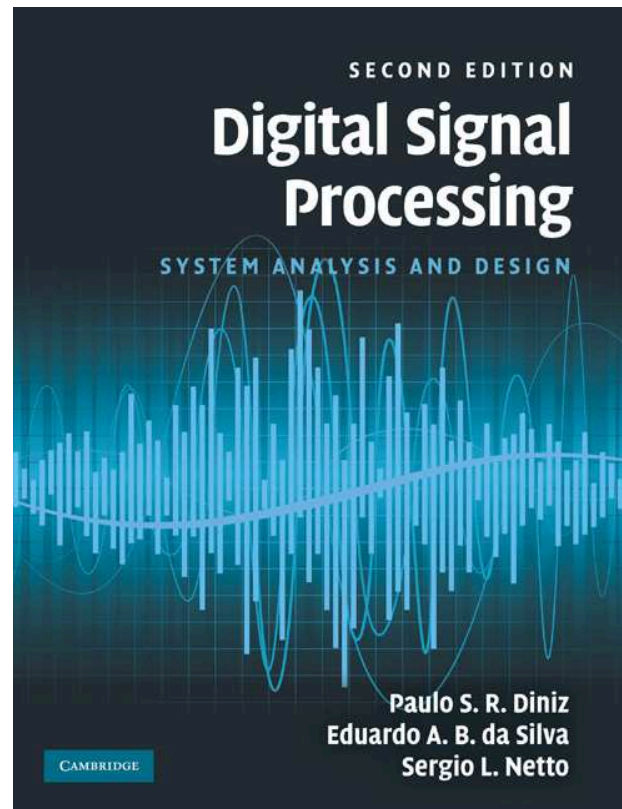


## Efficient IIR Structures



Paulo S. R. Diniz

Eduardo A. B. da Silva

Sergio L. Netto

`diniz,eduardo,sergioln@lps.ufrj.br`

September 2010

## Contents

- IIR parallel and cascade filters
- State-space sections
- Lattice filters
- Doubly complementary filters
- Wave filters
- Do-It-Yourself - Efficient IIR structures

## Introduction

The most widely used realizations for IIR filters are the cascade and parallel forms of second-order, and, sometimes, first-order, sections. The main advantages of these realizations come from their inherent modularity, which leads to efficient VLSI implementations, to simplified noise and sensitivity analyses, and to simple limit-cycle control.

We also deal with other interesting realizations such as the doubly-complementary filters, made from allpass blocks, and IIR lattice structures, whose synthesis method is presented. A related class of realizations are the wave digital filters, which have very low sensitivity and also allow the elimination of zero-input and overflow limit cycles.

## IIR parallel and cascade filters

- The  $N$ th-order IIR direct forms seen in Chapter 4 have roundoff-noise transfer functions  $G_i(z)$  and scaling transfer functions  $F_i(z)$  whose  $L_2$  or  $L_\infty$  norms assume significantly high values, because these transfer functions do not present the filter zeros to attenuate the gain introduced by the filter poles close to the unit circle.
- Also, in an  $N$ th-order IIR direct-form filter, a variation in a single coefficient causes variation on all the polynomial roots, leading to high sensitivity to coefficient quantization.
- To deal with these issues, it is wise to implement high-order transfer functions through the cascade or parallel connection of second-order building blocks, instead of using the direct-form realization. Such a structures also have the advantage of modularity, making them suitable for VLSI implementation.

## Parallel form

- A canonic parallel realization is shown in Figure 1, where the corresponding transfer function is given by

$$H(z) = h_0 + \sum_{i=1}^m H_i^p(z) = h_0 + \sum_{i=1}^m \frac{\gamma_{0i}z^2 + \gamma_{1i}z}{z^2 + m_{1i}z + m_{2i}} \quad (1)$$

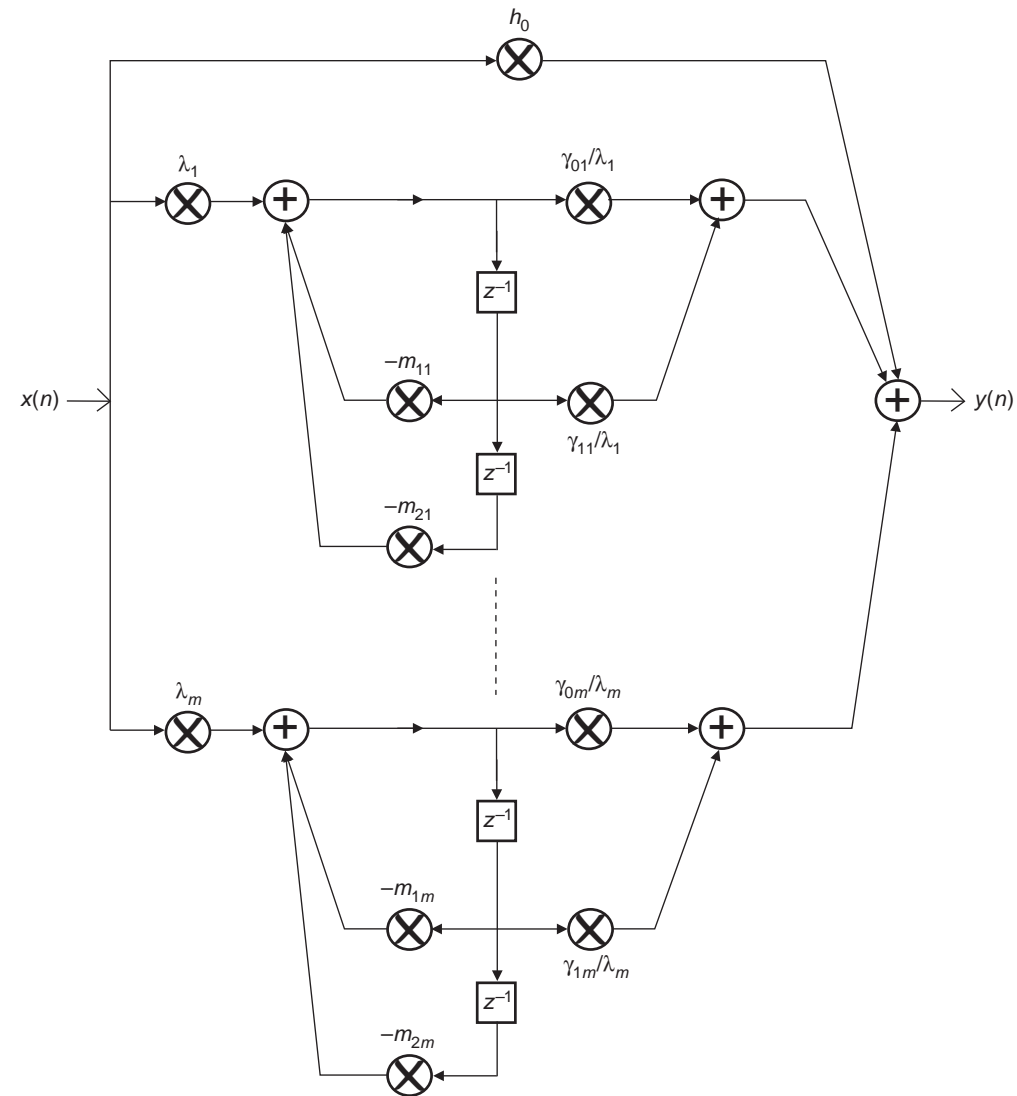


Figure 1: Parallel structure with direct-form sections.

## Parallel form

- It is easy to show that the scaling coefficients to avoid internal signal overflow in the form seen in Figure 1 are given by

$$\lambda_i = \frac{1}{\|F_i(z)\|_p} \quad (2)$$

where

$$F_i(z) = \frac{1}{D_i(z)} = \frac{1}{z^2 + m_{1i}z + m_{2i}} \quad (3)$$

- Naturally, the numerator coefficients of each section must be divided by  $\lambda_i$ , so that the overall filter transfer function remains unchanged.

## Parallel form

- The PSD of the output roundoff noise for the structure in Figure 1 is

$$\Gamma_Y(e^{j\omega}) = \sigma_e^2 \left( 2m + 1 + 3 \sum_{i=1}^m \frac{1}{\lambda_i^2} H_i^p(e^{j\omega}) H_i^p(e^{-j\omega}) \right) \quad (4)$$

when quantizations are performed before the additions.

