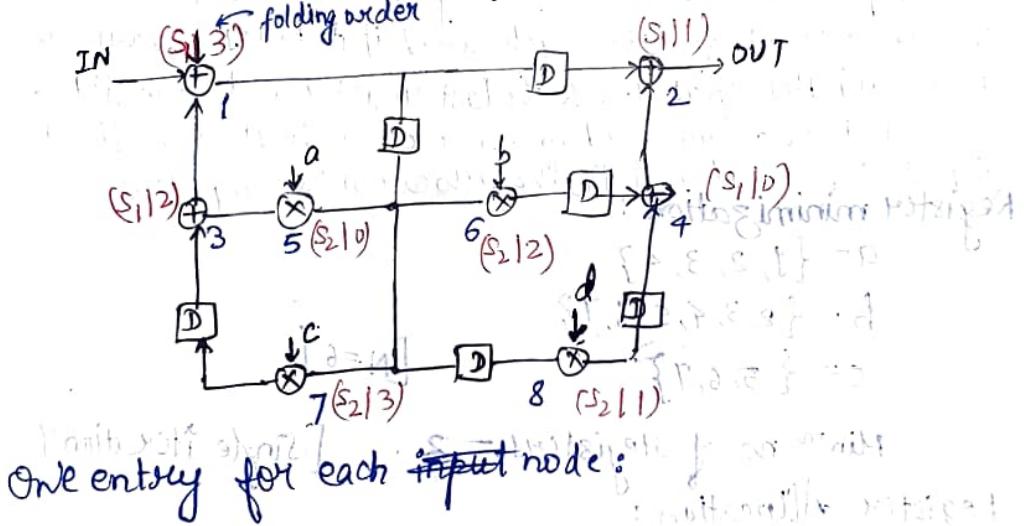
| cycle | input | R_1 | R_2 | R_3 | output |
|-------|-------|-------------------------|-------|-------|--------|
| 0 | a | a ready by 2nd cycle | | | |
| 1 | b | a | | | |
| 2 | | b ready by 1st cycle | a | | |
| 3 | | | a | b | a |
| 4 | c | | b | a | |
| 5 | | c | b | a | |
| 6 | | c | b | a | |
| 7 | | | b | c | b, c |

Prepare for next input (a₂)

Forward \rightarrow preference
↓
Backward



Register Minimization of Biquad filter



$$T_{\text{input}} = u + p_u$$

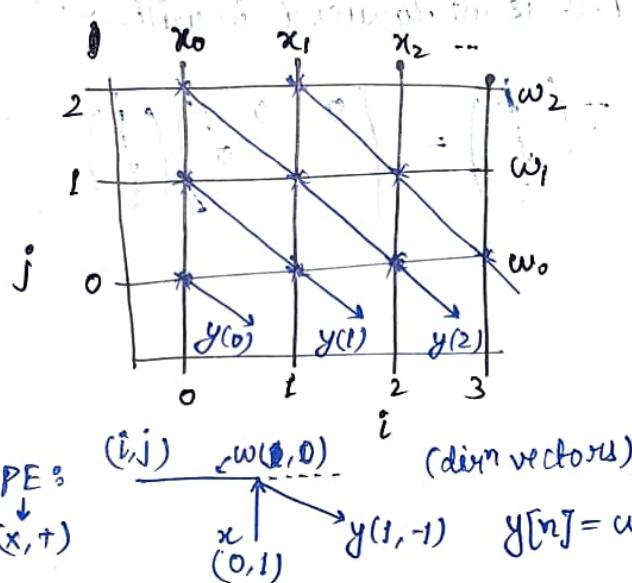
$$T_{\text{output}} = u + p_u + \max \{ D_F(u \rightarrow v) \}$$

Systolic Array

→ Systolic architectures are designed by using linear mapping on regular dependence graphs (DG).

Regular Dependence Graph: The presence of an edge in a certain direction at any node in the DG represents the presence of an edge in the direction at all nodes in the DG.

$$\text{Eg. } y[n] = w_0x[n] + w_1x[n-1] + w_2x[n-2]$$



$$\text{PE: } (i, j) \xrightarrow{w(i, j)} y(i, -j) \quad \text{dim vector } u$$

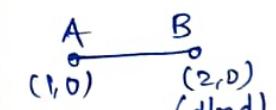
$$\text{Linear mapping: } (i, j) \xrightarrow{\quad} (j', t')$$

$\begin{cases} w \\ x \\ y \end{cases}$ } delay or no delay

$$j' = i + j, \quad t' = i$$

Basic vectors:

→ Projection vector / iteration vector : $d = \begin{pmatrix} d_1 \\ d_2 \end{pmatrix}$



→ Two nodes that are displaced by 'd' at multiples of 'd' are executed by the same processor.

→ Processor space vector : $P^T = (P_1 \ P_2)$

$$(0') \rightarrow (1, 0)$$

Any node with index $I^T = (i, j)$ would be executed by processor

$$P^T I = (P_1 \ P_2) \begin{pmatrix} i \\ j \end{pmatrix}$$

→ Scheduling vector: $S^T = (s_1, s_2)$

Any node with index i would be executed at time $s^T i$.

→ Hardware Utilization Efficiency, $HUE = \frac{1}{|STD|}$

If A is mapped as processor of B.

⇒ Can't be executed at same time.

$$S^T i_A \neq S^T i_B, \text{ i.e., } S^T d \neq 0.$$

Edge mapping: $P^T e$ is introduced in with $S^T e$ delay.

$$\begin{pmatrix} i' \\ j' \\ t' \end{pmatrix} = T \begin{pmatrix} i \\ j \\ t \end{pmatrix} = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} i \\ j \\ t \end{pmatrix}$$

FIR Filter Design B₁ (Broadcast inputs, Move results, weights stay)

Distance vector, $d^T = (1, 0)$

$$p^T = (0, 1) \quad [\text{choose}] \quad [\because p^T d = 0] \quad [\text{or } p^T = (0, k)]$$

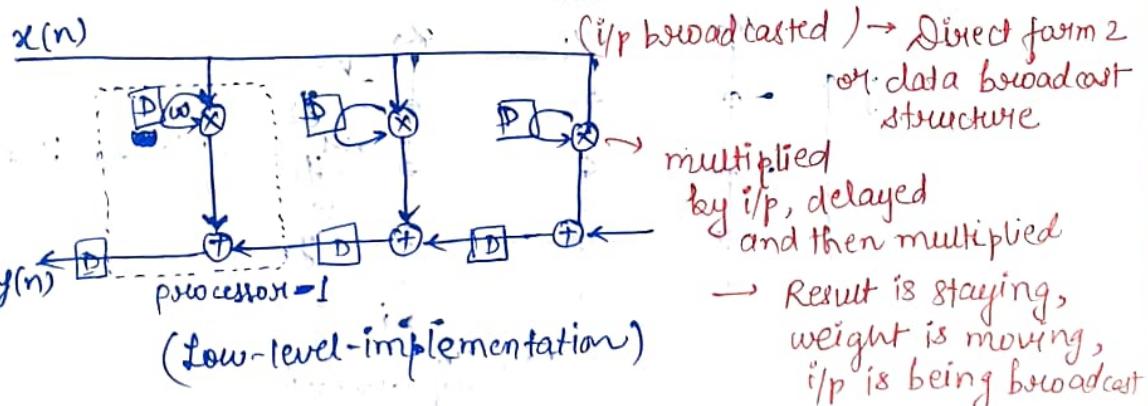
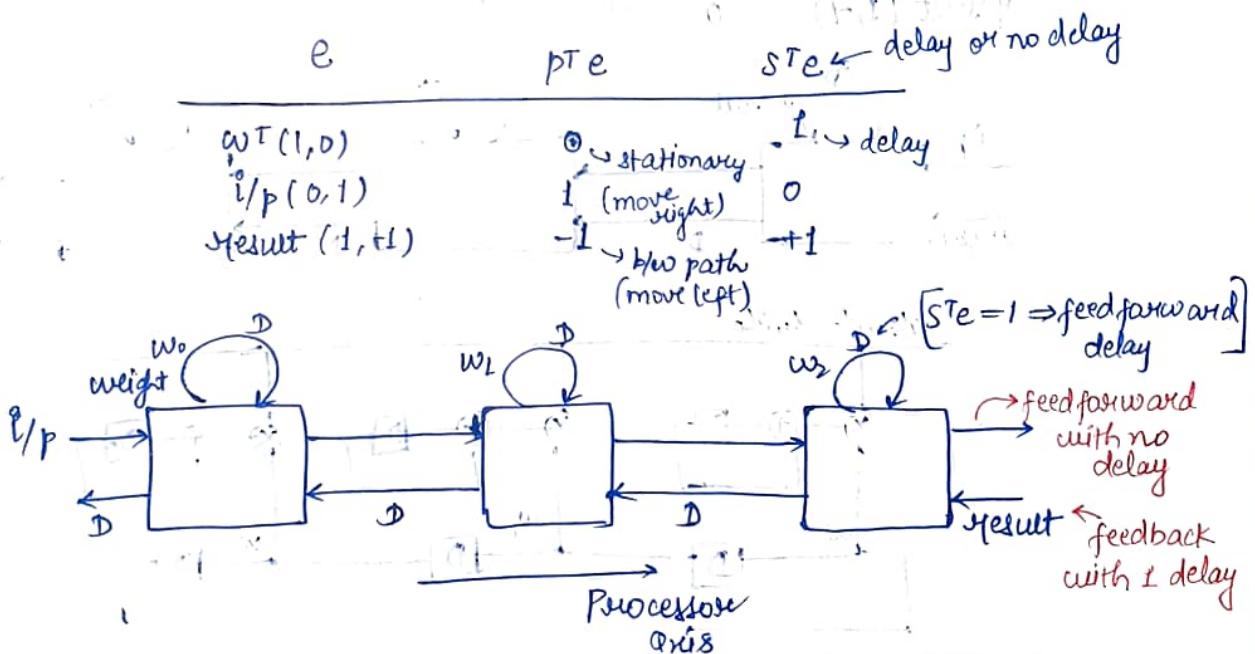
$$s^T = (1, 0) \quad [\because s^T d \neq 0] \quad [\text{or } s^T = (k_1, k_2), \\ k_1 \neq 0, \\ k_2 \rightarrow \text{any}]$$

Any node with index $I^T = (i, j)$

- is mapped to processor $p^T I = j$.
- is executed at time $s^T I = i$.

