

Structures For Discrete-Time Systems

- Realization (implementation) of digital filters
- Structures of IIR and FIR filters, their advantages and disadvantages – efficiency and error

Example: Given $H(z) = \frac{b_0 + b_1 z^{-1}}{1 - a z^{-1}} \quad |z| > |a|$

$\Rightarrow h[n] = b_0 a^n u[n] + b_1 a^{n-1} u[n-1] \quad \text{IIR}$

It is not possible to implement the system by discrete convolution!

$\Rightarrow y[n] = a y[n-1] + b_0 x[n] + b_1 x[n-1]$

Actually, an unlimited variety of computational structures result in the same relation between $y[n]$ and $x[n]$!

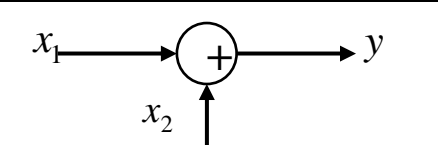
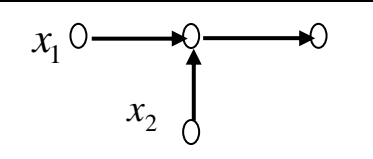
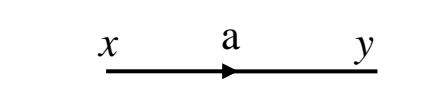
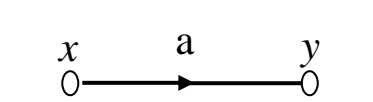
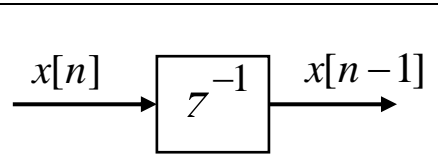
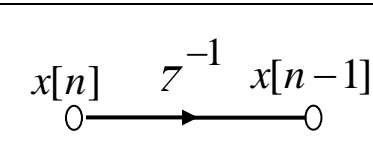
When the numerical precision is limited

\Rightarrow different structures may have vastly different behavior.

- (i) Finite-precision representation of the system coefficients
- (ii) Truncation or rounding of intermediate computations.

✧ Block Diagram and Signal Flow Graph

- Three elements in LTI discrete-time systems:

	Block diagram	Signal flow graph
Adder		
Scalar (Multiplication by a constant)		
Unit delay		

- Nodes and branches are keys in a signal flow graph
 - Source node:** No entering branches
 - Sink node:** Only entering branches

✧ Basic Structures for IIR Systems

- Direct Forms

- (1) Direct Form I

$$H(z) = \frac{\sum_{k=0}^M b_k z^{-k}}{1 - \sum_{k=1}^N a_k z^{-k}} \Leftrightarrow y[n] - \sum_{k=1}^N a_k y[n-k] = \sum_{k=0}^M b_k x[n-k]$$

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

$$Y(z) = H_2(z)V(z) \Leftrightarrow y[n] - \sum_{k=1}^N a_k y[n-k] = v[n]$$

$$X(z) = V(z)/H_1(z)$$

$$V(z) = H_1(z)X(z) \Leftrightarrow v[n] = \sum_{k=0}^M b_k x[n-k]$$

$$V(z) = \sum_{k=0}^M b_k X(z)z^{-k}$$

$$Y(z) - \sum_{k=1}^N a_k Y(z)z^{-k} = V(z)$$

$$Y(z) = V(z) + \sum_{k=1}^N a_k Y(z)z^{-k}$$

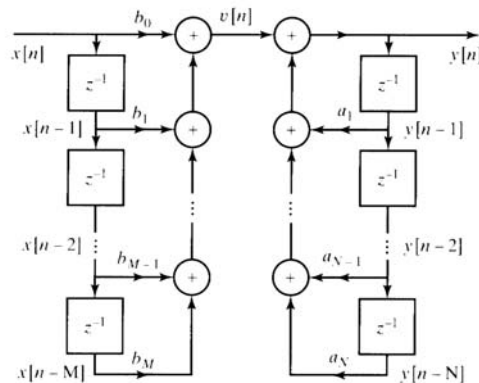


Figure 6.3 Block diagram representation for a general N th-order difference equation.

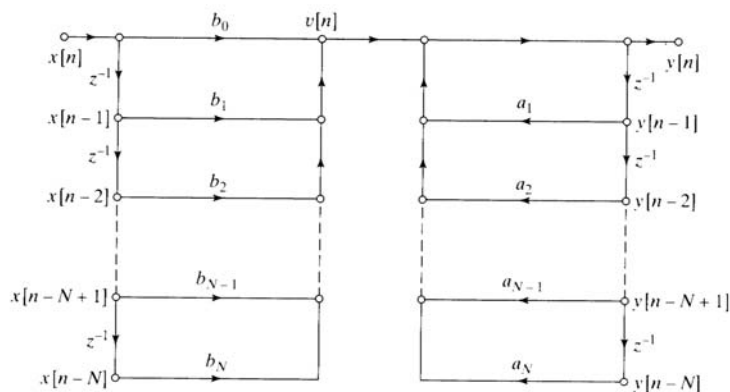


Figure 6.14 Signal flow graph of direct form I structure for an N th-order system.

(2) Direct Form II (Canonic form)

-- Interchange 1st and 2nd “segments” and merge the delay lines (z^{-1})

-- Number of delay = $\max(N, M) \leftarrow$ “Canonic”

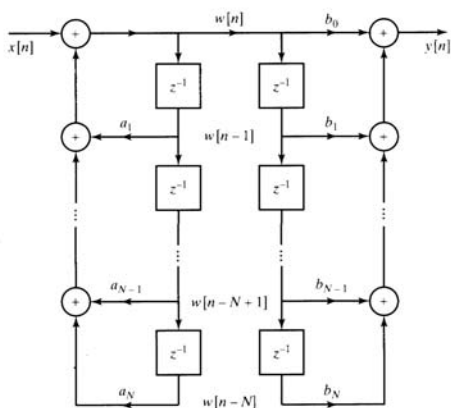


Figure 6.4 Rearrangement of block diagram of Figure 6.3. We assume for convenience that $N = M$. If $N \neq M$, some of the coefficients will be zero.