- In this case, the output-noise variance, or the average power of the output noise, is

$$\sigma_o^2 = \sigma_e^2 \left( 2m + 1 + 3 \sum_{i=1}^m \frac{1}{\lambda_i^2} \|H_i^p(e^{j\omega})\|_2^2 \right) \quad (5)$$

- And the relative noise variance becomes

$$\sigma^2 = \frac{\sigma_o^2}{\sigma_e^2} = \left( 2m + 1 + 3 \sum_{i=1}^m \frac{1}{\lambda_i^2} \|H_i^p(e^{j\omega})\|_2^2 \right) \quad (6)$$



## Parallel form

- For the cases where quantization is performed after the additions, the PSD becomes

$$\Gamma_Y(e^{j\omega}) = \sigma_e^2 \left( 1 + \sum_{i=1}^m \frac{1}{\lambda_i^2} H_i^p(e^{j\omega}) H_i^p(e^{-j\omega}) \right) \quad (7)$$

and then

$$\sigma^2 = \frac{\sigma_o^2}{\sigma_e^2} = \left( 1 + \sum_{i=1}^m \frac{1}{\lambda_i^2} \|H_i^p(e^{j\omega})\|_2^2 \right) \quad (8)$$

- Although only even-order structures have been discussed so far, expressions for odd-order structures (containing one first-order section) are obtained in a similar way.
- In the parallel forms, as the positions of the zeros depend on the summation of several polynomials, which involves all filter coefficients, the precise positioning of the filter zeros becomes a difficult task. Such high sensitivity of the zeros to coefficient quantization constitutes the main drawback of the parallel forms for most practical implementations.

## Cascade form

- The cascade connection of direct-form second-order sections, depicted in Figure 2, has a transfer function given by

$$H(z) = \prod_{i=1}^m H_i(z) = \prod_{i=1}^m \frac{\gamma_{0i}z^2 + \gamma_{1i}z + \gamma_{2i}}{z^2 + m_{1i}z + m_{2i}} \quad (9)$$

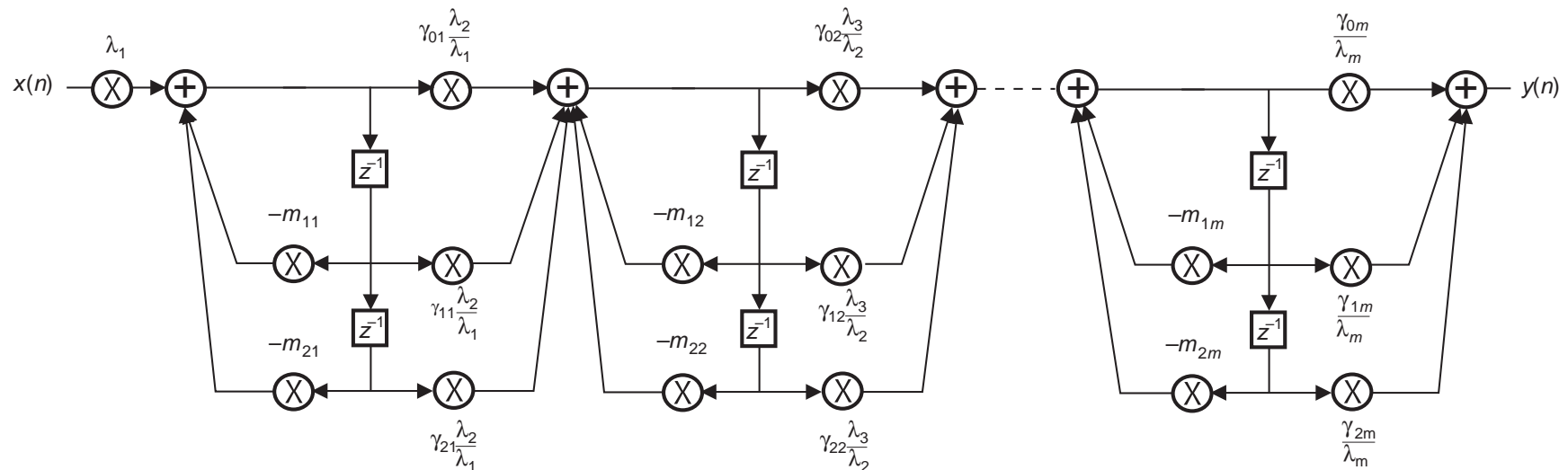


Figure 2: Cascade of direct-form sections.

## Cascade form

- In this structure, the scaling coefficients are calculated as

$$\lambda_i = \frac{1}{\| \prod_{j=1}^{i-1} H_j(z) F_i(z) \|_p} \quad (10)$$

with

$$F_i(z) = \frac{1}{D_i(z)} = \frac{1}{z^2 + m_{1i}z + m_{2i}} \quad (11)$$

as before.

- As illustrated in Figure 2, the scaling coefficient of each section can be incorporated with the output coefficients of the previous section. This strategy leads not only to a reduction in the multiplier count, but also to a possible decrease in the quantization noise at the filter output, since the number of nodes to be scaled is reduced.

## Cascade form

- Assuming that the quantizations are performed before the additions, the output PSD for the cascade structure of Figure 2 is given by

$$\Gamma_Y(e^{j\omega}) = \sigma_e^2 \left( 3 + \frac{3}{\lambda_1^2} \prod_{i=1}^m H_i(e^{j\omega}) H_i(e^{-j\omega}) + 5 \sum_{j=2}^m \frac{1}{\lambda_j^2} \prod_{i=j}^m H_i(e^{j\omega}) H_i(e^{-j\omega}) \right) \quad (12)$$

- The relative noise variance is then

$$\sigma^2 = \frac{\sigma_o^2}{\sigma_e^2} = \left( 3 + \frac{3}{\lambda_1^2} \left\| \prod_{i=1}^m H_i(e^{j\omega}) \right\|_2^2 + 5 \sum_{j=2}^m \frac{1}{\lambda_j^2} \left\| \prod_{i=j}^m H_i(e^{j\omega}) \right\|_2^2 \right) \quad (13)$$

## Cascade form

- For the case where quantizations are performed after additions,  $P_y(z)$  becomes

$$\Gamma_Y(e^{j\omega}) = \sigma_e^2 \left( 1 + \sum_{j=1}^m \frac{1}{\lambda_j^2} \prod_{i=j}^m H_i(e^{j\omega}) H_i(e^{-j\omega}) \right) \quad (14)$$

and then

$$\frac{\sigma_o^2}{\sigma_e^2} = \left( 1 + \sum_{j=1}^m \frac{1}{\lambda_j^2} \left\| \prod_{i=j}^m H_i(e^{j\omega}) \right\|_2^2 \right) \quad (15)$$

## Cascade form

- Two practical problems that need consideration in the design of cascade structures are:
  - Which pairs of poles and zeros will form each second-order section (the pairing problem).
  - The ordering of the sections.
- Both issues have a large effect on the output quantization noise. In fact, the roundoff noise and the sensitivity of cascade form structures can be very high if an inadequate choice for the pairing and ordering is made.

## Cascade form: Pole-zero pairing

- A rule of thumb for the pole-zero pairing in cascade form using second-order sections is to minimize the  $L_p$  norm of the transfer function of each section, for either  $p = 2$  or  $p = \infty$ .
- The pairs of complex conjugate poles close to the unit circle, if not accompanied by zeros which are close to them, tend to generate sections whose norms of  $H_i(z)$  are high. As a result, a natural rule is to pair the poles closest to the unit circle with the zeros that are closest to them.
- Then one should pick the poles second closest to the unit circle and pair them with the zeros, amongst the remaining, that are closest to them, and so on, until all sections are formed.
- Needless to say, when dealing with filters with real coefficients, most poles and zeros come in complex conjugate pairs, and in those cases the complex conjugate poles (and zeros) are jointly considered in the pairing process.

## Cascade form: Section ordering

- For section ordering, first we must notice that, for a given section of the cascade structure, the previous sections affect its scaling factor, whereas the following sections affect the noise gain.
- We then define a peaking factor that indicates how sharp the section frequency response is

$$P_i = \frac{\|H_i(z)\|_\infty}{\|H_i(z)\|_2} \quad (16)$$



## Cascade form: Section ordering

- We now consider two separate cases:
  - If we scale the filter using the  $L_2$  norm, then the scaling coefficients tend to be large, and thus the signal-to-noise ratio at the output of the filter is in general not problematic. In these cases, it is interesting to choose the section ordering so that the maximum value of the output-noise PSD,  $\|\text{PSD}\|_\infty$ , is minimized. Section  $i$  amplifies the  $\|\text{PSD}\|_\infty$  originally at its input by  $(\lambda_i \|H_i(e^{j\omega})\|_\infty)^2$ . Since in the  $L_2$  scaling,  $\lambda_i = \frac{1}{\|H_i(e^{j\omega})\|_2}$ , then each section amplifies the  $\|\text{PSD}\|_\infty$  by  $\left(\frac{\|H_i(e^{j\omega})\|_\infty}{\|H_i(e^{j\omega})\|_2}\right)^2 = P_i^2$ , the square of the peaking factor. Since the first sections affect the least number of noise sources, one should order the sections in decreasing order of peaking factors so as to minimize the maximum value of output-noise PSD.