Since, $s^T d = 1$, we have $HUE = \frac{1}{|s^T d|} = 1$

Edge mapping:



Design B₂:

$$d^T = (1, -1)$$

$$p^T = (1, 1)$$

$$s^T = (1, 0)$$

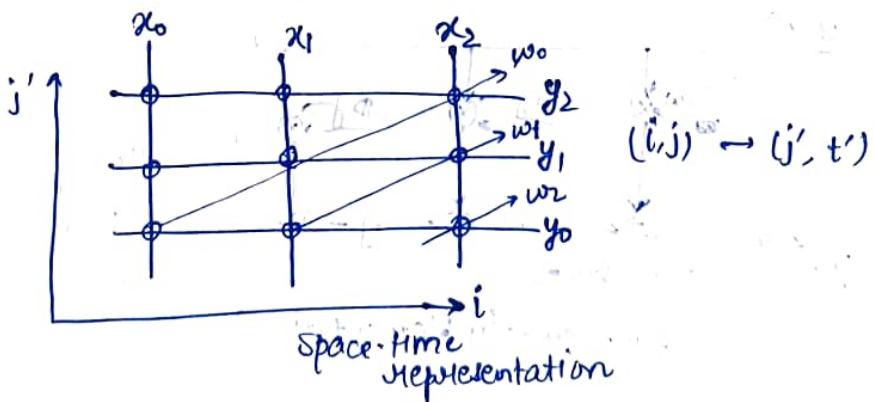
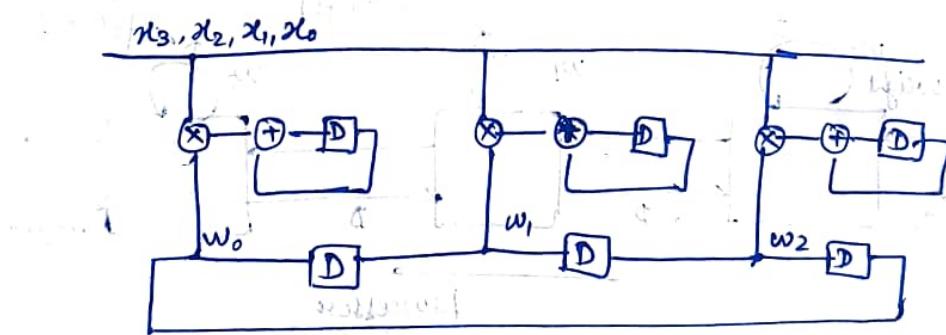
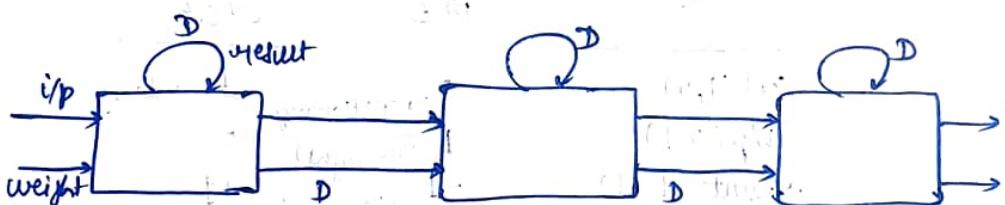
$$p^T I = i + j$$

$$s^T I = i$$

$$HUE = \int s^T dI = 1$$

03-04-2025

| e | $p^T e$ | $s^T e$ |
|------------------|-------------------------------------|---------|
| $w^T(1, 0)$ | $1 \rightarrow (\text{move right})$ | 1 |
| $i/p(0, 1)$ | $1 \rightarrow$ | 0 |
| Result $(1, -1)$ | $0 \rightarrow (\text{stays})$ | 1 |



Design F :

$$d^T = (1, 0), p^T = (0, 1), s^T = (1, 1)$$

$$|s^T d| = 1 \Rightarrow HDE = 1$$

e. $p^T e$ $s^T e$

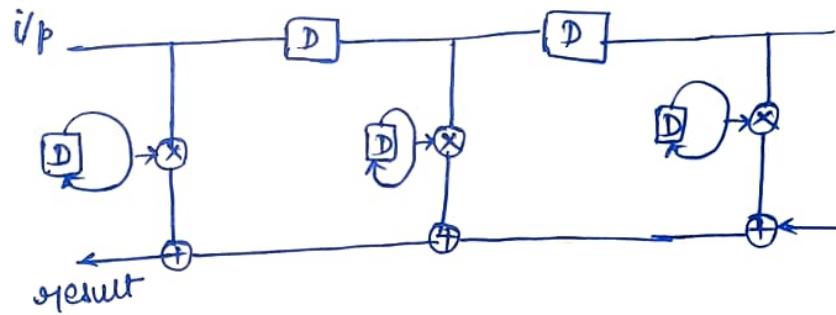
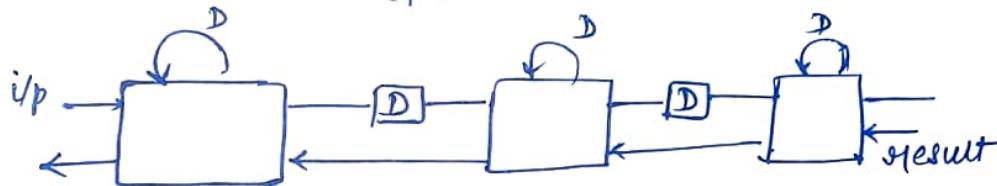
| | | |
|-----------------|----|---|
| $w^T(1, 0)$ | 0 | 1 |
| $i/p(0, 1)$ | 1 | 0 |
| $result(1, -1)$ | -1 | 0 |

↓
move
(left)

J  null space

$$d = \begin{bmatrix} i \\ j \\ k \end{bmatrix}$$

$$p^T \rightarrow 2 \times 3$$

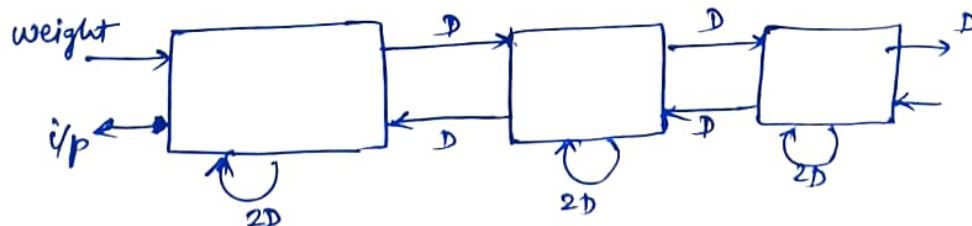


Design R₁:

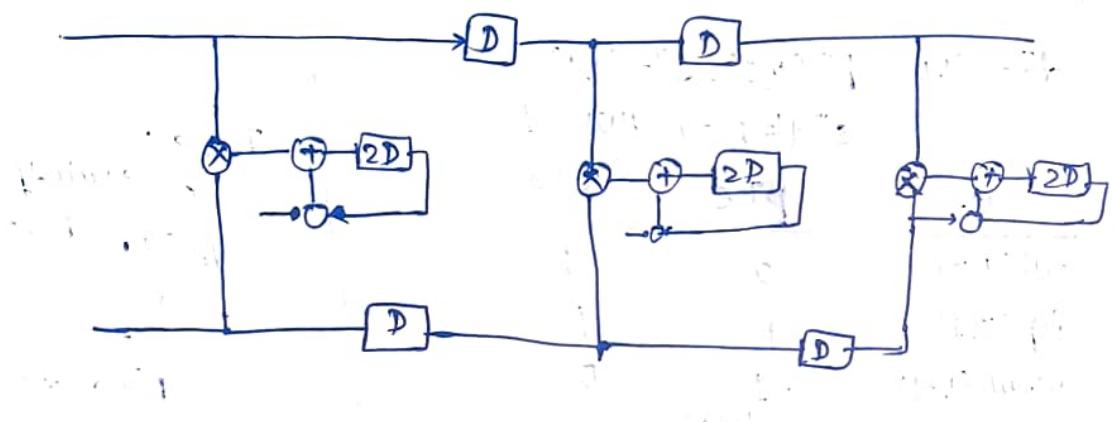
$$d^T = (1, -1), p^T = (1, 1), s^T = (1, -1)$$

$$|s^T d| = 2 \Rightarrow HDE = \frac{1}{2}.$$

| e | $p^T e$ | $s^T e$ |
|-----------------|---------|---------|
| $w^T(1, 0)$ | 1 | 1 |
| $i/p(0, -1)$ | -1 | 1 |
| $result(1, -1)$ | 0 | 2 |



$$\text{Hardware utilization} = \frac{1}{s^T d} = \frac{1}{2} = 50\% \text{ utilization.}$$



Block diagram of a signal processing system.



Block diagram of a signal processing system.



Block diagram of a signal processing system.

Block diagram of a signal processing system.

CONVOLUTION

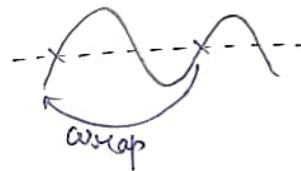
Linear Convolution

$$y[n] = x[n] * h[n] = \sum x[n] \cdot h[n-k]$$

$$= \sum h[n] \cdot x[n-k]$$

Circular Convolution

↳ For periodic signals
 $x[n] * h[n]$ modulo operations



Fast Convolution

Systolic
 ↳ Latency ↑
 ↓
 Throughput ↑

Operations:

Addition

Multiplication → computationally intensive (Repeated addition)

↳ More hardware resource

↳ More time

↳ Reduce multiplication at the expense of increasing addition.