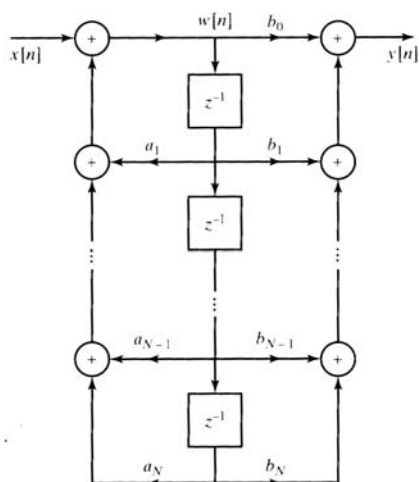


Figure 6.5 Combination of delays in Figure 6.4.

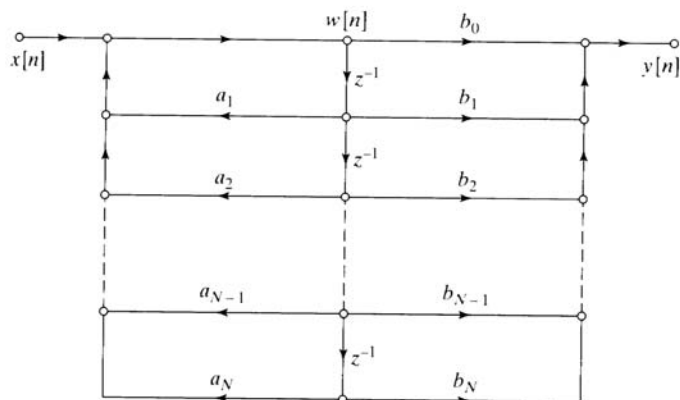


Figure 6.15 Signal flow graph of direct form II structure for an N th-order system.

• **Cascade Form**

-- Serial connection of 1st order and 2nd order factors

$$H(z) = \prod_{k=1}^{N_s} \frac{b_{0k} + b_{1k}z^{-1} + b_{2k}z^{-2}}{1 - a_{1k}z^{-1} - a_{2k}z^{-2}}$$

Remark: Each factor is a Direct Form II.

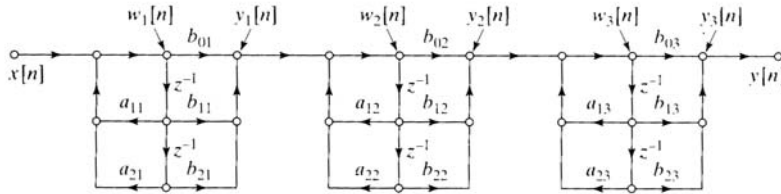


Figure 6.18 Cascade structure for a sixth-order system with a direct form II realization of each second-order subsystem.

If there are N_s second-order sections

$\Rightarrow (N_s!)^2$ different pairings and orderings!

- $\prod_{k=1}^{N_s} \frac{b_{0k} + b_{1k}z^{-1} + b_{2k}z^{-2}}{1 - a_{1k}z^{-1} - a_{2k}z^{-2}}$ needs 5 constant multipliers for each section.
- $b_0 \prod_{k=1}^{N_s} \frac{1 + \tilde{b}_{1k}z^{-1} + \tilde{b}_{2k}z^{-2}}{1 - a_{1k}z^{-1} - a_{2k}z^{-2}}$ needs 4 constant multipliers for each section.

The 5-multiplier sections are commonly used when implemented with fix-point arithmetic.

• **Parallel Form**

-- Parallel connection of 1st order and 2nd order factors

$$H(z) = \sum_{k=0}^{N_p} C_k z^{-k} + \sum_{k=1}^{N_s} \frac{e_{0k} + e_{1k}z^{-1}}{1 - a_{1k}z^{-1} - a_{2k}z^{-2}}$$

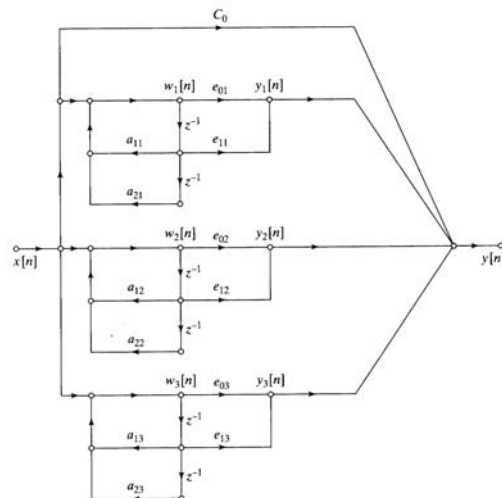


Figure 6.20 Parallel-form structure for sixth-order system ($M = N = 6$) with the real and complex poles grouped in pairs.

• Feedback in IIR Systems

-- Basic formula of a feedback system (negative feedback)

$$H(z) = \frac{F(z)}{1 + F(z)B(z)}$$

-- If a system has poles, a corresponding block diagram or signal flow graph will have feedback loops.

(BUT neither poles in the system function nor loops in the network are sufficient for the impulse response to be infinitely long.)

-- A delay element is necessary in the feedback loop; otherwise, it is *noncomputable*.
(The structure should be modified to eliminate the noncomputable loops.)