## Cascade form: Section ordering

- Second case:
  - If we scale the filter using the  $L_\infty$  norm, then the scaling coefficients tend to be small, and thus the maximum peak value of the output-noise PSD is in general not problematic. In these cases, it is interesting to choose the section ordering so that the output signal-to-noise ratio is maximized, that is, the output-noise variance  $\sigma_o^2$  is minimized. Section  $i$  amplifies the output-noise variance at its input by  $(\lambda_i \|H_i(e^{j\omega})\|_2)^2$ . Since in the  $L_\infty$  scaling,  $\lambda_i = \frac{1}{\|H_i(e^{j\omega})\|_\infty}$ , then each section amplifies the  $\sigma_o^2$  by  $\left(\frac{\|H_i(e^{j\omega})\|_2}{\|H_i(e^{j\omega})\|_\infty}\right)^2 = \frac{1}{p_i^2}$ , the inverse of the square of the peaking factor. Since the first sections affect the least number of noise sources, one should order the sections in increasing order of peaking factors so as to minimize  $\sigma_o^2$ .
- For other types of scaling, both ordering strategies are considered equally efficient.

## Example 13.1

- Design an elliptic bandpass filter satisfying the following specifications:

$$\left. \begin{aligned}
 A_p &= 0.5 \text{ dB} \\
 A_r &= \text{dB} \\
 \Omega_{r_1} &= 850 \text{ rad/s} \\
 \Omega_{p_1} &= 980 \text{ rad/s} \\
 \Omega_{p_2} &= 1020 \text{ rad/s} \\
 \Omega_{r_2} &= 1150 \text{ rad/s} \\
 \Omega_s &= 10\,000 \text{ rad/s}
 \end{aligned} \right\} \quad (17)$$

- Realize the filter using the parallel and cascade forms of second-order direct-form sections. Then scale the filters using  $L_2$  norm and quantize the resulting coefficients to 9 bits, including the sign bit, and verify the results.

## Example 13.1 - Solution

- Using the `ellipord` and `ellip` commands in tandem, one can readily obtain the direct-form filter in MATLAB. We may then use the `residuez` command, and combine the resulting first-order sections, to determine the parallel structure, whose coefficients are shown in Table 1.

Table 1: Parallel structure using direct-form second-order sections. Feedforward coefficient:  $h_0 = -0.00015$ .

Coefficient	Section 1	Section 2	Section 3
$\gamma_0$	$-0.0077$	$-0.0079$	$0.0159$
$\gamma_1$	$0.0049$	$0.0078$	$-0.0128$
$m_1$	$-1.6268$	$-1.5965$	$-1.6054$
$m_2$	$0.9924$	$0.9921$	$0.9843$

### Example 13.1 - Solution

- Using  $L_2$  norm, each block can be scaled by

$$\lambda_i = \frac{1}{\|F_i(z)\|_2} = \frac{1}{\left\|\frac{1}{D_i(z)}\right\|_2} \quad (18)$$

which can be determined in MATLAB using the command lines

```
D_i = [1 m1i m2i];
```

```
F_i = freqz(1,D_i,npoints);
```

```
lambda_i = 1/sqrt(sum(abs(F_i).^2)/npoints);
```

where `npoints` is the number of points used in the `freqz` command. Scaling the second-order blocks using these factors, the resulting  $\gamma_0$  and  $\gamma_1$  coefficients are as given in Table 2, whereas the denominator coefficients  $m_1$  and  $m_2$  for each block remain unchanged.

Table 2: Scaled parallel structure using direct-form second-order sections. Feedforward coefficient:  $h_0 = -0.00015$ .

Coefficient	Section 1	Section 2	Section 3
$\lambda$	0.0711	0.0750	0.1039
$\frac{\gamma_0}{\lambda}$	-0.1077	-0.1055	0.1528
$\frac{\gamma_1}{\lambda}$	0.0692	0.1036	-0.1236
$m_1$	-1.6268	-1.5965	-1.6054
$m_2$	0.9924	0.9921	0.9843

## Example 13.1 - Solution

- Quantization of a given coefficient  $x$  using  $B$  bits (including the sign bit), can be performed in MATLAB using the command line:

$$x_Q = \text{quant}(x, 2^{-(B-1)});$$

Using this approach with  $(B-1) = 8$  results in the coefficients shown in Table 3.

Table 3: Parallel structure using direct-form second-order sections quantized with 9 bits.

Feedforward coefficient:  $[h_0]_Q = 0.0000$ .

Coefficient	Section 1	Section 2	Section 3
$[\lambda]_Q$	0.0703	0.0742	0.1055
$\left[\frac{\gamma_0}{\lambda}\right]_Q$	-0.1094	-0.1055	0.1523
$\left[\frac{\gamma_1}{\lambda}\right]_Q$	0.0703	0.1055	-0.1250
$[m_1]_Q$	-1.6250	-1.5977	-1.6055
$[m_2]_Q$	0.9922	0.9922	0.9844

## Example 13.1 - Solution

- The cascade form can be obtained from the direct form in MATLAB using the `tf2sos` command. This yields the coefficients shown in Table 4.

Table 4: Cascade structure using direct-form second-order sections. Gain constant:

$$h_0 = 1.4362 \text{ E} - 04.$$

Coefficient	Section 1	Section 2	Section 3
$\gamma_0$	1.0000	1.0000	1.0000
$\gamma_1$	0.0000	−1.4848	−1.7198
$\gamma_2$	−1.0000	1.0000	1.0000
$m_1$	−1.6054	−1.5965	−1.6268
$m_2$	0.9843	0.9921	0.9924



## Example 13.1 - Solution

- After section reordering and coefficient scaling, (a detailed implementation of these procedures is given in Experiment 13.1) the cascade realization is characterized as given in Table 5.

Table 5: Reordered cascade structure after coefficient scaling. Gain constant:  $h'_0 = h_0\lambda_2 = 0.0750$ .

Coefficient	Section 1'	Section 2'	Section 3'
$\gamma'_0$	0.1605	0.1454	0.0820
$\gamma'_1$	-0.2383	-0.2501	0.0000
$\gamma'_2$	0.1605	0.1454	-0.0820
$m'_1$	-1.5965	-1.6268	-1.6054
$m'_2$	0.9921	0.9924	0.9843

### Example 13.1 - Solution

- The quantized coefficients are as shown in Table 6. Notice that in this case, the structure gain is not quantized to avoid it to become zero.

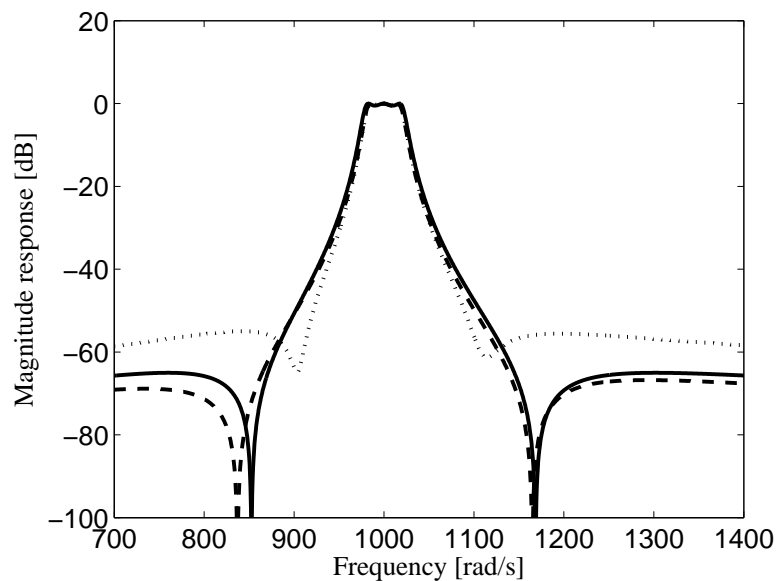
Table 6: Reordered cascade structure after coefficient quantization. Gain constant:

$$[h'_0]_Q = 0.0742.$$

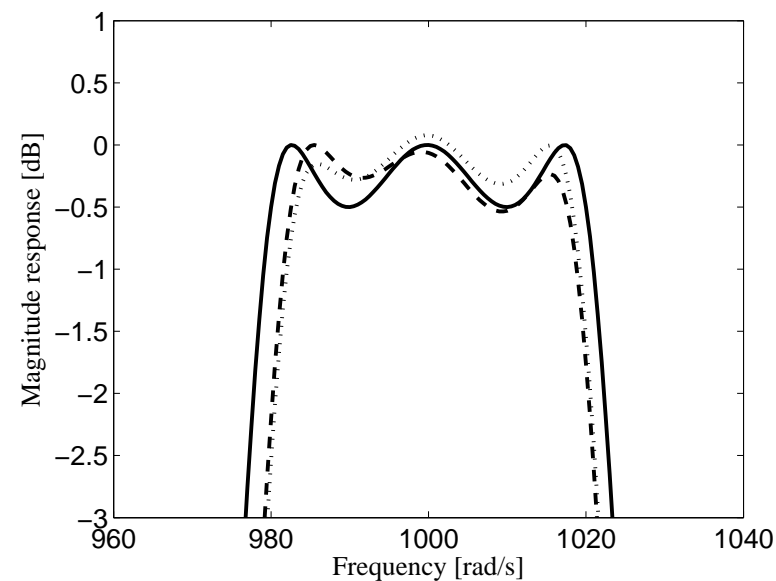
Coefficient	Section 1'	Section 2'	Section 3'
$[\gamma'_0]_Q$	0.1602	0.1445	0.0820
$[\gamma'_1]_Q$	-0.2383	-0.2500	0.0000
$[\gamma'_2]_Q$	0.1602	0.1445	-0.0820
$[m'_1]_Q$	-1.5977	-1.6250	-1.6055
$[m'_2]_Q$	0.9922	0.9922	0.9844

## Example 13.1 - Solution

- The magnitude responses for the ideal filter and the quantized parallel and cascade realizations are depicted in Figure 3. Note that despite the reasonably large number of bits used to represent the coefficients, the magnitude responses moved notably away from the ideal ones.



(a)



(b)

Figure 3: Coefficient-quantization effects in the cascade and parallel forms, using direct-form second-order sections: (a) overall magnitude response; (b) passband detail. (Solid line – initial design; dashed line – cascade of direct-form sections (9 bits); dotted line – parallel of direct-form sections (9 bits).)

## Error spectrum shaping

- We now have a technique to reduce the quantization noise effects on digital filters by feeding back the quantization error. This technique is known as error spectrum shaping (ESS) or error feedback.
- Consider every adder whose inputs include at least one nontrivial product which is followed by a quantizer. The ESS consists of replacing all these adders by a recursive structure, as illustrated in Figure 4, whose purpose is to introduce zeros in the output-noise PSD.

## Error spectrum shaping

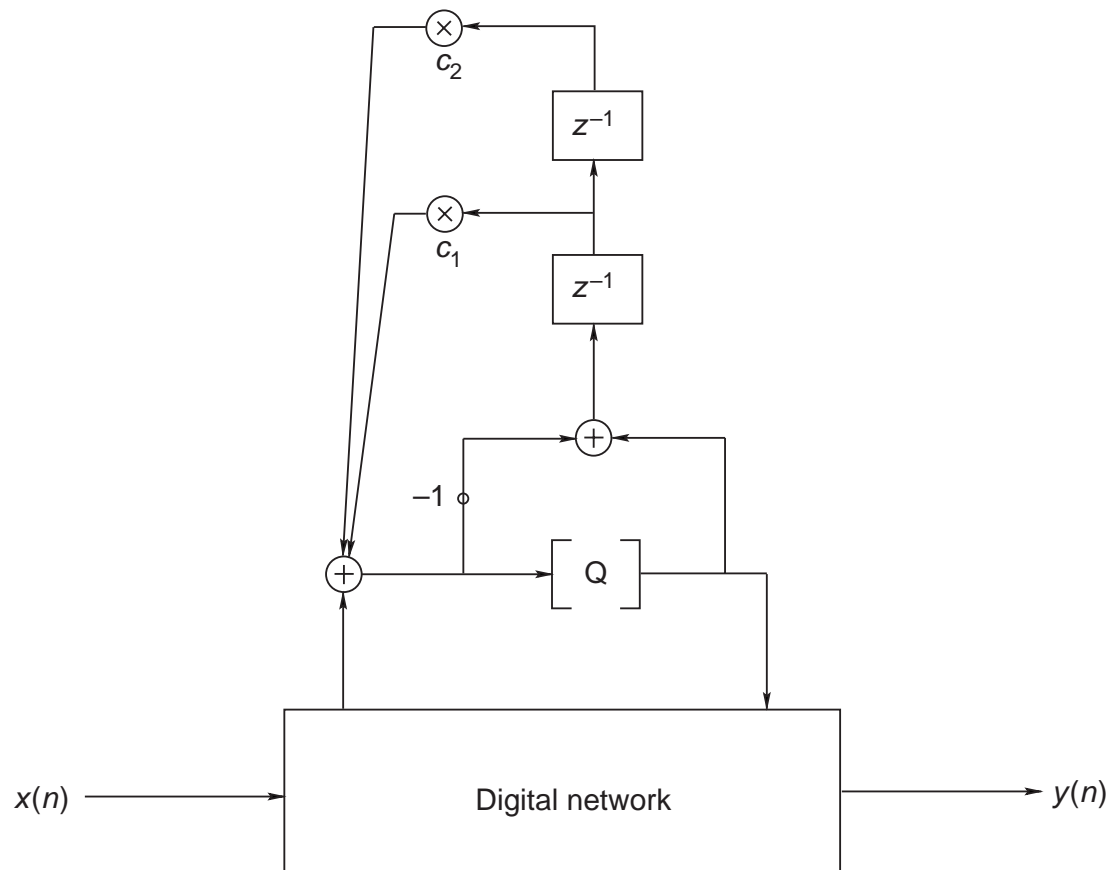


Figure 4: Error spectrum shaping structure (Q denotes quantizer).

## Error spectrum shaping

- Although Figure 4 depicts a second-order feedback network for the error signal, in practice the order of this network can assume any value.
- The ESS coefficients are chosen to minimize the output-noise PSD. In some cases, these coefficients can be made trivial and still achieve sufficient noise reduction.
- Overall, the ESS approach can be interpreted as a form of recycling the quantization error signal, thus reducing the effects of signal quantization after a particular adder.
- In theory, ESS can be applied to any internal quantization node of any digital filter. However, since it implies an implementation overhead, ESS should be applied only at selected internal nodes, whose noise gains to the filter output are high.
- Structures having reduced number of quantization nodes, for instance, are particularly suitable for ESS implementation. For example, the direct-form structure requires a single ESS substitution for the whole filter.

## Error spectrum shaping

- For the cascade structure of direct-form second-order sections, as in Figure 2, each section  $j$  requires an ESS substitution. Let each feedback network be of second order. The values of  $c_{1,j}$  and  $c_{2,j}$  that minimize the output noise are calculated by solving the following optimization problem:

$$\min_{c_{1,j}, c_{2,j}} \left\{ \left\| (1 + c_{1,j}z^{-1} + c_{2,j}z^{-2}) \prod_{i=j}^m H_i(e^{j\omega}) \right\|_2^2 \right\} \quad (19)$$



## Error spectrum shaping

- In this case, the optimal values of  $c_{1,j}$  and  $c_{2,j}$  are given by

$$c_{1,j} = \frac{t_1 t_2 - t_1 t_3}{t_3^2 - t_1^2}; \quad c_{2,j} = \frac{t_1^2 - t_2 t_3}{t_3^2 - t_1^2} \quad (20)$$

where

$$t_1 = \int_{-\pi}^{\pi} \left| \prod_{i=j}^m H_i(e^{j\omega}) \right|^2 \cos \omega d\omega \quad (21)$$

$$t_2 = \int_{-\pi}^{\pi} \left| \prod_{i=j}^m H_i(e^{j\omega}) \right|^2 \cos(2\omega) d\omega \quad (22)$$

$$t_3 = \int_{-\pi}^{\pi} \left| \prod_{i=j}^m H_i(e^{j\omega}) \right|^2 d\omega \quad (23)$$

## Error spectrum shaping

- For a first-order ESS, the optimal value of  $c_{1,j}$  would be

$$c_{1,j} = \frac{-t_1}{t_3}. \quad (24)$$

- Using ESS, with the quantization performed after summation, the output relative power spectrum density (RPSD), which is independent of  $\sigma_e^2$  for the cascade design, is given by

$$\text{RPSD} = 1 + \sum_{j=1}^m \left| \frac{1}{\lambda_j} (1 + c_{1,j}z^{-1} + c_{2,j}z^{-2}) \prod_{i=j}^m H_i(e^{j\omega}) \right|^2 \quad (25)$$

where  $\lambda_j$  is the scaling factor of section  $j$ . This expression explicitly shows how the ESS technique introduces zeros in the RPSD, thus allowing its subsequent reduction.