"Arithmetic Strength Reduction"

Polynomial Multiplication

$$x[n] = [3, 5] \quad h[n] = [1, 2, 4]$$

↓

$$3x+5$$

$$\begin{array}{r} & 1 & 2 & 4 \\ \times & 3 & & \\ \hline & 3 & 6 & 12 \\ \hline & 5 & & \\ \hline & 5 & 10 & 20 \end{array}$$

$$y[n] = x(n) * h(n) = [3, 6, 15, 10 + 12, 20]$$

$$y[n] = [3, 11, 22, 20]$$

Polynomial wise: $(3x+5)(x^2+2x+4)$

$$= 3x^3 + 6x^2 + 12x + 5x^2 + 10x + 20$$

$$= 3x^3 + 11x^2 + 22x + 20$$

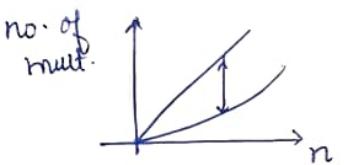
$$x[n] \longleftrightarrow X[z]$$

$$h[n] \longleftrightarrow H[z]$$

$$Y[z] = X[z] \cdot H[z]$$

No. of multiplications: $3 \times 2 = 6$

No. of additions: $2 \times 2 = 4$



Algorithms:

① Cooks and Toom Algorithm

② Winograd Winograd

↳ Chinese Remainder theorem

Lagrange interpolation

Complex Multiplication:

$$(a+jb)(c+jd) = e+jf$$

$$\begin{bmatrix} e \\ f \end{bmatrix} = \begin{bmatrix} c & -d \\ d & c \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

↳ 4 multiplications
↳ 2 additions.

$$\begin{aligned} ac - bd &= ac - ad + ad - bd \\ &= a(\overset{(M_1)}{c-d}) + d(\overset{(M_2)}{a-b}) \end{aligned}$$

$$\begin{aligned} ad + bc &= bc + bd - bd + ad \\ &= b(c+d) + d(a-b) \end{aligned}$$

$\overset{(M_3)}{}$ $\overset{(M_4)}{}$ $\overset{(M_5)}{}$

↳ 3 multiplications
+
5 additions

$$\begin{bmatrix} e \\ f \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}_{2 \times 3} \begin{bmatrix} (c-d)a \\ (c+d)b \\ d(a-b) \end{bmatrix}_{3 \times 1}$$

$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} c-d & 0 & 0 \\ 0 & c+d & 0 \\ 0 & 0 & d \end{bmatrix} \begin{bmatrix} a \\ b \\ a-b \end{bmatrix}$$

\downarrow
Post-computation
matrix

\downarrow
Pre-computed

\downarrow
computed
on-the-fly

no need to
multiply with it,
just choose
operands

\hookrightarrow 3 multiplications + 3 additions
(reduced from 5)

Lagrange interpolation polynomial:

$$y = f(x)$$

$$P_0(x) = \frac{(x-x_1)(x-x_2) \dots (x-x_N)}{(x_0-x_1)(x_0-x_2) \dots (x_0-x_N)}$$

if x_0 is known

Substitute $x = x_0$,

$$\begin{aligned} P_0(x) &= 1, & x = x_0 \\ &= 0, & x \neq x_0. \end{aligned}$$

$$\text{In general, } P_i(x) = \prod_{\substack{j=0 \\ i \neq j}}^N \frac{x-x_j}{x_i-x_j}$$

$$\begin{aligned} P(x_j) &= 1, & i=j \\ &= 0, & i \neq j \end{aligned}$$

$$\therefore f(x) = y_0 P_0(x) + y_1 P_1(x) + \dots + y_n P_n(x)$$

Lagrange interpolation formula: $f(x) = \sum_{i=0}^n y_i \cdot P_i(x)$

Lagrange interpolation formula

$$y[n] = h(n) \cdot x(n)$$

$$S(p) = h(p) \cdot x(p)$$

2x2

$$h = [h_0 \ h_1], \quad x = [x_0 \ x_1]$$

Cook Toom Algorithm:

- ① Choose $L+N-1$ different real numbers $\beta_0, \dots, \beta_{L+N-2}$.
(small values)
 - ② Compute $h(\beta_i) + x(\beta_i)$.
 - ③ Compute $s(\beta_i) = h(\beta_i) \cdot x(\beta_i)$.
 - ④ Compute $s(p)$.

Eg.

$$h(p) = h_0 + h_1 p$$

$$\mathcal{H}((P) = x_0 + x_1 P$$

$$S(p) = S_0 + S_1 p + S_2 p^2$$

$$= h_0 x_0 + \underset{M_1}{(h_1 x_0 + h_0 x_1)} p + \underset{M_2}{h_1 x_1} p^2 \rightarrow 4 \text{ Multip.} \\ + 3 \text{ Additions}$$

$$\begin{bmatrix} S_0 \\ S_1 \\ S_2 \end{bmatrix} = \begin{bmatrix} h_0 x_0 \\ h_1 x_0 + h_0 x_1 \\ h_1 x_1 \end{bmatrix}$$

$$= \begin{bmatrix} h_0 & 0 \\ h_1 & h_0 \\ 0 & h_1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix}_{3 \times 2}^{2 \times 1}$$

$$\text{or, } S = \boxed{T X}$$

Cook-Toom Algorithm:

$$\textcircled{1} \quad L+N-1 = 2+2-1 = 3$$

Given that $\beta_0 = 0$

$$\beta_1 = 1$$

$$b_2 = -1$$

② $h(p_i)$

| β_i | $h(\beta_i)$ | $x(\beta_i)$ | $s(\beta_i) = h(\beta_i) \cdot x(\beta_i)$ |
|-----------|--------------|--------------|--|
| β_0 | h_0 | x_0 | (M1) $h_0 x_0$ |
| β_1 | $h_0 + h_1$ | $x_0 + x_1$ | (M2) $(h_0 + h_1)(x_0 + x_1)$ (A1) \rightarrow on-the-fly add. |
| β_2 | $h_0 - h_1$ | $x_0 - x_1$ | (M3) $(h_0 - h_1)(x_0 - x_1)$ (A2) \rightarrow |

$$\# y = p_0 \cdot y_0 + p_1 \cdot y_1$$

$$s(\beta_0) \downarrow s(\beta_1)$$

Lagrange Interp.

$$\begin{aligned} & \hookrightarrow x \rightarrow P \\ & x_i \rightarrow \beta_i, y \rightarrow s(\beta) \end{aligned}$$

$$\begin{aligned} s(P) &= s(\beta_0) \cdot \frac{(P-\beta_1)(P-\beta_2)}{(\beta_0-\beta_1)(\beta_0-\beta_2)} + s(\beta_1) \cdot \frac{(P-\beta_0)(P-\beta_2)}{(\beta_1-\beta_0)(\beta_1-\beta_2)} \\ &\quad + s(\beta_2) \cdot \frac{(P-\beta_0)(P-\beta_1)}{(\beta_2-\beta_0)(\beta_2-\beta_1)} \\ &= s(\beta_0) \cdot \frac{(P-1)(P+1)}{2} + s(\beta_1) \cdot \frac{(P)(P+1)}{2} + s(\beta_2) \cdot \frac{(P)(P-1)}{2} \\ &= s(\beta_0) + \left(\frac{s(\beta_1) - s(\beta_2)}{2} \right) \cdot P + \left(-s(\beta_0) + \frac{s(\beta_1) + s(\beta_2)}{2} \right) P^2 \end{aligned}$$

$$\therefore s(P) = s_0 + s_1 P + s_2 P^2$$

$$\begin{bmatrix} s_0 \\ s_1 \\ s_2 \end{bmatrix} = \begin{bmatrix} s(\beta_0) \\ \frac{s(\beta_1) - s(\beta_2)}{2} \\ -s(\beta_0) + \frac{s(\beta_1) + s(\beta_2)}{2} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} h(\beta_0) \cdot x(\beta_0) \\ \frac{h(\beta_1) \cdot x(\beta_1)}{2} \\ h(\beta_2) \cdot x(\beta_2) \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} h(\beta_0) & 0 & 0 \\ 0 & \frac{h(\beta_1)}{2} & 0 \\ 0 & 0 & \frac{h(\beta_2)}{2} \end{bmatrix} \begin{bmatrix} x_0 \\ x_0 + x_1 \\ x_0 - x_1 \end{bmatrix}$$

(3x3) (3x1)

$$= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} h(B_0) & 0 & 0 \\ 0 & \frac{h(B_1)}{2} & 0 \\ 0 & 0 & \frac{h(B_2)}{2} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 1 & -1 \\ (3 \times 2) \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$$

(3x3)

↓

Post comp.
(add.)

Pre-computation
(addition)

$$\therefore S = A H \underbrace{B X}_{\text{diagonal matrix}}$$

(post-computation) pre-computation (preeaddition)

$$\rightarrow H_0 = h_0$$

$$H_1 = (h_0 + h_1)/2$$

$$H_2 = (h_0 - h_1)/2$$

} (pre-computed)

$$X_0 = x_0$$

$$X_1 = x_0 + x_1$$

$$X_2 = x_0 - x_1$$

$$S_0 = H_0 X_0$$

$$S_1 = H_1 X_1$$

$$S_2 = H_2 X_2$$

$$\therefore S_1 = S_0$$

$$S_1 = S_1 - S_2$$

$$S_2 = -S_0 + S_1 + S_2$$

} + ③ Additions

\therefore Total $\rightarrow 3M + 5A \rightarrow 1M$ reduction, $4A$ increment.