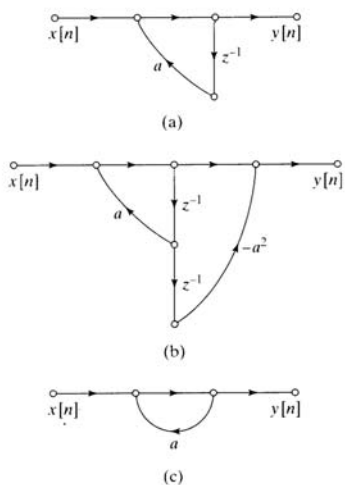


Figure 6.23 (a) System with feedback loop. (b) FIR system with feedback loop. (c) Noncomputable system.

• Transpose Forms

-- *Transposition* of a flow graph is reversing the *directions* of *all* branches in the network while keeping the branch transmittances (as they were) and reversing the roles of the input and output (so that source nodes become sink nodes and vice versa).

-- **Flow Graph Reversal Theorem**

For single-input, single-output systems, the transposed flow graph has the same system function as the original graph if the input nodes and output nodes are interchanged.

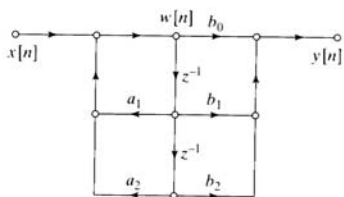


Figure 6.27 Direct form II structure for Example 6.8.

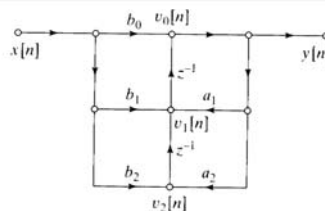


Figure 6.28 Transposed direct form II structure for Example 6.8.

✧ Basic Structures for FIR Systems

- **Direct Form**

-- Transversal filter or tapped delay line

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

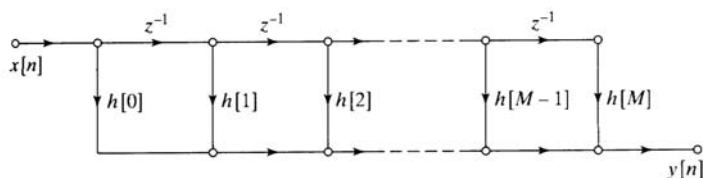


Figure 6.31 Direct-form realization of an FIR system.

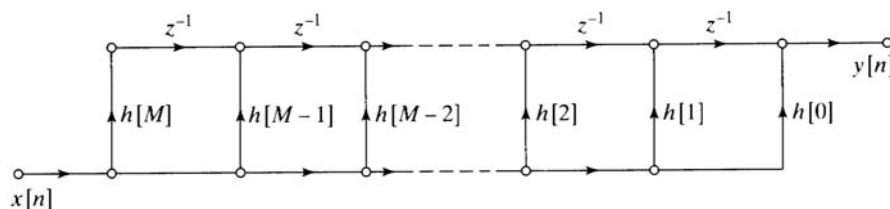


Figure 6.32 Transposition of the network of Figure 6.31.

- **Cascade Form**

-- Serial connection of 1st order and 2nd order factors

$$H(z) = \prod_{k=1}^{M_s} (b_{0k} + b_{1k}z^{-1} + b_{2k}z^{-2})$$

Remark: Each factor is a Direct Form.

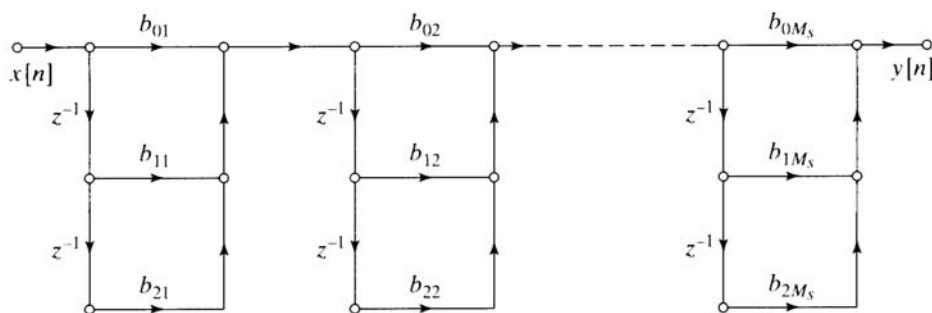


Figure 6.33 Cascade-form realization of an FIR system.

• **Linear Phase FIR Filters**

- Take the advantage of the symmetry property of the impulse response

$$h[M - n] = h[n]$$

$$h[M - n] = -h[n]$$

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

- **Type I or III: M even (order odd)**

$$y[n] = \sum_{k=0}^{M/2-1} h[k]x[n - k] + h\left[\frac{M}{2}\right]x\left[n - \frac{M}{2}\right] + \sum_{k=\frac{M}{2}+1}^M h[k]x[n - k]$$

$$= \sum_{k=0}^{M/2-1} h[k](x[n - k] \pm x[n - M + k]) + k\left[\frac{M}{2}\right]x\left[n - \frac{M}{2}\right]$$

- **Type II or IV: M odd (order even)**

$$y[n] = \sum_{k=0}^{M-1} h[k](x[n - k] \pm x[n - M + k])$$

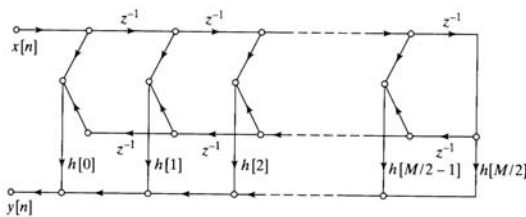


Figure 6.34 Direct-form structure for an FIR linear-phase system when M is an even integer.