## Closed-form scaling

- In most types of digital filter implemented with fixed-point arithmetic, scaling is based on the  $L_2$  and  $L_\infty$  norms of the transfer functions from the filter inputs to the inputs of the multipliers.
- Usually, the  $L_2$  norm is computed through summation of a large number of sample points (of the order of 200 or more) of the squared magnitude of the scaling transfer function.
- For the  $L_\infty$  norm, a search for the maximum magnitude of the scaling transfer function is performed over about the same number of sample points.
- It is possible, however, to derive simple closed-form expressions for the  $L_2$  and  $L_\infty$  norms of second-order transfer functions. Such expressions are useful for scaling the sections independently, and greatly facilitate the design of parallel and cascade realizations of second-order sections.

## Closed-form scaling

- Consider, for example,

$$H(z) = \frac{\gamma_1 z + \gamma_2}{z^2 + m_1 z + m_2} \quad (26)$$

- Using, for instance, the pole-residue approach for solving circular integrals, the corresponding  $L_2$  norm is given by

$$\|H(e^{j\omega})\|_2^2 = \frac{\gamma_1^2 + \gamma_2^2 - 2\gamma_1\gamma_2 \frac{m_1}{m_2 + 1}}{(1 - m_2^2) \left[ 1 - \left( \frac{m_1}{m_2 + 1} \right)^2 \right]} \quad (27)$$

## Closed-form scaling

- For the  $L_\infty$  norm, we have to find the maximum of  $|H(z)|^2$ . By noticing that  $|H(e^{j\omega})|^2$  is a function of  $\cos \omega$ , and  $\cos \omega$  is limited to the interval  $[-1, 1]$ , we have that the maximum is either at the extrema  $\omega = 0$  ( $z = 1$ ),  $\omega = \pi$  ( $z = -1$ ), or at  $\omega_0$  such that  $-1 \leq \cos \omega_0 \leq 1$ .
- Therefore, the  $L_\infty$  norm is given by

$$\|H(e^{j\omega})\|_\infty^2 = \max \left\{ \left( \frac{\gamma_1 + \gamma_2}{1 + m_1 + m_2} \right)^2, \left( \frac{-\gamma_1 + \gamma_2}{1 - m_1 + m_2} \right)^2, \frac{\gamma_1^2 + \gamma_2^2 + 2\gamma_1\gamma_2\zeta}{4m_2[(\zeta - \eta)^2 + v]} \right\}, \quad (28)$$

## Closed-form scaling

- Where

$$\eta = \frac{-m_1(1+m_2)}{4m_2}; v = \left(1 - \frac{m_1^2}{4m_2}\right) \frac{(1-m_2)^2}{4m_2}; \gamma = \frac{\gamma_1^2 + \gamma_2^2}{2\gamma_1\gamma_2} \quad (29)$$

$$\zeta = \begin{cases} \text{sat}(\eta), & \text{for } \gamma_1\gamma_2 = 0 \\ \text{sat} \left\{ \gamma \left[ \sqrt{\left(1 + \frac{\eta}{\gamma}\right)^2 + \frac{v}{\gamma^2}} - 1 \right] \right\}, & \text{for } \gamma_1\gamma_2 \neq 0 \end{cases} \quad (30)$$

with  $\text{sat}(\cdot)$  being defined as

$$\text{sat}(x) = \begin{cases} 1, & \text{for } x > 1 \\ -1, & \text{for } x < -1 \\ x, & \text{for } -1 \leq x \leq 1 \end{cases} \quad (31)$$

## Example 13.2

- Given the transfer function below:

$$H(z) = \frac{(0.5z^2 - z + 1)}{(z^2 - z + 0.5)(z + 0.5)} \quad (32)$$

1. Show cascade and parallel decompositions using Type 1 direct-form sections.
2. Scale the filters using  $L_2$  norm.
3. Calculate the output noise variances.

## Example 13.2 - Solution

- The cascade decomposition is

$$H(z) = \frac{(0.5z^2 - z + 1)}{(z^2 - z + 0.5)} \frac{1}{z + 0.5} \quad (33)$$

- The second-order section of the cascade design is an allpass so that we only have to scale the internal nodes of the section.
- The result is obtained by employing equation (27), that is

$$\|F_1(z)\|_2^2 = \left\| \frac{1}{D(z)} \right\|_2^2 = \frac{1}{(1 - 0.25) \left( 1 - \left( \frac{-1}{1.5} \right)^2 \right)} = 1.44 \quad (34)$$

so that

$$\lambda_1 = \sqrt{\frac{1}{1.44}} = \frac{1}{1.2} = 0.8333 \quad (35)$$



### Example 13.2 - Solution

- For the second section the scaling factor should be

$$\lambda_2 = \sqrt{0.75} = \sqrt{1 - (0.5)^2} \quad (36)$$

- The relative output noise variance for the cascade design is given by

$$\frac{\sigma_y^2}{\sigma_e^2} = 3 \frac{1}{\lambda_1^2} \frac{1}{0.75} + 4 \frac{1}{\lambda_2^2} \frac{1}{0.75} + 1 = 13.88 \quad (37)$$

### Example 13.2 - Solution

- The parallel decomposition is

$$H(z) = \frac{\left(-\frac{8}{5}z + \frac{7}{5}\right)}{(2z^2 - 2z + 1)} + \frac{\frac{13}{10}}{z + 0.5} \quad (38)$$

- The scaling factors for the parallel realization are given by

$$\lambda_1 = \sqrt{\frac{1}{1.44}} = \frac{1}{1.2} = 0.8333 \quad (39)$$

and

$$\lambda_2 = \frac{10}{13}\sqrt{0.75} = 0.6667 \quad (40)$$

respectively.

### Example 13.2 - Solution

- Using the result of equation (27) we can compute the  $L_2$  norm of the second-order section in the parallel solution.

$$\|H_1(z)\|_2^2 = \frac{1}{4} \left[ \frac{\frac{64}{25} + \frac{49}{25} + \frac{2 \times 8 \times 7}{25} \times \frac{-1}{1.5}}{(1 - 0.25) \left(1 - \left(\frac{-1}{1.5}\right)^2\right)} \right] = \frac{1}{4} \left[ \frac{\frac{38.333}{25}}{0.41666} \right] = 0.92 \quad (41)$$

so that the relative output noise variance for the parallel design is given by

$$\frac{\sigma_y^2}{\sigma_e^2} = 3 \frac{1}{\lambda_1^2} \|H_1(z)\|_2^2 + \frac{1}{\lambda_2^2} \frac{1}{0.75} + 4 = 10.98 \quad (42)$$

## State-space sections

- The state-space approach allows the formulation of a design method for IIR digital filters with minimum roundoff noise. The theory behind this elegant design method was originally proposed by Mullis and Roberts. For a filter of order  $N$ , the minimum noise method leads to a realization entailing  $(N + 1)^2$  multiplications. This multiplier count is very high for most practical implementations, which induced investigators to search for realizations which could approach the minimum noise performance while employing a reasonable number of multiplications. A good tradeoff is achieved if we realize high-order filters using parallel or cascade forms, where the second-order sections are minimum-noise state-space structures. In this section, we study two commonly used second-order state-space sections suitable for such approaches.

## Optimal state-space sections

- The second-order state-space structure shown in Figure 5 can be described by

$$\left. \begin{aligned} \mathbf{x}(n+1) &= \mathbf{A}\mathbf{x}(n) + \mathbf{B}u(n) \\ y(n) &= \mathbf{C}^T\mathbf{x}(n) + \mathbf{D}u(n) \end{aligned} \right\} \quad (43)$$

where  $\mathbf{x}(n)$  is a column vector representing the outputs of the delays,  $y(n)$  is a scalar, and

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}; \mathbf{B} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}; \mathbf{C}^T = \begin{bmatrix} c_1 & c_2 \end{bmatrix}; \mathbf{D} = \begin{bmatrix} d \end{bmatrix}. \quad (44)$$