Conventionally, $4M + 1A$

(ignoring additions in pre-computation)

Generally,

$$S = T \alpha = A H B \alpha$$

Since $T = A H B$, CT algorithm provides a way to factorize the convolution matrix T into 3 multiplying matrices and the total no. of multiplications is determined by the non-zero elements on the main diagonal of the matrix H .

[C and D contain only small integers.]

→ Still, the no. of additions has increased.

⇒ Use modified CT algorithm.

Modified Cook-Toom Algorithm

$$S'(p) = S(p) - S_{L+N-2} \cdot p^{L+N-2}$$

$\downarrow L+N-3$

- ① choose $L+N-2$ real nos. $\beta_0, \dots, \beta_{L+N-3}$.
- ② compute $h(\beta_i)$ and $x(\beta_i)$, $i=0, \dots, L+N-3$.
- ③ compute $s(\beta_i) = h(\beta_i) \cdot x(\beta_i)$ for $i=0, \dots, L+N-3$.
- ④ compute $s'(\beta_i) = s(\beta_i) - S_{L+N-2} \beta_i^{L+N-2}$, $i=0, \dots, L+N-3$.
- ⑤ $S'(p) = \sum_{i=0}^{L+N-3} s'(\beta_i) \cdot \frac{\prod_{j \neq i} (p - \beta_j)}{\prod_{j \neq i} (\beta_i - \beta_j)}$
- ⑥ $S(p) = S'(p) + S_{L+N-2} \cdot p^{L+N-2}$.

2×2 convolution Algorithm, using Modified CT with

$$\beta_0 = 0, \beta_1 = -1.$$

$$S'(p) = S(p) - h_1 x_1 p^2 \quad \text{at } \beta_0 = 0, \beta_1 = -1.$$

$$s'(\beta_i) = h(\beta_i) \cdot x(\beta_i) - h_1 x_1 \beta_i^2$$

$$\beta_0 = 0 \Rightarrow h(\beta_0) = h_0, x(\beta_0) = x_0$$

$$\beta_1 = 1 \Rightarrow h(\beta_1) = h_0 - h_1, x(\beta_1) = x_0 - x_1$$

$$s'(\beta_0) = h(\beta_0) x(\beta_0) - h_1 x_1 \beta_0^2 = h_0 x_0$$

$$s'(\beta_1) = h(\beta_1) x(\beta_1) - h_1 x_1 \beta_1^2 = (h_0 - h_1)(x_0 - x_1) - h_1 x_1$$

$$S'(p) = s'(\beta_0) + p(s'(\beta_0) - s'(\beta_1))$$

$$\therefore S(p) = S'(p) + h_1 x_1 p^2$$

$$= S_0 + S_1 p + S_2 p^2 = S_0 + S_1 p$$

$$\begin{bmatrix} S_0 \\ S_1 \\ S_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & -1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} h_0 & 0 & 0 \\ 0 & h_0 - h_1 & 0 \\ 0 & 0 & h_1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$$

$$S_0 = H_0 X_0, S_1 = H_1 X_1, S_2 = H_2 X_2$$

$$\Rightarrow S_0 = S_0, S_1 = S_0 - S_1, S_2 = S_2$$

Chinese Remainder Theorem (CRT) for integer:

Unknown: x

$x \% m_0 \rightarrow a_0 \rightarrow$ remainder

$x \% m_1 \rightarrow a_1$

$x \% m_n \rightarrow a_n \rightarrow$ remainder

↳ non-negative integer \rightarrow uniquely

↳ Remainder \rightarrow w.r.t. modulus under the condition

$m_0, m_1, m_2, \dots, m_n \rightarrow$ relatively prime
and, $x < m_0 \cdot m_1 \cdot m_2 \cdots m_n$ (coprime)

$x \equiv a_1 \pmod{m_1}$
 $x \equiv a_2 \pmod{m_2}$
 \vdots
 $x \equiv a_n \pmod{m_n}$

} $\rightarrow x$ is smaller than the product of moduli
 $x < m_1 \cdot m_2 \cdots m_n$
 set of congruence \rightarrow solve simultaneously to get the value of x .

CRT for integer :

$c_i = R_{m_i}[c]$, for $i=0, 1, 2, \dots, k$.
 Remainder moduli \downarrow Unknown integer \uparrow $R_{m_i}[c] \rightarrow$ Remainder when $c \div m_i$

$\text{GCD}(m_0, m_1, \dots, m_n) = 1$, i.e., m_i s are relatively prime.
 $c \ll m_0, m_1, \dots, m_n$.

$$c = \sum_{i=0}^k c_i N_i M_i \pmod{M}, \quad M = \prod_{i=0}^k m_i,$$

$$\therefore M_i = \frac{M}{m_i},$$

N_i is the solution of the equation $N_i M_i + n_i m_i = \text{GCD}(M_i, m_i) = 1$.

Solve using Euclidean GCD algorithm.

If $A \nmid B$,

$$\text{GCD}(A, B) = G = A R + B S = 1$$

Eg.

Given $m_0 = 3, m_1 = 4, m_2 = 5$

$$M = 3 \times 4 \times 5 = 60$$

$$M_0 = \frac{M}{m_0} = 20, \quad M_1 = \frac{M}{m_1} = 15, \quad M_2 = \frac{M}{m_2} = 12.$$

$$N_i M_i + n_i m_i = \text{GCD}(M_i, m_i).$$

$$i=0: \quad N_0(20) + n_0(3) = \text{GCD}(20, 3) = 1$$

$$\Rightarrow N_0 = -1, n_0 = 7$$

$$i=1: \quad N_1(15) + n_1(4) = \text{GCD}(15, 4) = 1$$

$$\Rightarrow N_1 = -1, n_1 = 4$$

$$i=2: \quad N_2(12) + n_2(5) = 1$$

$$\Rightarrow N_2 = -2, n_2 = 5.$$

$$c = \sum_{i=0}^2 (c_i N_i M_i) \pmod{M}$$

$$\text{Let } c = 17: \quad 17 \pmod{3} = 2 \rightarrow c_0$$

$$17 \pmod{4} = 1 \rightarrow c_1$$

$$17 \pmod{5} = 2 \rightarrow c_2$$

$$\therefore c = (c_0 N_0 M_0 + c_1 N_1 M_1 + c_2 N_2 M_2) \pmod{M}$$

$$= (2 \times 20 \times -1 + 1 \times -1 \times 15 + 2 \times -2 \times 12) \pmod{60}$$

$$= (40 - 15 - 48) \pmod{60} = (-103) \pmod{60} = \boxed{17}.$$

CRT algorithm for polynomial:

11-04-2025

$$c_i(p) = R_{m^i}(p) [C(p)], \quad i=0, 1, \dots, K.$$

↓
remainder
polynomial
moduli
polynomial
unknown
polynomial

$$R_{m^i}(p) \rightarrow \text{Remainder } \left[\frac{C(p)}{m^i(p)} \right].$$

$$M(p) = \prod_{i=0}^K m^i(p) \rightarrow \text{relatively prime}, \quad M^i(p) = \frac{M(p)}{m^i(p)}.$$

$$c(p) = \sum_{i=0}^K (c_i(p) N^i(p) M^i(p)) \bmod M(p)$$

$$s(p) = h(p) \times x(p)$$

$$N^i(p) \cdot M^i(p) + n^i(p) \cdot m^i(p) = \text{GCD}[M^i(p), m^i(p)] = 1,$$

provided that the degree of $c(p)$ should be less than the degree of $M(p)$.

→ The remainder of polynomial w.r.t. moduli polynomial,
 $p_i + f(p)$ when

$$\text{degree of } f(p) \leq i-1$$

e.g. poly. $5x^2 + 3x + 5$.

④ moduli $\rightarrow (x+2)$
 $\text{remainder} = 5(-2)^2 + 3(-2) + 5$

can be obtained by substituting p_i by $-f(p)$ within polynomial.

⑤ moduli $\rightarrow (x^2+2)$, $\text{Rem.} = 5(-2) + 3x + 5 = 3x - 5$

11-04-2025

Eg. Consider 2x3 linear convolution.

Construct an effective realization using Winograd algorithm with

$$h(p) = h_0 + h_1 p$$

$$x(p) = x_0 + x_1 p + x_2 p^2$$

$$\text{Given: } m(p) = p(p-1)(p^2+1)$$

$$m^{(0)}(p) = p, \quad m^{(1)}(p) = p-1, \quad m^{(2)}(p) = p^2+1$$

$$\therefore m(p) = m^{(0)}(p) + m^{(1)}(p)p + m^{(2)}(p)p^2$$

$$M^{(i)}(p) = \frac{m(p)}{m^{(i)}(p)}$$

quotient remainder

$$\Rightarrow M^{(0)}(p) = (p-1)(p^2+1), = p^3 - p^2 + p - 1 = p \underbrace{(p^2 - p + 1)}_{\text{quotient}} - 1$$

$$M^{(1)}(p) = p(p^2+1),$$

$$M^2(p) = p(p-1).$$

and,

$$N^{(i)}(p)M^{(i)}(p) + n^{(i)}(p)m^{(i)}(p) = 1, \quad i=0,1,2$$

1st iteration ($i=0$):

$$\Rightarrow N^0(p)M^0(p) + n^0(p)m^0(p) = 1$$

$$\Rightarrow N^0(p)(p-1)(p^2+1) + n^0(p).p = 1$$

$$\Rightarrow N^0(p) = -1,$$

$$n^0(p) = p^2 - p + 1$$

| i | $m^i(p)$ | $M^i(p)$ | $n^i(p)$ | $N^i(p)$ |
|---|-----------|---------------------|-----------------------------|--------------------|
| 0 | p | $p^3 - p^2 + p - 1$ | $p^2 - p + 1$ | -1 |
| 1 | $p-1$ | $p^3 + p$ | $\frac{-1}{2}(p^2 + p + 2)$ | $\frac{1}{2}$ |
| 2 | $p^2 + 1$ | $p^2 - p$ | $\frac{-1}{2}(p-2)$ | $\frac{1}{2}(p-1)$ |