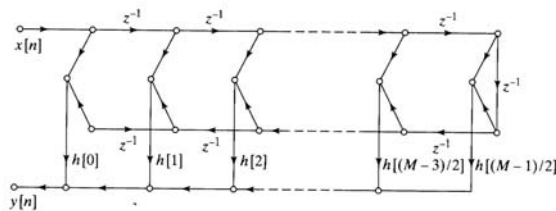


Figure 6.35 Direct-form structure for an FIR linear-phase system when M is an odd integer.

- Linear-phase FIR filters can also be implemented as a cascade of 1st-order, 2nd-order, and 4th-order real-coefficient systems. (The 4th-order system is formed by grouping the conjugate and the conjugate reciprocal zeros together.)

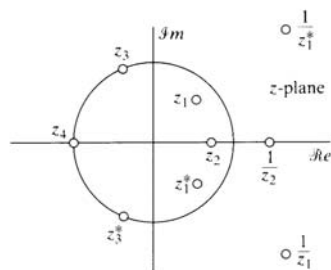


Figure 6.36 Symmetry of zeros for a linear-phase FIR filter.

✧ Finite-precision Numerical Effects

- Due to finite-precision (finite-word length) of computational and/or storage devices.
 - Parameter quantization
 - Round-off error
 - Limit cycle (IIR) ← zero input!

- Number Representation -- **Two's complement** number representation

$$\hat{x} = X_m \left(-b_0 + \underbrace{\sum_{i=1}^B b_i 2^{-i}}_{\hat{x}_B = b_0 b_1 b_2 \dots b_B} \right)$$

$$(-1 \leq \hat{x}_B < 1)$$

$$\therefore -X_m \leq \hat{x} < X_m$$

$$\Delta = X_m 2^{-B}$$

$\hat{x} = Q_B[x]$: quantized value of x

\hat{x}_B : normalized quantized value of x ; normalized value of \hat{x}

Δ : quantization stepsize

b_0 : sign bit; $b_0=0$, if x is non-negative; $b_0=1$, if x is negative.

- Quantization error $e = Q_B[x] - x$

- (1) Overflow: if $x > X_m$. This can be a serious problem if, for example, 0111 → 1000,

and we don't check it first. (This is *natural overflow*.) We first clip the input. It becomes *saturation*.

- (2) -- Rounding: nearest integer $-\Delta/2 < e \leq \Delta/2$

-- Truncation: smaller integer $-\Delta < e \leq 0$

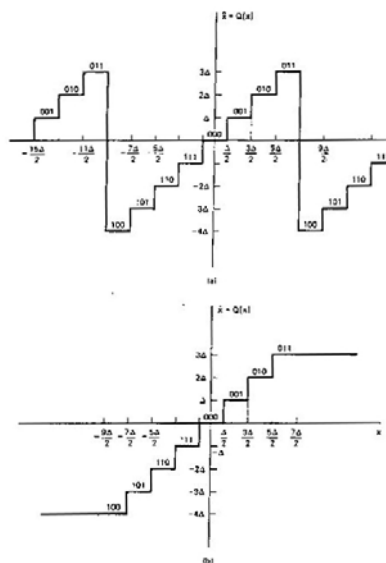


Figure 6.44 Two's-complement rounding: (a) Natural overflow (b) Saturation.

• Quantization in implementing systems

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

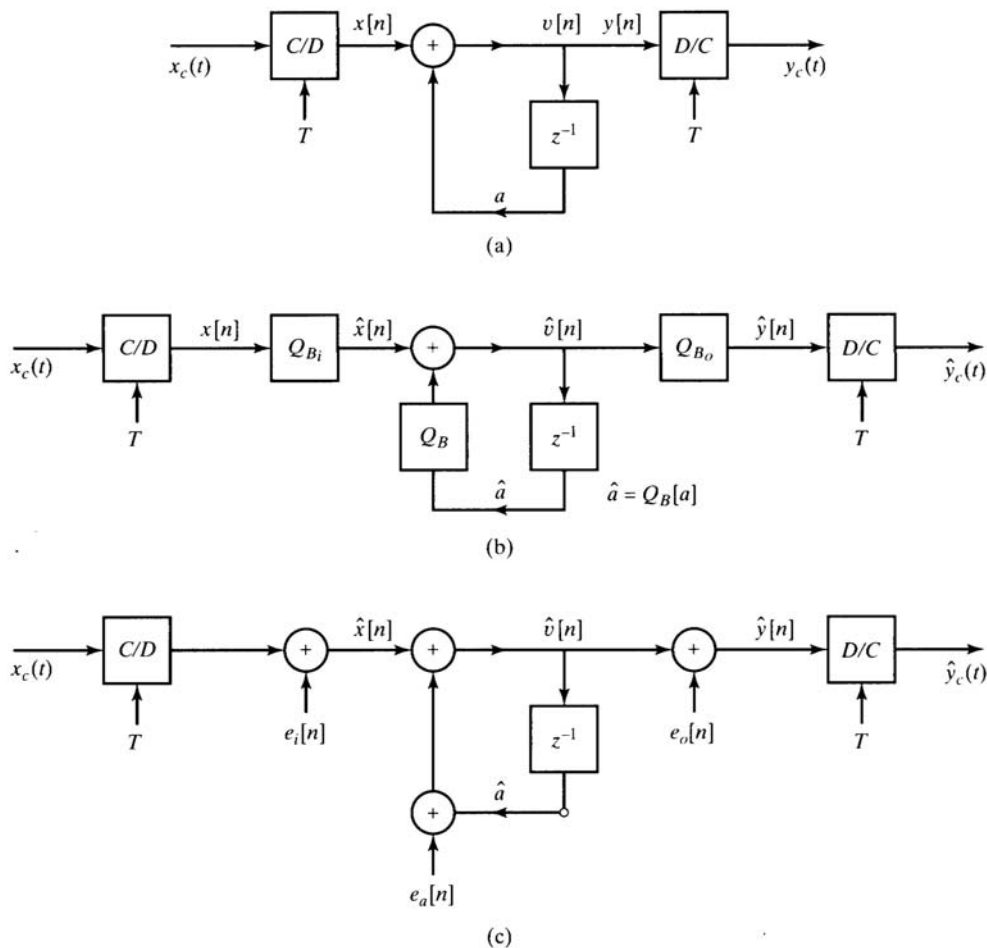


Figure 6.39 Implementation of discrete-time filtering of an analog signal. (a) Ideal system. (b) Nonlinear model. (c) Linearized model.