## Optimal state-space sections

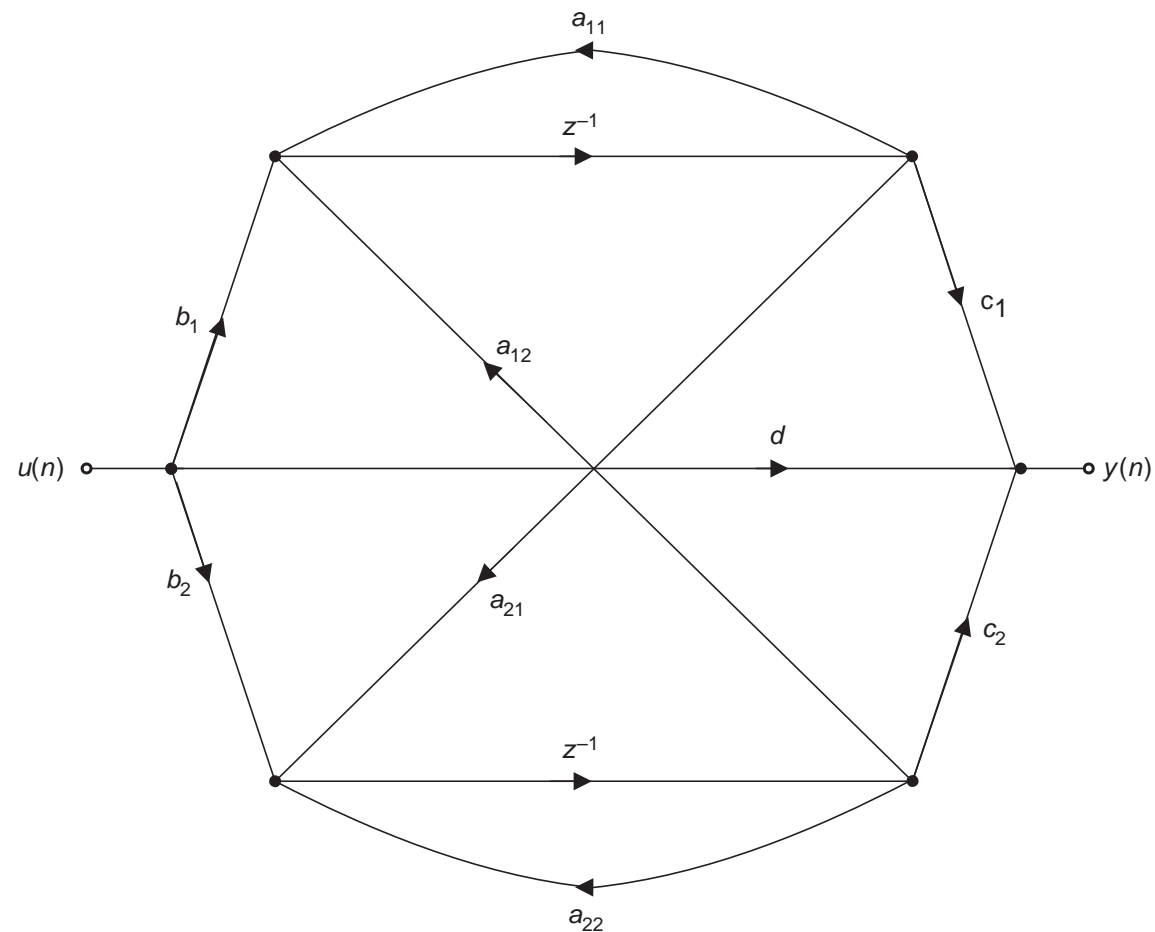


Figure 5: Second-order state-space structure.

## Optimal state-space sections

- The overall transfer function, described as a function of the matrix elements related to the state-space formulation, is given by

$$H(z) = \mathbf{C}^T [\mathbf{I}z - \mathbf{A}]^{-1} \mathbf{B} + \mathbf{D}. \quad (45)$$

- The second-order state-space structure can realize transfer functions described by

$$H(z) = d + \frac{\gamma_1 z + \gamma_2}{z^2 + m_1 z + m_2}. \quad (46)$$

- Given  $H(z)$  in the form of equation (46), an optimal design in the sense of minimizing the output roundoff noise can be derived, since the state-space structure has more coefficients than the minimum required.

## Optimal state-space sections

- To explore this feature, we examine, without proof, a theorem first proposed by Mullis and Roberts. The design procedure resulting from the theorem generates realizations with a minimum-variance output noise, provided that the  $L_2$  norm is employed to determine the scaling factor.
- It is interesting to notice that, despite being developed for filters using  $L_2$  scaling, the minimum-noise design also leads to low-noise filters scaled using the  $L_\infty$  norm.
- Note that, in the remainder of this subsection, primed variables will indicate filter parameters after scaling.



## Optimal state-space sections

- Theorem 13.1: The necessary and sufficient conditions to obtain an output noise with minimum variance in a state-space realization are given by

$$\mathbf{W}' = \mathbf{R}\mathbf{K}'\mathbf{R} \quad (47)$$

$$\mathbf{K}'_{ii}\mathbf{W}'_{ii} = \mathbf{K}'_{jj}\mathbf{W}'_{jj} \quad (48)$$

for  $i, j = 1, 2, \dots, N$ , where  $N$  is the filter order,  $\mathbf{R}$  is an  $N \times N$  diagonal matrix, and

$$\mathbf{K}' = \sum_{k=0}^{\infty} \mathbf{A}'^k \mathbf{B}' \mathbf{B}'^H (\mathbf{A}'^k)^H \quad (49)$$

$$\mathbf{W}' = \sum_{k=0}^{\infty} (\mathbf{A}'^k)^H \mathbf{C}'^H \mathbf{C}' \mathbf{A}'^k \quad (50)$$

where the  $H$  indicates the conjugate and transpose operations.

## Optimal state-space sections

- It can be shown that

$$K'_{ii} = \|F'_i(e^{j\omega})\|_2^2 \quad (51)$$

$$W'_{ii} = \|G'_i(e^{j\omega})\|_2^2 \quad (52)$$

for  $i = 1, 2, \dots, N$ , where  $F'_i(z)$  is the transfer function from the scaled filter input to the state variable  $x_i(k+1)$ , and  $G'_i(z)$  is the transfer function from the state variable  $x_i(k)$  to the scaled filter output.

- Then, from equations (51) and (52), we have that, in the frequency domain, equation (48) is equivalent to

$$\|F'_i(e^{j\omega})\|_2^2 \|G'_i(e^{j\omega})\|_2^2 = \|F'_j(e^{j\omega})\|_2^2 \|G'_j(e^{j\omega})\|_2^2 \quad (53)$$

## Optimal state-space sections

- In the case of second-order filters, if the  $L_2$  scaling is performed, then

$$K'_{11} = K'_{22} = \|F'_1(e^{j\omega})\|_2^2 = \|F'_2(e^{j\omega})\|_2^2 = 1 \quad (54)$$

and then, from Theorem 13.1, the following equality must hold

$$W'_{11} = W'_{22} \quad (55)$$

- Similarly, we can conclude that we must have

$$\|G'_1(e^{j\omega})\|_2^2 = \|G'_2(e^{j\omega})\|_2^2 \quad (56)$$

indicating that the contributions of the internal noise sources to the output-noise variance are identical.

## Optimal state-space sections

- The conditions  $K'_{ii} = \|F'_i(e^{j\omega})\|_2^2 = 1$  and  $W'_{ii} = W'_{jj}$ , for all  $i$  and  $j$ , show that equation (47) can only be satisfied if

$$\mathbf{R} = \alpha \mathbf{I} \quad (57)$$

and, as a consequence, the optimality condition of Theorem 13.1 is equivalent to

$$\mathbf{W}' = \alpha^2 \mathbf{K}' \quad (58)$$

- For a second-order filter, since  $\mathbf{W}'$  and  $\mathbf{K}'$  are symmetric and their respective diagonal elements are identical, equation (58) remains valid if we rewrite it as

$$\mathbf{W}' = \alpha^2 \mathbf{J} \mathbf{K}' \mathbf{J} \quad (59)$$

where  $\mathbf{J}$  is the reverse identity matrix defined as

$$\mathbf{J} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad (60)$$

## Optimal state-space sections

- By employing the definitions of  $\mathbf{W}'$  and  $\mathbf{K}'$  in equations (49) and (50), equation (59) is satisfied when

$$\mathbf{A}'^T = \mathbf{J}\mathbf{A}'\mathbf{J} \quad (61)$$

$$\mathbf{C}'^T = \alpha\mathbf{J}\mathbf{B}' \quad (62)$$

or, equivalently

$$a'_{11} = a'_{22} \quad (63)$$

$$\frac{b'_1}{b'_2} = \frac{c'_2}{c'_1} \quad (64)$$

- Then, the following procedure can be derived for designing optimal second-order state-space sections.