obtained using:

$$\therefore M(p) = PQ + R$$

$$P \left[\begin{array}{c} p^2 - p + 1 \\ p^3 - p^2 + p - 1 \\ p^3 \end{array} \right]$$

$$\begin{array}{c} -p^2 \\ -p^2 \\ \hline p-1 \\ \hline p \\ \hline -1 \end{array}$$

$$M^0(p) = m^0(p)(p^2 - p + 1) + (-1)$$

$$\Rightarrow m^0(p)(p^2 - p + 1) - M^0(p) = 1,$$

$$\text{and, } N^0(p) \cdot (p^3 - p^2 + p - 1) + n^0(p) \cdot p = 1$$

$$\Rightarrow m^0(p)(p^2 - p + 1) - M^0(p) = N^0(p)(p^3 - p^2 + p - 1) + n^0(p) \cdot p$$

$$\therefore N^{(0)}(p) = -1$$

$$n^0(p) = p^2 - p + 1$$

2nd iteration : (i=1)

$$M^1(p) = m^1(p) \cdot (p^2 + p + 2) + 1$$

$$\Rightarrow \frac{M^1(p)}{2} = \frac{m^1(p)}{2} \cdot (p^2 + p + 2) + 1$$

$$\Rightarrow J = \frac{M^1(p)}{2} - \frac{m^1(p)}{2} \cdot (p^2 + p + 2)$$

$$\text{and, } N^1(p) \cdot M^1(p) + n^1(p) \cdot m^1(p) = 1$$

$$\Rightarrow N^1(p) \cdot M^1(p) + n^1(p) \cdot m^1(p) = \frac{M^1(p)}{2} - \frac{m^1(p)}{2} (p^2 + p + 2)$$

$$\therefore N^1(p) = \frac{1}{2}$$

$$n^1(p) = \frac{1}{2} (p^2 + p + 2)$$

$$\begin{array}{r} p^2 + p + 2 \\ \hline p-1 | p^3 + p \\ \quad \quad \quad p^3 - p^2 \\ \hline \quad \quad \quad + p^2 + p \\ \quad \quad \quad p^2 - p \\ \hline \quad \quad \quad 2p - 2 \\ \quad \quad \quad \hline 2p - 2 \end{array}$$

3rd iteration : (i=2)

$$M^2(p) = m^2(p) (1) + (-p-1)$$

$$\Rightarrow -M^2(p) + m^2(p) -p = 1 \quad \cancel{+ p^2 + p + 2 + 2p}$$

$$\text{and, } N^2(p) M^2(p) + n^2(p) m^2(p) = 1$$

\Rightarrow since remainder is a fm of P.

$$\therefore m^2(p) = E(p+1) Q + r$$

$$\Rightarrow m^2(p) = E(p+1) (-p+1) + 2 \quad \begin{matrix} \downarrow \text{not a fm of P} \\ \text{stop!} \end{matrix}$$

$$\Rightarrow \frac{m^2(p)}{2} + \frac{(p+1)(-p+1)}{2} = 1.$$

$$\Rightarrow + (p+1) = m^2(p) \cdot 1 - M^2(p)$$

$$\therefore \frac{m^2(p)}{2} + (-p+1) \left[\frac{m^2(p) - M^2(p)}{2} \right] = 1 \quad \}$$

$$\text{and, } N^2(p) M^2(p) + n^2(p) m^2(p) = 1 \quad \}$$

$$\begin{array}{r} 1 \\ \hline p^2 - p \\ \quad \quad \quad p^2 + 1 \\ \hline \quad \quad \quad -p - 1 \end{array}$$

$$\begin{array}{r} -p + 1 \\ \hline p^2 + 1 \\ \quad \quad \quad p^2 + p \\ \hline \quad \quad \quad -p + 1 \\ \quad \quad \quad -p - 1 \\ \hline \quad \quad \quad 2 \end{array}$$

$$\Rightarrow N^2(p) = \frac{1}{2} (p-1), \quad n^2(p) = \frac{1}{2} (-p+2).$$

Now,

$$\left. \begin{array}{l} h^i(p) = h(p) \bmod m^i(p) \\ x^i(p) = x(p) \bmod m^i(p) \end{array} \right\} \text{for } i=0,1,2.$$

$$\Rightarrow h^0(p) = (h_0 + h_1 p) \bmod p = h_0$$

$$h^1(p) = (h_0 + h_1 p) \bmod (p-1) = h_0 + h_1$$

$$h^2(p) = (h_0 + h_1 p) \bmod (p^2+1) = (h_0 + h_1 p)$$

degree > degree of $(h_0 + h_1 p)$

and,

$$\Rightarrow \text{Rem.} = (h_0 + h_1 p)$$

$$x^0(p) = (x_0 + x_1 p + x_2 p^2) \bmod p = x_0$$

$$x^1(p) = (x_0 + x_1 p + x_2 p^2) \bmod (p-1) = x_0 + x_1 + x_2$$

$$x^2(p) = (x_0 + x_1 p + x_2 p^2) \bmod (p^2+1) = x_0 + x_1 p - x_2$$

$$= (x_0 - x_2) + x_1 p$$

and,

$$s_i(p) = h^i(p) \cdot x^i(p), \text{ for } i=0,1,2.$$

$$\Rightarrow s^0(p) = h_0 x_0 = s_0^{(0)} \rightarrow 1 \text{ mult.}$$

$$s^1(p) = (h_0 + h_1)(x_0 + x_1 + x_2) = s_0^{(1)} \rightarrow 1 \text{ mult.}$$

$$s^2(p) = (h_0 + h_1 p)(x_0 - x_2 + x_1 p)$$

$$= \underbrace{h_0(x_0 - x_2)}_{s_0^{(2)}} - h_1 x_1 + \underbrace{(h_0 x_1 + h_1(x_0 - x_2))p}_{s_1^{(2)}} \rightarrow 4 \text{ mult.}$$

Rewrite, \rightarrow not a fn of p

$$s_0^{(2)} = h_0(x_0 + x_1 - x_2) - h_1 x_1$$

$$s_1^{(2)} = h_0(x_0 + x_1 - x_2) - (h_1 - h_0)(x_0 - x_2) \quad \} 3 \text{ mult.}$$

\therefore Total mult. required = 5.

In matrix form,

$$\begin{bmatrix} s_0^{(2)} \\ s_1^{(2)} \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 \\ 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} h_0 & 0 & 0 \\ 0 & h_1 - h_0 & 0 \\ 0 & 0 & h_0 + h_1 \end{bmatrix} \begin{bmatrix} x_0 + x_1 - x_2 \\ x_0 - x_2 \\ x_1 \end{bmatrix}$$

Now,

$$s(p) = \sum_{i=0}^2 s^{(i)}(p) N^{(i)}(p) m^{(i)}(p) \bmod m(p)$$

$$s_0 + s_1 p + s_2 p^2 + s_3 p^3 = [-s^{(0)}(p)(p^3 - p^2 + p - 1) + s^{(1)}(p) \cdot (p^3 + p)]$$

$$+ \frac{s^{(2)}(p)(p^3 - 2p^2 + p)}{2} \bmod [p^4 - p^3 + p^2 - p]$$

$$\Rightarrow S_0 + S_1 p + S_2 p^2 + S_3 p^3 = p^0 S_0^{(0)} + p^1 \left[-S_0^{(0)} + \frac{S_0^{(1)}}{2} + \frac{S_0^{(2)}}{2} + \frac{S_1^{(2)}}{2} \right] \\ + p^2 \left[S_0^{(0)} - S_0^{(2)} \right] \\ + p^3 \left[-S_0^{(0)} + \frac{S_0^{(1)}}{2} + \frac{S_0^{(2)}}{2} - \frac{S_1^{(0)}}{2} \right]$$

In matrix form,

$$\begin{bmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -1 & 1 & 1 & 1 \\ 1 & 0 & -2 & 0 \\ -1 & 1 & 1 & -1 \end{bmatrix} \begin{bmatrix} S_0^{(0)} \\ \frac{1}{2} S_0^{(1)} \\ \frac{1}{2} S_0^{(2)} \\ \frac{1}{2} S_1^{(2)} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ -1 & 1 & 2 & 1 & -1 \\ 1 & 0 & -2 & 0 & 2 \\ -1 & 1 & 0 & -1 & -1 \end{bmatrix} \begin{bmatrix} h_0 & 0 & 0 & 0 & 0 \\ 0 & \frac{h_0+h_1}{2} & 0 & 0 & 0 \\ 0 & 0 & \frac{h_0}{2} & 0 & 0 \\ 0 & 0 & 0 & (h_1 - \frac{h_0}{2}) & 0 \\ 0 & 0 & 0 & 0 & \frac{h_0+h_1}{2} \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 1 & -1 \\ 1 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \quad (\text{final form})$$

Using Winograd algorithm,
for 2x3 convolution,

multiplication required = 5

addition required = 11.

whereas, in direct form,

mult. req. = 6

add. req. = 2.