✧ Effects of Coefficient Quantization

- Coefficient Quantization in IIR Systems

-- depends on the filter structure

- Direct form

$$H(z) = \frac{\sum_{k=0}^M b_k z^{-k}}{1 - \underbrace{\sum_{k=1}^N a_k z^{-k}}_{A(z)}} \rightarrow \hat{H}(z) = \frac{\sum_{k=0}^M \hat{b}_k z^{-k}}{1 - \sum_{k=1}^N \hat{a}_k z^{-k}}$$

Note: $\hat{a}_k = a_k + \Delta a_k$; $\hat{b}_k = b_k + \Delta b_k$

Effect on pole locations

(→ affect frequency response and stability)

$$\text{Compare } \begin{cases} A(z) = 1 - \sum_{k=1}^N a_k z^{-k} = \prod_{j=1}^N (1 - z_j z^{-1}) \\ \hat{A}(z) = 1 - \sum_{k=1}^N \hat{a}_k z^{-k} = \prod_{j=1}^N (1 - \hat{z}_j z^{-1}) \end{cases}$$

The change of pole location: $\hat{z}_j = z_j + \Delta z_j, \quad j = 1, \dots, N$

Δz_k is affected by all $\{\Delta a_k\}$.

$$\Delta z_i \approx \sum_{k=1}^N \left(\frac{\partial z_i}{\partial a_k} \right) \Delta a_k, \quad i = 1, 2, \dots, N$$

Remark: This formula is approximately true when Δa_k and Δz_k are small.

Note that if $\frac{\partial z_i}{\partial a_k}$ is large, then a small Δa_k leads to a large Δz_k . If so, this is a *sensitive* system. (Undesirable)

$$\text{One step further, } \frac{\partial z_i}{\partial a_k} = \frac{z_i^{N-k}}{\prod_{j=1, j \neq i}^N (z_i - z_j)} \cdot \left(\begin{array}{ccc} \left(\frac{\partial A(z)}{\partial z_i} \right)_{z=z_i} & \frac{\partial z_i}{\partial a_k} = & \left. \frac{\partial A(z)}{\partial a_k} \right|_{z=z_i} \\ \Downarrow & & \Downarrow \\ \prod_{i \neq j} (1 - z_i^{-1} \cdot z_j) & & z_i^{-k} \end{array} \right)$$

That is, if $(z_i - z_j)$ is small, then $\frac{\partial z_i}{\partial a_k}$ is large; for example, narrow bandwidth

lowpass and bandpass filters which have clustered poles.

Remark: The preceding analysis can be applied to zeros.

■ **Parallel and cascade forms**

-- consists of 1st-order and 2nd-order sections.

Errors in a particular pole pairs (section) are independent of the other poles (sections).

This is also true for zeros in cascade form. → In general, both the *cascade form* and the *parallel form* are less sensitive to coefficient quantization (because zeros are often widely distributed the unit circle).

Example: Bandpass IIR elliptic filter

$$0.99 \leq |H(e^{j\omega})| \leq 1.01 \quad 0.3\pi \leq \omega \leq 0.4\pi$$

$$|H(e^{j\omega})| \leq 0.01 \quad (-40\text{dB}) \quad \omega \leq 0.29\pi$$

$$|H(e^{j\omega})| \leq 0.01 \quad (-40\text{dB}) \quad 0.41\pi \leq \omega \leq \pi$$

TABLE 6.1 UNQUANTIZED CASCADE FORM COEFFICIENTS FOR A 12TH-ORDER ELLIPTIC FILTER

k	a _{1k}	a _{2k}	b _{0k}	b _{1k}	b _{2k}
1	0.738409	-0.850835	0.135843	0.026265	0.135843
2	0.960374	-0.860000	0.278901	-0.444500	0.278901
3	0.629449	-0.931460	0.535773	-0.249249	0.535773
4	1.116458	-0.940429	0.697447	-0.891543	0.697447
5	0.605182	-0.983693	0.773093	-0.425920	0.773093
6	1.173078	-0.986166	0.917937	-1.122226	0.917937

TABLE 6.2 SIXTEEN-BIT QUANTIZED CASCADE-FORM COEFFICIENTS FOR A 12TH-ORDER ELLIPTIC FILTER

k	a _{1k}	a _{2k}	b _{0k}	b _{1k}	b _{2k}
1	24196 × 2 ⁻¹⁵	-27880 × 2 ⁻¹⁵	17805 × 2 ⁻¹⁷	3443 × 2 ⁻¹⁷	17805 × 2 ⁻¹⁷
2	31470 × 2 ⁻¹⁵	-28180 × 2 ⁻¹⁵	18278 × 2 ⁻¹⁶	-29131 × 2 ⁻¹⁶	18278 × 2 ⁻¹⁶
3	20626 × 2 ⁻¹⁵	-30522 × 2 ⁻¹⁵	17556 × 2 ⁻¹⁵	-8167 × 2 ⁻¹⁵	17556 × 2 ⁻¹⁵
4	18292 × 2 ⁻¹⁴	-30816 × 2 ⁻¹⁵	22854 × 2 ⁻¹⁵	-29214 × 2 ⁻¹⁵	22854 × 2 ⁻¹⁵
5	19831 × 2 ⁻¹⁵	-32234 × 2 ⁻¹⁵	25333 × 2 ⁻¹⁵	-13957 × 2 ⁻¹⁵	25333 × 2 ⁻¹⁵
6	19220 × 2 ⁻¹⁴	-32315 × 2 ⁻¹⁵	15039 × 2 ⁻¹⁴	-18387 × 2 ⁻¹⁴	15039 × 2 ⁻¹⁴