## Optimal state-space sections

- Step 1: For filters with complex conjugate poles, choose an antisymmetric matrix **A** such that

$$\left. \begin{aligned} a_{11} &= a_{22} = \text{real part of the poles} \\ -a_{12} &= a_{21} = \text{imaginary part of the poles} \end{aligned} \right\} \quad (65)$$

Note that the first optimality condition (equation (63)) is satisfied by this choice of **A**. The coefficients of matrix **A** can be calculated as a function of the coefficients of the transfer function  $H(z)$  using

$$\left. \begin{aligned} a_{11} &= -\frac{m_1}{2} \\ a_{12} &= -\sqrt{m_2 - \frac{m_1^2}{4}} \\ a_{21} &= -a_{12} \\ a_{22} &= a_{11} \end{aligned} \right\} \quad (66)$$

## Optimal state-space sections

- Also, compute the parameters  $b_1$ ,  $b_2$ ,  $c_1$ , and  $c_2$ , using

$$\left. \begin{aligned} b_1 &= \sqrt{\frac{\sigma + \gamma_2 + a_{11}\gamma_1}{2a_{21}}} \\ b_2 &= \frac{\gamma_1}{2b_1} \\ c_1 &= b_2 \\ c_2 &= b_1 \end{aligned} \right\} \quad (67)$$

where

$$\sigma = \sqrt{\gamma_2^2 - \gamma_1\gamma_2 m_1 + \gamma_1^2 m_2} \quad (68)$$

## Optimal state-space sections

- For real poles, the matrix  $\mathbf{A}$  must be of the form

$$\mathbf{A} = \begin{bmatrix} a_1 & a_2 \\ a_2 & a_1 \end{bmatrix} \quad (69)$$

where

$$\left. \begin{aligned} a_1 &= \frac{1}{2}(p_1 + p_2) \\ a_2 &= \pm \frac{1}{2}(p_1 - p_2) \end{aligned} \right\} \quad (70)$$

with  $p_1$  and  $p_2$  denoting the real poles.



- The elements of vectors **B** and **C** are given by

$$\left. \begin{aligned} b_1 &= \pm \sqrt{\frac{\pm\sigma + \gamma_2 + a_1\gamma_1}{2a_2}} \\ b_2 &= \frac{\gamma_2}{2b_1} \\ c_1 &= b_2 \\ c_2 &= b_1 \end{aligned} \right\} \quad (71)$$

with  $\sigma$  as before.

## Optimal state-space sections

- Step 2: Scale the filter using  $L_2$  norm, through the following similarity transformation

$$(\mathbf{A}', \mathbf{B}', \mathbf{C}', d) = (\mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \mathbf{T}^{-1} \mathbf{B}, \mathbf{C} \mathbf{T}, d) \quad (72)$$

where

$$\mathbf{T} = \begin{bmatrix} \|F_1(e^{j\omega})\|_2 & 0 \\ 0 & \|F_2(e^{j\omega})\|_2 \end{bmatrix} \quad (73)$$

## Optimal state-space sections

- For the  $L_\infty$ -norm scaling, use the following scaling matrix

$$\mathbf{T} = \begin{bmatrix} \|F_1(e^{j\omega})\|_\infty & 0 \\ 0 & \|F_2(e^{j\omega})\|_\infty \end{bmatrix} \quad (74)$$

In this case the resulting second-order section is not optimal in the  $L_\infty$  sense.

Nevertheless, practical results indicate that the solution is close to the optimal one.

## Optimal state-space sections

- Defining the vectors  $\mathbf{f}(z) = [F_1(z), F_2(z)]^T$  and  $\mathbf{g}(z) = [G_1(z), G_2(z)]^T$ , the effects of the scaling matrix on these vectors are

$$\mathbf{f}'(z) = [z\mathbf{I} - \mathbf{A}']^{-1} \mathbf{B}' = \mathbf{T}^{-1} \mathbf{f}(z) \quad (75)$$

$$\mathbf{g}'(z) = [z\mathbf{I} - \mathbf{A}'^T]^{-1} \mathbf{C}'^T = (\mathbf{T}^{-1})^T \mathbf{g}(z) \quad (76)$$

## Optimal state-space sections

- The transfer functions  $F_i(z)$  from the input node, where  $u(n)$  is inserted, to the state variables  $x_i(n)$  of the system  $(\mathbf{A}, \mathbf{B}, \mathbf{C}, d)$  are given by

$$F_1(z) = \frac{b_1 z + (b_2 a_{12} - b_1 a_{22})}{z^2 - (a_{11} + a_{22})z + (a_{11} a_{22} - a_{12} a_{21})} \quad (77)$$

$$F_2(z) = \frac{b_2 z + (b_1 a_{21} - b_2 a_{11})}{z^2 - (a_{11} + a_{22})z + (a_{11} a_{22} - a_{12} a_{21})} \quad (78)$$

- The expressions for the transfer functions from the internal nodes, that is, from the signals  $x_i(n+1)$  to the section output node are

$$G_1(z) = \frac{c_1 z + (c_2 a_{21} - c_1 a_{22})}{z^2 - (a_{11} + a_{22})z + (a_{11} a_{22} - a_{12} a_{21})} \quad (79)$$

$$G_2(z) = \frac{c_2 z + (c_1 a_{12} - c_2 a_{11})}{z^2 - (a_{11} + a_{22})z + (a_{11} a_{22} - a_{12} a_{21})} \quad (80)$$

## Optimal state-space sections

- The output roundoff-noise PSD, considering quantization before the adders, for the section-optimal state-space structure in cascade form can be expressed as

$$\Gamma_Y(e^{j\omega}) = 3\sigma_e^2 \sum_{j=1}^m \prod_{l=j+1}^m H_l(e^{j\omega}) H_l(e^{-j\omega}) \left( 1 + \sum_{i=1}^2 G'_{ij}(e^{j\omega}) G'_{ij}(e^{-j\omega}) \right) \quad (81)$$

where  $G'_{ij}$ , for  $i = 1, 2$ , are the noise transfer functions of the  $j$ th scaled section, and we consider  $\prod_{l=m+1}^m H_l(z) H_l(z^{-1}) = 1$ .

## Optimal state-space sections

- The scaling in the state-space sections of the cascade form is performed internally, using the transformation matrix  $\mathbf{T}$ . In order to calculate the elements of matrix  $\mathbf{T}$ , we can use the same procedure as in the cascade direct form, taking the effect of previous blocks into consideration.
- In the case of the parallel form, the expression for the output roundoff-noise PSD, assuming quantization before additions, is

$$\Gamma_Y(e^{j\omega}) = \sigma_e^2 \left( 2m + 1 + 3 \sum_{j=1}^m \sum_{i=1}^2 G'_{ij}(e^{j\omega}) G'_{ij}(e^{-j\omega}) \right) \quad (82)$$

## State-space sections without limit cycles

- This section presents a design procedure for a second-order state-space section that is free from constant-input limit cycles.
- The transition matrix related to the section-optimal structure, described in the previous subsection (see equation (63)), has the following general form

$$\mathbf{A} = \begin{bmatrix} \alpha & -\frac{\zeta}{\sigma} \\ \zeta\sigma & \alpha \end{bmatrix} \quad (83)$$

where  $\alpha$ ,  $\zeta$ , and  $\sigma$  are constants. This form is the most general for  $\mathbf{A}$  that allows the realization of complex conjugate poles and the elimination of zero-input limit cycles.



## State-space sections without limit cycles

- As studied in Subsection 7.6.3, one can eliminate zero-input limit cycles on a recursive structure if there is a positive-definite diagonal matrix  $\mathbf{G}$ , such that  $(\mathbf{G} - \mathbf{A}^T \mathbf{G} \mathbf{A})$  is positive semidefinite.
- For second-order sections, this condition is satisfied if

$$a_{12}a_{21} \geq 0 \quad (84)$$

or

$$a_{12}a_{21} < 0 \text{ and } |a_{11} - a_{22}| + \det\{\mathbf{A}\} \leq 1 \quad (85)$$

- In the section-optimal structures, the elements of matrix  $\mathbf{A}$  automatically satisfy equation (84), since  $a_{11} = a_{22}$  and  $\det(\mathbf{A}) \leq 1$ , for stable filters.

## State-space sections without limit cycles

- Naturally, the quantization performed at the state variables still must be such that

$$|[\mathbf{x}_i(k)]_Q| \leq |\mathbf{x}_i(k)|, \text{ for all } k \quad (86)$$

where  $[\mathbf{x}]_Q$  denotes the quantized value of  $\mathbf{x}$ . This condition can be easily guaranteed by using, for example, magnitude truncation and saturation arithmetic to deal with overflow.