Linear Winograd:Winograd algorithm also called Minimal Filtering algorithm: g : filter d : input y : output

$$Y = A^T [(Gg) \odot (B^T d)]$$

dot/pointwise mult.
each dimension

$$G = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & -1 & -1 \end{bmatrix}$$

$$B^T = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix}$$

1 D CNN:

$$F [2, 3]$$

\downarrow filter
 ip

$$g: [g_0 \ g_1 \ g_2 \ g_3]$$

\downarrow Rearrange

$$F = \begin{bmatrix} g_0 & g_1 & g_2 \\ g_1 & g_2 & g_3 \end{bmatrix} \begin{bmatrix} f_0 \\ f_1 \\ f_2 \end{bmatrix}$$

$$= \begin{bmatrix} f_0 g_0 + f_1 g_1 + f_2 g_2 \\ f_0 g_1 + f_1 g_2 + f_2 g_3 \end{bmatrix} \rightarrow \begin{array}{l} 6 \text{ multiplications} \\ + \\ 4 \text{ additions} \end{array}$$

$$m_1 = (g_0 - g_2) f_0$$

$$m_2 = (g_1 + g_2) \left(\frac{f_0 + f_1 + f_2}{2} \right)$$

$$m_3 = (g_2 - g_1) \left(\frac{f_0 - f_1 + f_2}{2} \right)$$

$$m_4 = (g_1 - g_3) f_2$$

$$F = \begin{bmatrix} m_1 + m_2 + m_3 \\ m_2 - m_3 - m_4 \end{bmatrix} \rightarrow \begin{array}{l} 4 \text{ multiplications (reduced)} \\ 8 \text{ additions (increased)} \\ \text{[excluding precomputed } f_i's] \end{array}$$

$$\text{2D CNN: } A^T [(G, g G^T) \odot (B^T d B)]$$

$$f(x, y) \circledast g(x, y)$$

$$= \sum_{+k} \sum_{+l} f(x-k, y-l) \cdot g(x, y)$$

↳ Same as 1D operated 2 times; column-wise and row-wise and then added.

$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} e \\ f \end{bmatrix}$: inner product

outer product : reduce the ~~total~~ power consumption and computation complexity.

$$ae + bf$$

$$ce + df$$

17-04-2025

Iterative Fast convolution Method
[in book (Parhi)].

~~Quadratic~~ Heuristic based fast convolution
(in book)

Arithmetics :

- ① Canonical Sign Digit Arithmetic
- ② Distributed Arithmetic

Distributed Arithmetic (DA)

→ ~~No direct multiplication~~ No direct multiplication.

→ For inner product and MAC calculation.

→ A multiply and Accumulate (MAC) operation:

$$y = A_1x_1 + A_2x_2 + \dots + A_Kx_K$$

$$= \sum_{K=1}^K A_Kx_K$$

↓ const. no. ↓ → Tertiary bit/
 ↓ Redundant bit

→ DA technique is bit-serial in nature. ⇒ Slow

→ When no. of elements in a vector is nearly same as the wordsize ⇒ Fast. → Replaces explicit multiplications by ROM look-ups.

Formulation of DA:

$$y = \sum_{K=1}^K A_Kx_K$$

Let x_K be an N-bit scaled 2's complement no.

$$|x_K| < 1$$

$x_K: \{b_{K0}, b_{K1}, b_{K2}, \dots, b_{K(N-1)}\} \rightarrow 2^8$ complement.

where b_{K0} is a sign bit

$$\therefore x_K = -b_{K0} + \sum_{n=1}^{N-1} b_{Kn} 2^{-n}$$

$$\Rightarrow y = \sum_{K=1}^K A_K \left[-b_{K0} + \sum_{n=1}^{N-1} b_{Kn} 2^{-n} \right]$$

$$= - \sum_{K=1}^K (b_{K0} \cdot A_K) + \sum_{K=1}^K \underbrace{\sum_{n=1}^{N-1} (A_K \cdot b_{Kn}) 2^{-n}}$$

$$= - \sum_{K=1}^K (b_{K0} \cdot A_K) + \sum_{K=1}^K \left[(A_K \cdot b_{K1}) 2^{-1} + (A_K \cdot b_{K2}) 2^{-2} + \dots + (A_K \cdot b_{K(N-1)}) 2^{-(N-1)} \right]$$

$$= - [b_{10} \cdot A_1 + b_{20} \cdot A_2 + \dots + b_{K0} \cdot A_K]$$

$$+ [(b_{11} \cdot A_1) 2^{-1} + (b_{12} \cdot A_1) 2^{-2} + \dots + (b_{1(N-1)} \cdot A_1) 2^{-(N-1)}]$$

$$+ [(b_{21} \cdot A_2) 2^{-1} + (b_{22} \cdot A_2) 2^{-2} + \dots + (b_{2(N-1)} \cdot A_2) 2^{-(N-1)}]$$

⋮

$$+ [(b_{K1} \cdot A_K) 2^{-1} + (b_{K2} \cdot A_K) 2^{-2} + \dots + (b_{K(N-1)} \cdot A_K) 2^{-(N-1)}]$$

$$\begin{aligned}
&= -[b_{10} \cdot A_1 + b_{20} \cdot A_2 + \dots + b_{K0} \cdot A_K] \\
&\quad + [(b_{11} \cdot A_1) + (b_{21} \cdot A_2) + \dots + (b_{K1} \cdot A_K)] 2^{-1} \\
&\quad + [(b_{12} \cdot A_1) + (b_{22} \cdot A_2) + \dots + (b_{K2} \cdot A_K)] 2^{-2} \\
&\quad \vdots \\
&\quad + [(b_{1(N-1)} \cdot A_1) + (b_{2(N-1)} \cdot A_2) + \dots + (b_{K(N-1)} \cdot A_K)] 2^{-(N-1)} \\
&= -\sum_{k=1}^K (b_{k0}) \cdot A_k + \sum_{n=1}^{N-1} [b_{1n} \cdot A_k + b_{2n} \cdot A_2 + \dots + b_{Kn} \cdot A_K] 2^{-n}
\end{aligned}$$

$$\therefore y = -\sum_{k=1}^K A_k \cdot (b_{k0}) + \sum_{n=1}^{N-1} \left[\sum_{k=1}^K A_k \cdot b_{kn} \right] 2^{-n}$$

$$\text{Or, } y = \sum_{n=1}^{N-1} \left[\sum_{k=1}^K A_k b_{kn} \right] 2^{-n} + \sum_{k=1}^K A_k (-b_{k0})$$

↓ ↓
 2^K possible values

With sign bit as input, we can store it in ROM of size $= 2 \times 2^K$.

Eg. No. of taps, $K=4$

Fixed coeff, $A_1=0.72, A_2=-0.3, A_3=0.95, A_4=0.11$

We need $2^K = 2^4 = 16$ -word ROM.

ROM: Address and Content:

| b_{1n} | b_{2n} | b_{3n} | b_{4n} | Content |
|----------|----------|----------|----------|--------------------------------|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | $A_4 = 0.11$ |
| 0 | 1 | 0 | 1 | $A_2 + A_4 = -0.19$ |
| 1 | 1 | 1 | 1 | $A_1 + A_2 + A_3 + A_4 = 1.48$ |

$$\sum_{k=1}^4 A_k b_{kn} = A_1 b_{1n} + A_2 b_{2n} + A_3 b_{3n} + A_4 b_{4n}$$

↳ MAC without explicit multiplier

↳ Size of ROM grows exponentially with each added I/P address line.

↳ For each element in a vector, we have an address line.

For $K=16 \Rightarrow 2^{16}$ (i.e., $6+K$) of ROM $\rightarrow 16$ add. lines.

↳ Offset binary coding to reduce ROM size

$$Y_1 = \sum_{k=1}^K A_k b_{kn}$$

For $K=3$,

$$Y_1 = A_1 b_{1n} + A_2 b_{2n} + A_3 b_{3n}$$

If bits $\equiv 0, 0, 0$: $Y_1 = 0$

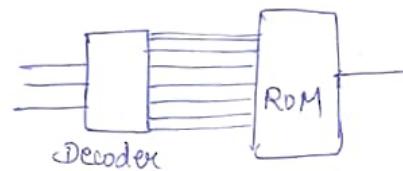
bits $\equiv 0, 0, 1$: $Y_1 = A_3$

bits $\equiv 1, 1, 1$: $Y_1 = A_1 + A_2 + A_3$.



ROM can be used to store the bits.

$$\text{Size} = 2 \times 2^K$$

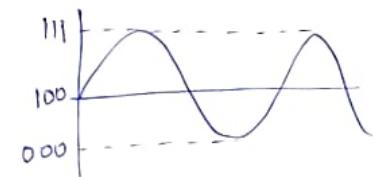
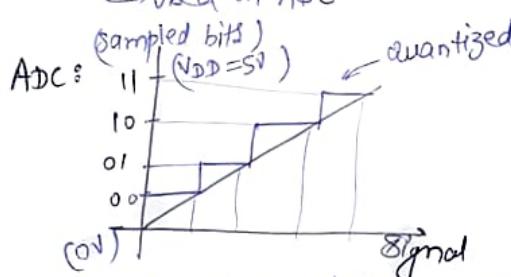


Lookup Table

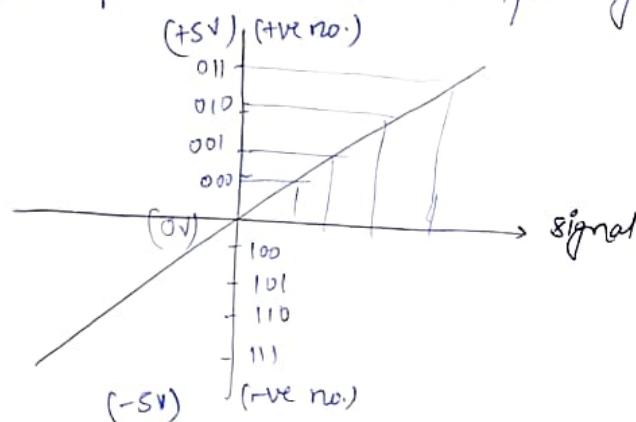
↳ Used in FPGA

Binary offset Coding Technique

↳ Used in ADC



Signed representation : 3rd bit for sign.