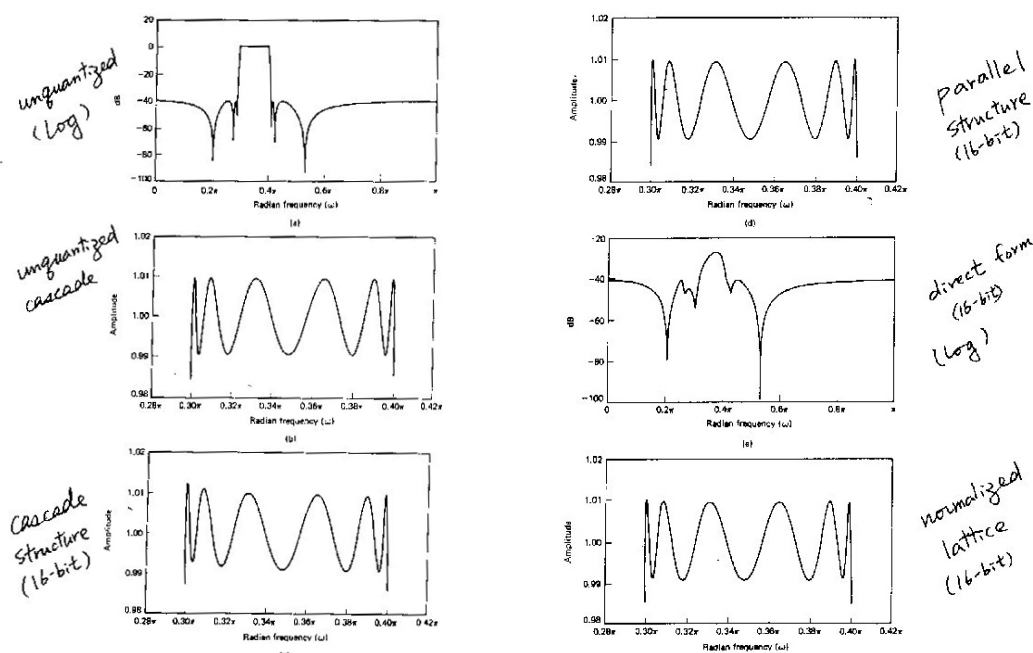


Figure 6.47 IIR coefficient quantization example. (a) Log magnitude for unquantized cascade bandpass filter. (b) Passband for unquantized cascade case. (c) Passband for cascade structure with 16-bit coefficients.

Figure 6.47 (continued) (d) Passband for parallel structure with 16-bit coefficients. (e) Log magnitude for direct form with 16-bit coefficients. (f) Passband for normalized lattice with 16-bit coefficients.

• **Second-order Section**

(1) Direct form

Poles: $(re^{j\theta}, re^{-j\theta})$

$$\frac{1}{(1 - re^{j\theta}z^{-1})(1 - re^{-j\theta}z^{-1})} = \frac{1}{1 - 2r\cos\theta z^{-1} + r^2z^{-2}}$$

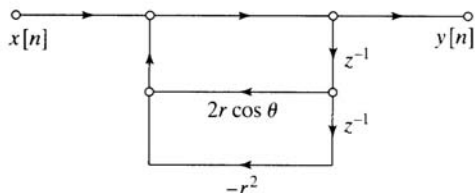


Figure 6.41 Direct-form implementation of a complex-conjugate pole pair.

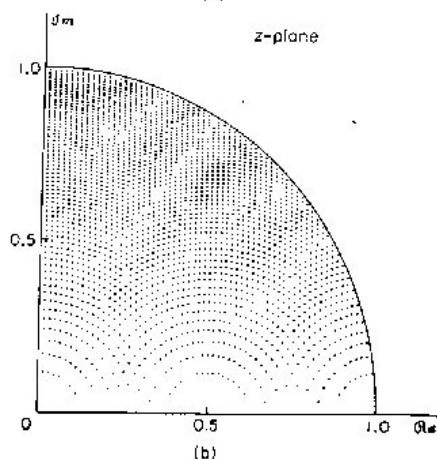
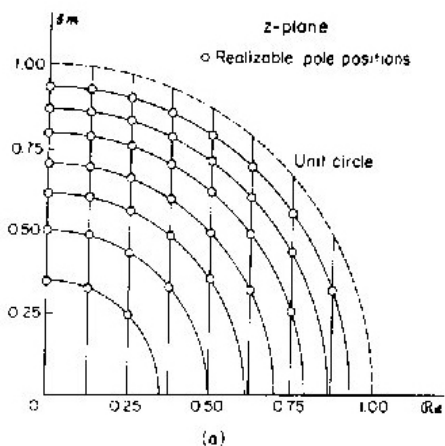


Figure 6.49 Pole locations for the second-order IIR direct form system of Fig. 6.48. (a) 4-bit quantization of coefficients. (b) 7-bit quantization.

(2) Coupled form

Poles: $(re^{j\theta}, re^{-j\theta})$

$$\begin{cases} Y_1 = X + r \cos \theta z^{-1} Y_1 - r \sin \theta z^{-1} Y \\ Y = r \sin \theta z^{-1} Y_1 + r \cos \theta z^{-1} Y \end{cases}$$

$$\Rightarrow Y_1 = (z - r \sin \theta z^{-1} Y) / (1 - r \cos \theta z^{-1})$$

$$\frac{Y}{X} = \frac{(r \sin \theta) z^{-1}}{1 - 2r \cos \theta z^{-1} + r^2 z^{-2}}$$

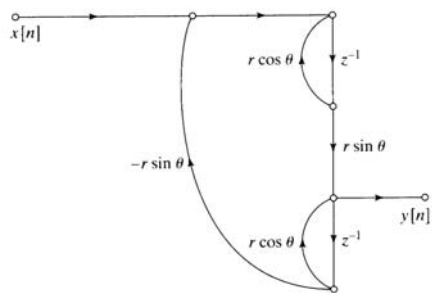


Figure 6.43 Coupled-form implementation of a complex-conjugate pole pair.