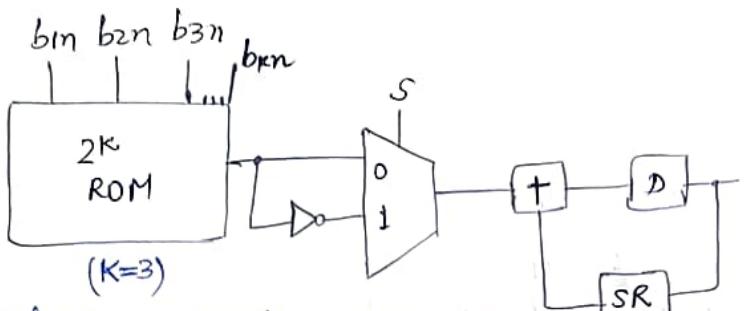


Binary offset coding
↓
Symmetric

| | |
|-----|---------|
| 111 | (+5V) |
| 110 | (VDD/2) |
| 010 | (0V) |
| 001 | (-5V) |
| 000 | (-VDD) |

→ Size of ROM : 2^2
(2-to-4 decoder)
(not 2^3)

Bit serial : DA for FIR filter



$-b_0, a + \dots = 2^8$ complement

$$\begin{pmatrix} \text{sign} \\ + \\ (\text{sum}) 2^{-1} \\ + \\ (\text{sum}) 2^{-2} \\ \vdots \end{pmatrix}$$

$$FT = \int f(x) e^{-j2\pi ft}$$

↳ Gives frequency component magnitude and phase.

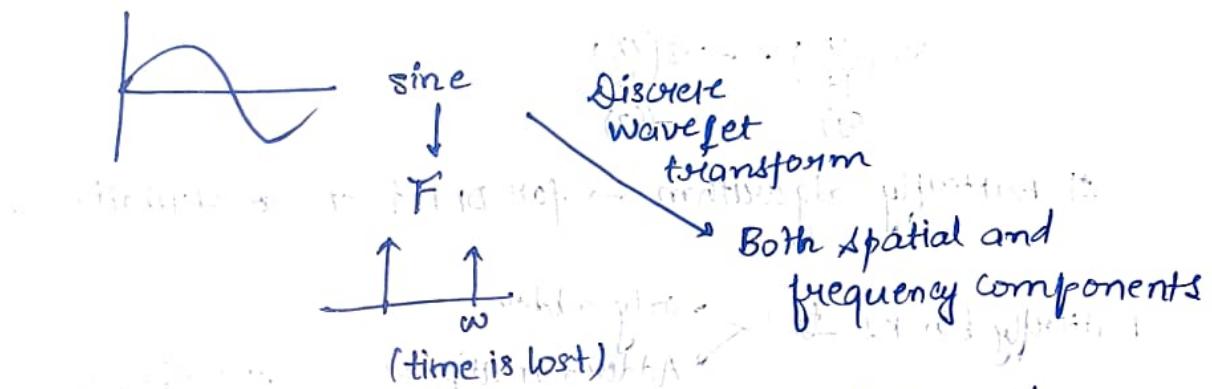
$e^{j2\pi ft}$
↳ orthogonal basis fn

Reduction in Filters and Transformations

Fourier transform: orthogonal transform

cosine transform → Real transform

↳ Discrete cosine transform (DCT)
(eg, jpeg compression)



Different transform → gives additional information
(e.g., noise)

$$\text{DCT: } X(k) = e(k) \sum_{n=0}^{N-1} x(n) \left[\cos \left[\frac{(2n+1)k}{2N} \pi \right] \right], \quad k=0, 1, \dots, N-1$$

$$\text{IDCT: } x(n) = \frac{2}{N} \sum_{k=0}^{N-1} e(k) X(k) \cos \left[\frac{(2n+1)k}{N} \pi \right], \quad n=0, 1, \dots, N-1.$$

$$e(k) = \begin{cases} 1/\sqrt{2}, & k=0 \\ 1, & \text{otherwise} \end{cases}$$

$$C_i = \cos \left(\frac{i\pi}{16} \right)$$

$$\therefore X(k) = e(k) \sum_{n=0}^7 x(n) \cos \left(\frac{(2n+1)k}{16} \pi \right), \quad k=0, 1, 2, \dots, 7.$$

$$\begin{pmatrix} x(0) \\ x(1) \\ \vdots \\ x(7) \end{pmatrix} = \begin{pmatrix} C_0 & C_1 & C_2 & C_3 & C_4 & C_5 & C_6 & C_7 \\ C_0 & C_1 & C_2 & C_3 & C_4 & C_5 & C_6 & C_7 \\ C_1 & C_2 & C_3 & C_4 & C_5 & C_6 & C_7 & C_0 \\ C_2 & C_3 & C_4 & C_5 & C_6 & C_7 & C_0 & C_1 \\ C_3 & C_4 & C_5 & C_6 & C_7 & C_0 & C_1 & C_2 \\ C_4 & C_5 & C_6 & C_7 & C_0 & C_1 & C_2 & C_3 \\ C_5 & C_6 & C_7 & C_0 & C_1 & C_2 & C_3 & C_4 \\ C_6 & C_7 & C_0 & C_1 & C_2 & C_3 & C_4 & C_5 \\ C_7 & C_0 & C_1 & C_2 & C_3 & C_4 & C_5 & C_6 \end{pmatrix} \begin{pmatrix} x(0) \\ x(1) \\ \vdots \\ x(7) \end{pmatrix}$$

↳ Reduce $8 \times 8 = 64$ coeff. to 8 coeff. (for 8-point DFT) $C_4 = \cos \left(\frac{4\pi}{16} \right)$

$$X(0) = \frac{1}{\sqrt{2}} \sum_n x_n = \frac{1}{\sqrt{2}} \cos\left(\frac{0\pi}{16}\right) \sum_n x_n = c_0 \sum_n x_n$$

Reduction: Using Trigonometric identities.

$$\cos\left(\frac{i\pi}{16} \pm 2n\pi\right) = \cos\frac{i\pi}{16} = c_i$$

$$\cos\left(\frac{i\pi}{16} \pm (k\pi + 1)\pi\right) = -\cos\frac{i\pi}{16} = -c_i$$

$$\cos\left(\frac{9\pi}{16}\right) = -\cos\left(\frac{\pi}{16}\right)$$

(c₉) -(c₁)

Butterfly algorithm → for DFT & CFT calculations.

24-04-2025

Butterfly blocks type → only adder

→ Adder + multiplier

Orthogonal transform

with unity gain filter to help in analysis which makes inverse possible

$$X(1) = M_0 C_1 + M_1 C_7 + M_2 C_3 + M_3 C_5$$

$$M_0 = x_0 - x_7$$

$$M_1 = x_3 - x_4$$

$$M_2 = x_1 - x_6$$

$$M_3 = x_2 - x_5$$

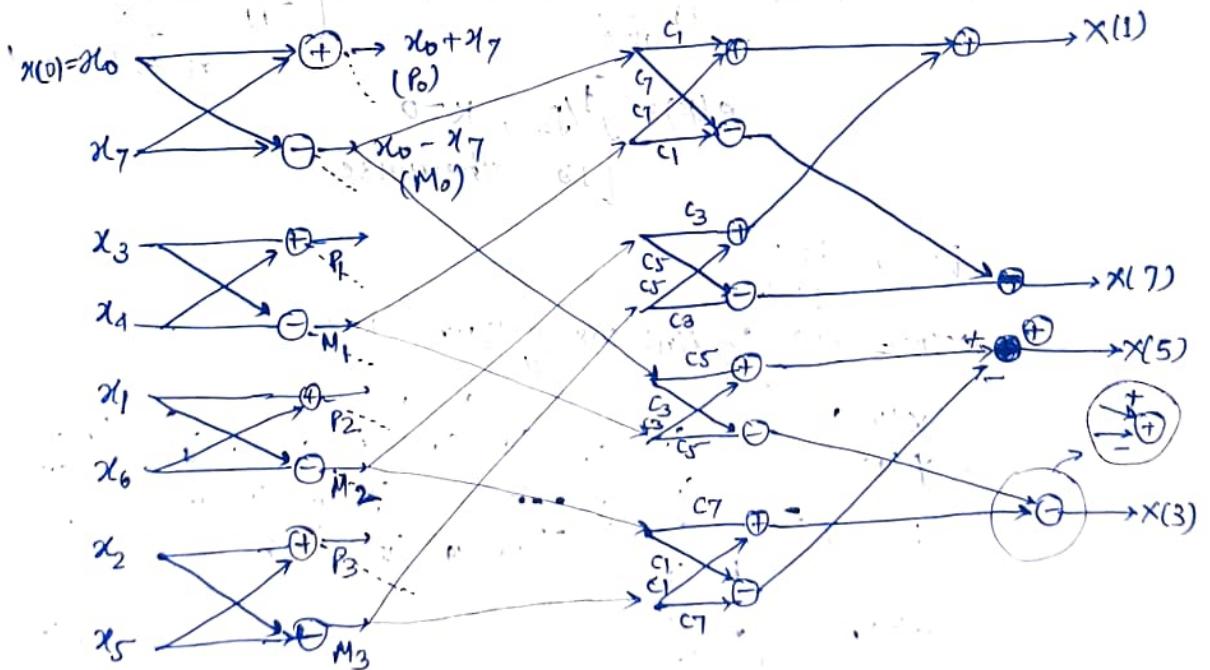
$$P_0 = x_0 + x_7$$

$$P_1 = x_3 + x_4$$

$$P_2 = x_1 + x_6$$

$$P_3 = x_2 + x_5$$

★ Decimation in time FFT algorithm
↳ for endsem



$$M_{10} = P_0 - P_1$$

$$M_{11} = P_2 - P_3$$

$$M_{100} = P_{10} - P_{11}$$

$$P_{10} = P_0 + P_1$$

$$P_{11} = P_2 + P_3$$

$$P_{100} = P_{10} + P_{11}$$

$$X(7) = M_0 C_7 - M_1 C_6 - M_2 C_5 + M_3 C_4$$

$$X(3) = M_0 C_3 - M_1 C_2 - M_2 C_1 - M_3 C_0$$

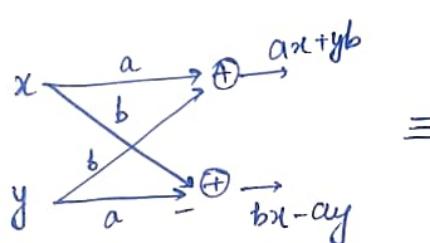
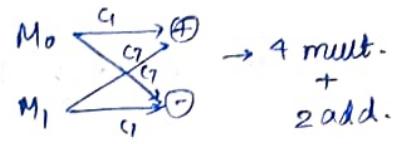
$$X(5) = M_0 C_5 + M_1 C_4 - M_2 C_3 + M_3 C_2$$

$$X(2) = M_{10} C_2 + M_{11} C_1$$

$$X(6) = M_{10} C_6 + M_{11} C_5$$

$$X(4) = M_{100} C_4$$

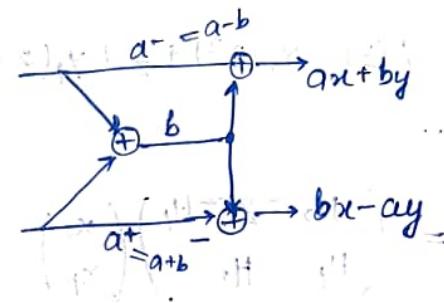
$$X(0) = P_{100} C_4$$



↳ Rotation matrix

(e.g., $\cos\theta, \sin\theta, -\cos\theta, -\sin\theta$)

↳ 4 mult. + 2 add.