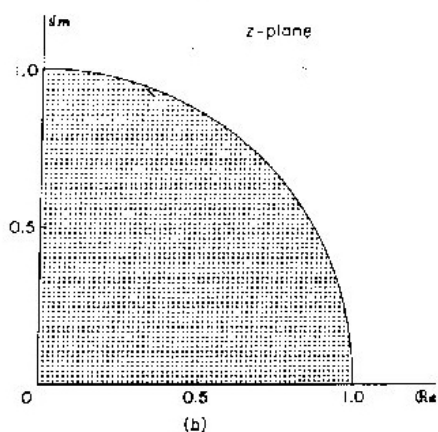
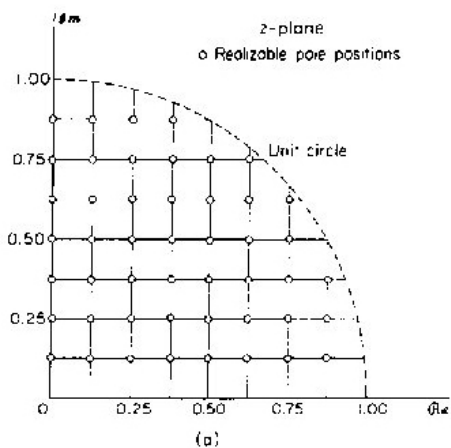


Figure 6.51 Pole locations for coupled form second-order IIR system of Fig. 6.50. (a) 4-bit quantization of coefficients. (b) 7-bit quantization.

- ➔ The pole location distribution is even.
- ⬅ The price: The number of multiplications is doubled!

• **Coefficient Quantization in FIR Systems**

■ **Direct form**

$$H(z) = \sum_{n=0}^M h[n]z^{-n} \rightarrow \hat{H}(z) = \sum \hat{h}[n]z^{-n}$$

$$= H(z) + \underbrace{\Delta H(z)}_{\sum \Delta h[n]z^{-n}}$$

$$\Delta H(z) = \sum_{n=0}^M \Delta h[n]z^{-n}$$

Effect on the zero locations

The sensitivity function of this form is similar to that of the direct form I IIR filter. That is, if the zeros are tightly clustered, their locations will then be highly sensitive to quantization errors. However, for most linear phase FIR systems, the zeros are more or less uniformly spread in the z-plane.

Effect on $H(e^{j\omega})$

After scaling, each $h[n]$ is represented by $(B+1)$ bits 2's complement number; i.e.,

$$-2^{-(B+1)} < \Delta h[n] \leq 2^{-(B+1)}.$$

$$\Delta H(e^{j\omega}) = \sum_{n=0}^M \Delta h[n]e^{-j\omega n}$$

$$|\Delta H(e^{j\omega})| = \left| \sum_{n=0}^M \Delta h[n]e^{-j\omega n} \right| \leq \sum_{n=0}^M |\Delta h[n]|e^{-j\omega n}$$

$$\leq \underbrace{(M+1)2^{-(B+1)}}_{\text{worst case!}}$$

Effect on linear phase

Not affect the linear phase property as long as $\hat{h}[n] = \hat{h}[M - n]$.

Example: Linear Phase Lowpass Filter

$$0.99 \leq |H(e^{j\omega})| \leq 1.01 \quad 0 \leq \omega \leq 0.4\pi$$

$$|H(e^{j\omega})| \leq 0.001 \quad (-60\text{dB}) \quad 0.6\pi \leq \omega \leq \pi$$

TABLE 6.3 UNQUANTIZED AND QUANTIZED COEFFICIENTS FOR AN OPTIMUM FIR LOWPASS FILTER ($M = 27$)

Coefficient	Unquantized	16 bits	14 bits	13 bits	8 bits
$h[0] = h[27]$	1.359657×10^{-3}	45×2^{-15}	11×2^{-13}	6×2^{-12}	0×2^{-7}
$h[1] = h[26]$	-1.616993×10^{-3}	-53×2^{-15}	-13×2^{-13}	-7×2^{-12}	0×2^{-7}
$h[2] = h[25]$	-7.738032×10^{-3}	-254×2^{-15}	-63×2^{-13}	-32×2^{-12}	-1×2^{-7}
$h[3] = h[24]$	-2.686841×10^{-3}	-88×2^{-15}	-22×2^{-13}	-11×2^{-12}	0×2^{-7}
$h[4] = h[23]$	1.255246×10^{-2}	411×2^{-15}	103×2^{-13}	51×2^{-12}	2×2^{-7}
$h[5] = h[22]$	6.591530×10^{-3}	216×2^{-15}	54×2^{-13}	27×2^{-12}	1×2^{-7}
$h[6] = h[21]$	-2.217952×10^{-2}	-727×2^{-15}	-182×2^{-13}	-91×2^{-12}	-3×2^{-7}
$h[7] = h[20]$	-1.524663×10^{-2}	-500×2^{-15}	-125×2^{-13}	-62×2^{-12}	-2×2^{-7}
$h[8] = h[19]$	3.720668×10^{-2}	1219×2^{-15}	305×2^{-13}	152×2^{-12}	5×2^{-7}
$h[9] = h[18]$	3.233332×10^{-2}	1059×2^{-15}	265×2^{-13}	132×2^{-12}	4×2^{-7}
$h[10] = h[17]$	-6.537057×10^{-2}	-2142×2^{-15}	-536×2^{-13}	-268×2^{-12}	-8×2^{-7}
$h[11] = h[16]$	-7.528754×10^{-2}	-2467×2^{-15}	-617×2^{-13}	-308×2^{-12}	-10×2^{-7}
$h[12] = h[15]$	1.560970×10^{-1}	5115×2^{-15}	1279×2^{-13}	639×2^{-12}	20×2^{-7}
$h[13] = h[14]$	4.394094×10^{-1}	14399×2^{-15}	3600×2^{-13}	1800×2^{-12}	56×2^{-7}