[$a+, -a^- \rightarrow$ pre-computed]

$$\text{Rotation matrix: } \begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Formulation of Parallel Filters Using Polyphase Decomposition: → area ↑ computation time ↓

Polyphase structure:

$$Y(z) = H(z) X(z) = \left(\sum_n h(n) z^{-n} \right) \left(\sum_n x(n) z^{-n} \right)$$

$$\begin{aligned} X(z) &= \sum_{n=0}^{\infty} x(n) z^{-n} \\ &= x_0 + x_1 z^{-1} + x_2 z^{-2} + \dots \\ &= (x_0 + x_2 z^{-2} + x_4 z^{-4} + \dots) \\ &\quad + (x_1 z^{-1} + x_3 z^{-3} + \dots) \\ &= x_0 (z^{+2}) + z^{-1} x_1 (z^{+2}) \\ &\quad \downarrow \text{polyphase components} \end{aligned}$$

Similarly,

$$H(z) = H_0(z^2) + z^{-1} H_1(z^2)$$

$$Y(z) = X_0(z^{+2}) H_0(z^2) + z^{-1} [X_0(z^2) H_1(z^2) + X_1(z^2) H_0(z^2)] \\ + z^{-2} X_1(z^2) H_1(z^2).$$

$$Y_0(z^2) = X_0(z^2) H_0(z^2) + z^{-2} X_1(z^2) H_1(z^2)$$

$$Y_1(z^2) = [X_0(z^2) H_1(z^2) + X_1(z^2) H_0(z^2)]$$

$$\therefore Y(z) = Y_0(z^2) + z^{-1} Y_1(z^2)$$

$$Y = HX$$

$$\Rightarrow \begin{pmatrix} Y_0 \\ Y_1 \end{pmatrix} = \begin{pmatrix} H_0 & z^{-2} H_1 \\ H_1 & H_0 \end{pmatrix} \begin{pmatrix} X_0 \\ X_1 \end{pmatrix}$$

$$x(2k) \rightarrow H_0 \oplus y(2k) \equiv Y_0$$

$$x(2k+1) \rightarrow H_1 \oplus y(2k+1) \equiv Y_1$$

$$x(2k+1) \rightarrow H_0 \oplus y(2k+1) \equiv Y_1$$

Rank-order filter
→ X

Coordinate Rotation Digital Computer (CoRDIC)

↳ To find trigonometric functions, hyperbolic functions, square roots.

Taylor series:

$$\cos x = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots$$

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

DSP chip:
 ↳ TMS 6455
 ↳ Tiger shark

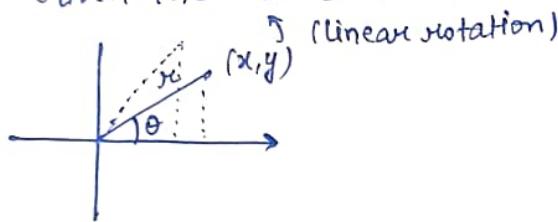
Iterative algorithm:

↳ Binary search

↳ Strength reduction.

Cordic Rotation Mode (vectoring Algorithm)

Given the coordinates, how to find θ ?



$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

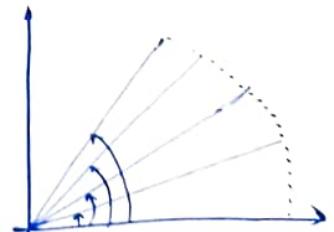
Eg. For angle, $\theta = 70^\circ$

$$\text{Rotate by } -45^\circ \Rightarrow -45^\circ + 70^\circ = 25^\circ$$

$\theta \rightarrow +ve$

$$\begin{aligned} \text{Rotate by } -22.5^\circ &\Rightarrow -22.5 + 25^\circ = 2.5^\circ \\ -11.25^\circ &\Rightarrow -11.25 + 2.5^\circ = -8.75^\circ \\ +5.625^\circ &\Rightarrow +5.625 - 8.75^\circ = -3.125^\circ \\ +2.8125^\circ &\Rightarrow 2.8125 - 3.125^\circ = -0.3125^\circ \\ +1.40625^\circ &\Rightarrow 1.40625 - 0.3125^\circ = 1.09375^\circ \end{aligned}$$

$$45^\circ + 22.5^\circ + 11.25^\circ - 5.625^\circ - 2.8125^\circ =$$



$$\tan \theta = 2^{-i}, \text{ if } \theta$$

$$\theta^i = 2^{-i}$$

$$\tan \theta^i = -2^{-i}$$

$$\hat{x}_2 = x_1 - y_1 \tan \theta = x_1 - y_1 2^{-i}$$

$$\hat{y}_2 = y_1 + x_1 \tan \theta = y_1 - x_1 2^{-i}$$

$$\text{iteration: } x^{i+1} = x^i - d_i 2^{-i} \cdot y^i \quad d = \pm 1$$

$$y^{i+1} = y^i + d_i (2^{-i} x^i)$$

$$z^{i+1} = z^i - d_i \theta(i)$$

Canonical Signed Digit Representation (CSD)

↪ Rep. no. of non-zero bits minimal.

shift add

Ternary Representation: $\{0, 1, -1\}$

Binary: $\{0, 1\}$

Eg. $x \times 30 = x \times 3 \times 10^1 + 0$ (decimal)

$$y = x \times 11110 \quad (\text{binary})$$

$$y = 2^4 x + 2^3 x + 2^2 x + 2^1 x$$

↪ 4 mult. + 3 add.

↪ actually shifting $\Rightarrow \begin{array}{r} 0010 \times 2 \\ \downarrow \\ 0100 \end{array}$

$$c = 32 - 2$$

$$= 100000 - 10$$

$$y = cx$$

$$y = x \times 2^5 - x \times 2^1$$

↪ 2 mult. + 1 sub. (or add.)

[2^8 complement]

Property:

① No 2 consecutive bits in CSD are non-zero.

② Minimal no. of non-zero bits.

③ It is unique.

④ Range of no.: $\left\{-\frac{4}{3}, \frac{4}{3}\right\}$

⑤ w -bit CSD $\{-1, 1\}$,

the avg. no. of non-zero bits = $\frac{w}{3} + \frac{1}{9} + O(2^{-w})$

15-bit no. $\rightarrow \frac{1}{3}$ ($\sim 30^\circ$) $\xrightarrow{\text{20° Reduction}}$ $\frac{1}{2}$

↑ on avg.
on avg.

possibility
(probability in)
(CSD)

of giving non-
zero bits

probability
(probability in binary)

Ex. Binary \rightarrow 1010 ~~111~~
 not allowed (consecutive non-zero)
 LSB

0111



100~~1~~ \rightarrow (-1)

$\therefore 10\boxed{1}00\bar{1}$

not allowed

↓
110100~~1~~

↓
 $\therefore 10\overline{1}0\overline{1}00\overline{1}$
add subtract

$$2^7 - 2^5 - 2^3 - 2^0$$

Binary CSD converter

Represent negative by 2's complement.

↳ MSB = sign bit

Ex. Convert 11110 to CSD.

↓

011110

↓
100010

↓
1000~~1~~0 \rightarrow 30

$$2^5 - 2^1 = 30$$

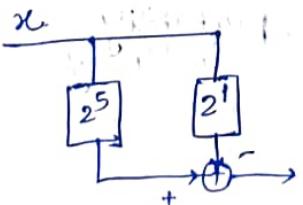
1000~~1~~0 \rightarrow -30

→ 011110

↓
1000~~1~~0 \rightarrow 30

1000~~1~~0 (if leading 0 is sign bit)

$$-30^{\circ}$$



Finite wordlength:

Hartree's rule for precision system

$$x \cdot 2^{-5} + x^{-3}$$

\rightsquigarrow right-shift by 5.

$x \rightarrow 8 \text{ bit} : 0111\ 1111$

$$0111\ 1111 \quad \boxed{0\ 0\ 0\ 0\ 0\ 1\ 1\ 1111} \quad (2^{-5})$$

$$+ \quad 0\ 0\ 0\ 0\ 1\ 1\ 1\ 111 \quad (2^{-3})$$

$$(x(2^{-2}) + x)2^{-3}$$

$$x \times (0.10100\bar{1}0010\bar{1}00\bar{1})$$

$$= x \times (2^{-1} + 2^{-3} - 2^{-6} + 2^{-9} - 2^{-11} - 2^{-14})$$

$$= 2^{-1} [(x + 2^{-2}x)] + 2^{-5} [(-x + 2^{-3}x) \\ - 2^{-5}(x + 2^{-3})] \quad \cancel{- 2^{-7}x - 2^{-10}x - 2^{-13}x}$$

→ Precision gets improved.

Sub-expression elimination