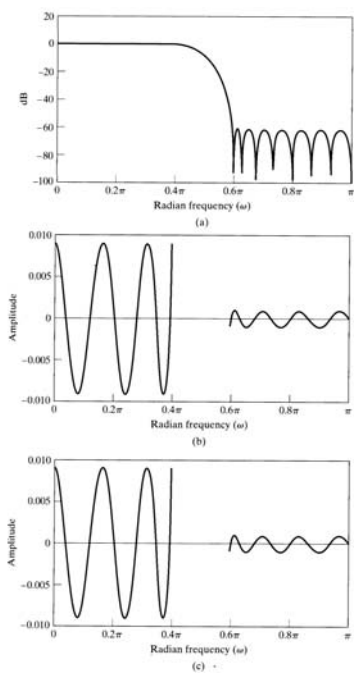


Figure 6.46 FIR quantization example. (a) Log magnitude for unquantized case. (b) Approximation error for unquantized case. (Error not defined in transition band.) (c) Approximation error for 16-bit quantization.

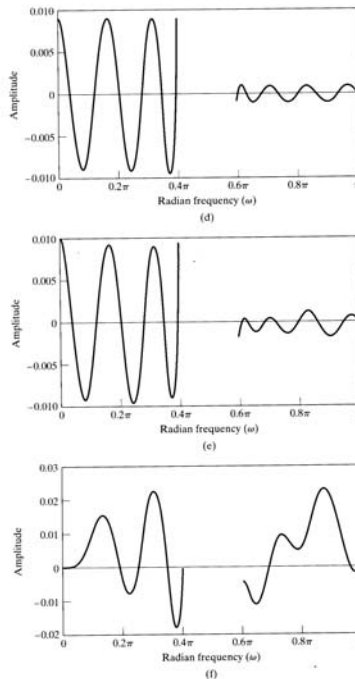


Figure 6.46 (continued) (d) Approximation error for 14-bit quantization. (e) Approximation error for 13-bit quantization. (f) Approximation error for 8-bit quantization.

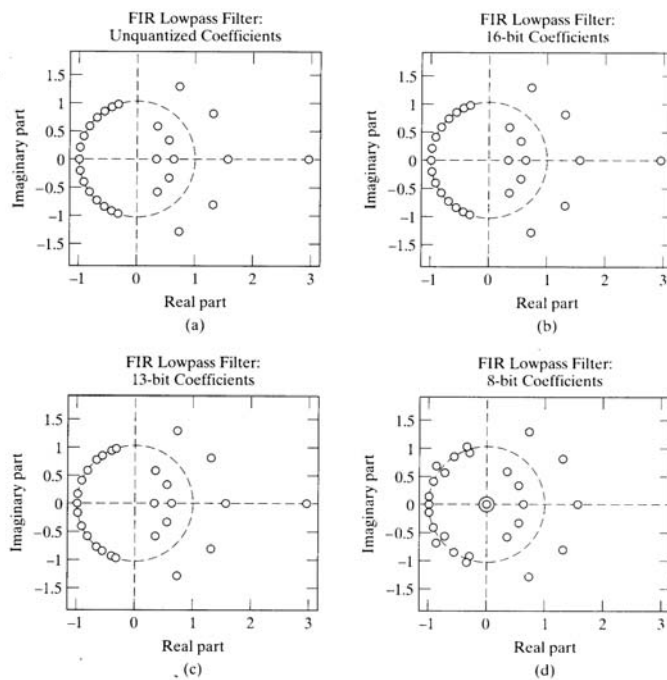


Figure 6.47 Effect of impulse response quantization on zeros of $H(z)$. (a) Unquantized. (b) Sixteen-bit quantization. (c) Thirteen-bit quantization. (d) Eight-bit quantization.

■ **Cascade form**

- less sensitive because it isolate the quantization errors from the other sections.
- To preserve linear phase each section is linear phase.
 - (a) Conjugate 2nd-order sections for conjugate zero pairs on the unit circle. $(1 + az^{-1} + z^{-2})$
 - (b) Real zero 2nd-order sections for a real zero inside the unit circle and its reciprocal (outside the unit circle).
 - (c) Zeros at ± 1 .
 - (d) 4th-order sections for conjugate zero pairs inside the unit circle and their associated reciprocals (outside the unit circle).

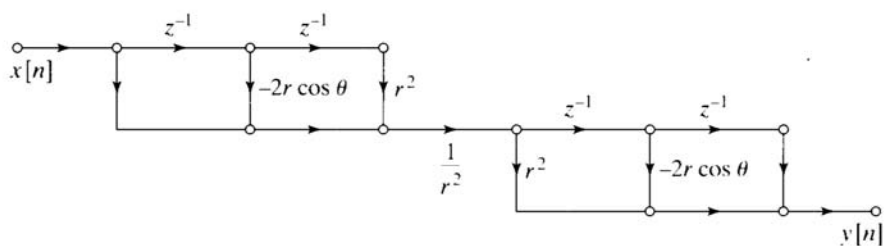


Figure 6.48 Subnetwork to implement fourth-order factors in a linear-phase FIR system such that linearity of the phase is maintained independently of parameter quantization.