- If we also want to eliminate constant-input limit cycles, according to Theorem 11.3, the values of the elements of  $\mathbf{p}\mathbf{u}_0$ , where  $\mathbf{p} = (\mathbf{I} - \mathbf{A})^{-1}\mathbf{b}$ , must be machine representable. In order to guarantee this condition independently of  $\mathbf{u}_0$ , the column vector  $\mathbf{p}$  must assume one of the forms

$$\mathbf{p} = \begin{cases} [\pm 1 \ 0]^T & \text{(Case I)} \\ [0 \ \pm 1]^T & \text{(Case II)} \\ [\pm 1 \ \pm 1]^T & \text{(Case III)} \end{cases} \quad (87)$$

## State-space sections without limit cycles

- For each case above, vector **B** of the state-space structure must be appropriately chosen to ensure the elimination of constant-input limit cycles, as given by

– Case I:

$$\left. \begin{aligned} b_1 &= \pm(1 - a_{11}) \\ b_2 &= \mp a_{21} \end{aligned} \right\} \quad (88)$$

– Case II:

$$\left. \begin{aligned} b_1 &= \mp a_{12} \\ b_2 &= \pm(1 - a_{22}) \end{aligned} \right\} \quad (89)$$

– Case III:

$$\left. \begin{aligned} b_1 &= \mp a_{12} \pm (1 - a_{11}) \\ b_2 &= \mp a_{21} \pm (1 - a_{22}) \end{aligned} \right\} \quad (90)$$

## State-space sections without limit cycles

- Based on the values of  $b_1$  and  $b_2$ , for each case, it is possible to generate three structures, henceforth referred to as Structures I, II, and III. Figure 6 depicts Structure I, where it can be seen that  $b_1$  and  $b_2$  are formed without actual multiplications. As a consequence, the resulting structure is more economical than the optimal second-order state-space structure. Similar results apply to all three structures. In fact, Structures I and II have the same complexity, whereas Structure III requires 5 extra additions, if we consider the adders needed for the elimination of constant-input limit cycles. For that reason, in what follows, we present the design for Structure I.

## State-space sections without limit cycles

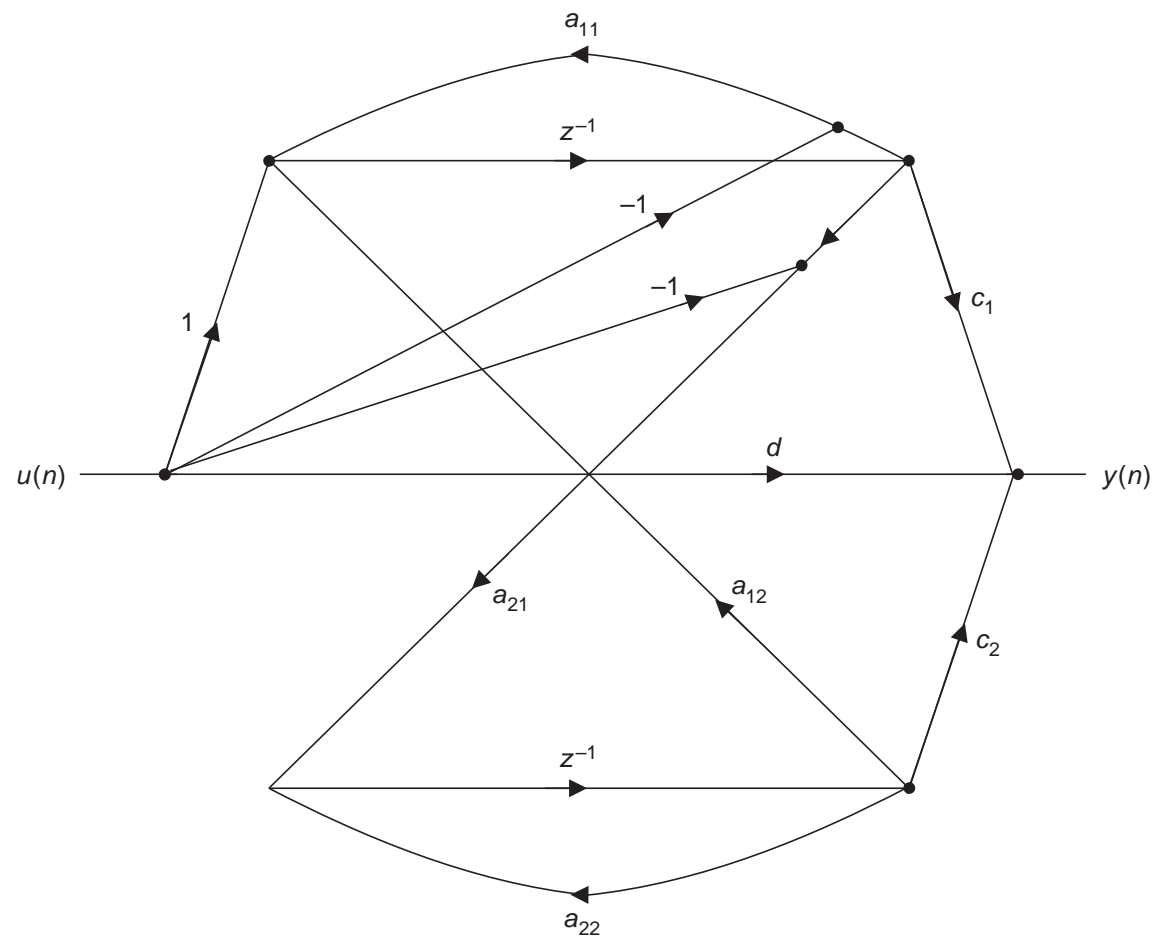


Figure 6: State-space structure free from limit cycles.

## State-space sections without limit cycles

- For Structure I, we have that

$$a_{11} = a; a_{12} = -\frac{\zeta}{\sigma}; a_{21} = \sigma\zeta; a_{22} = a_{11} \quad (91)$$

$$b_1 = 1 - a_{11}; b_2 = -a_{21} \quad (92)$$

$$c_1 = \frac{\gamma_1 + \gamma_2}{1 + m_1 + m_2}; c_2 = -\frac{(m_1 + 2m_2)\gamma_1 + (2 + m_1)\gamma_2}{2\sigma\zeta(1 + m_1 + m_2)} \quad (93)$$

where

$$a = -\frac{m_1}{2}; \zeta = \sqrt{\left(m_2 - \frac{m_1^2}{4}\right)} \quad (94)$$

and  $\sigma$  is a free parameter whose choice is explained below.

## State-space sections without limit cycles

- From the above equations,

$$\frac{b_1}{b_2} = -\frac{2 + m_1}{2\sigma\zeta} \quad (95)$$

$$\frac{c_2}{c_1} = -\frac{(m_1 + 2m_2)\gamma_1 + (2 + m_1)\gamma_2}{2\sigma\zeta(\gamma_1 + \gamma_2)} \quad (96)$$

- Therefore, in this case, the optimality condition derived from Theorem 13.1 (equations (63) and (64)) is only met if

$$\frac{\gamma_1}{\gamma_2} = \frac{m_1 + 2}{m_2 - 1} \quad (97)$$

## State-space sections without limit cycles

- Usually, this condition is violated, showing that the state-space structure which is free from constant-input limit cycles does not lead to minimum output noise.
- In practice, however, it is found that the performance of Structure I is very close to the optimal. More specifically, in the case where the zeros of  $H(z)$  are placed at  $z = 1$ , the values of  $\gamma_1$  and  $\gamma_2$  are

$$\left. \begin{aligned} \gamma_1 &= -\gamma_0(2 + m_1) \\ \gamma_2 &= \gamma_0(1 - m_2) \end{aligned} \right\} \quad (98)$$

which satisfy equation (97). Hence, in the special case of filters with zeros at  $z = 1$ , Structure I also leads to minimum output noise.



## State-space sections without limit cycles

- Parameter  $\sigma$  is usually used to optimize the dynamic range of the state variables.
- For Structure I, the transfer functions from the input node  $u(k)$  to the state variables  $x_i(k)$  are given by

$$F_1(z) = \frac{(1-a)z + (\zeta^2 - a + a^2)}{z^2 - 2az + (a^2 + \zeta^2)}; \quad F_2(z) = \sigma F_2''(z) \quad (99)$$

where

$$F_2''(z) = \frac{-\zeta z + \zeta}{z^2 - 2az + (a^2 + \zeta^2)} \quad (100)$$

- Equalization of the maximum signal level at the state variables is achieved by forcing

$$\|F_1(z)\|_p = \|\sigma F_2''(z)\|_p \quad (101)$$

where  $p = \infty$  or  $p = 2$ . Consequently, we must have

$$\sigma = \frac{\|F_1(z)\|_p}{\|F_2''(z)\|_p} \quad (102)$$

## State-space sections without limit cycles

- The transfer functions from the state variables  $x_i(k+1)$  to the output in the structure of Figure 6 can be expressed as

$$G_1(z) = \frac{c_1}{2} \frac{2z + (m_1 + 2\xi)}{z^2 + m_1z + m_2} \quad (103)$$

$$G_2(z) = \frac{c_2}{2} \frac{2z + (\alpha_1 + \frac{2\xi^2}{\xi})}{z^2 + \alpha_1z + \alpha_2} \quad (104)$$

where

$$\xi = \frac{-(\alpha_1 + 2\alpha_2)\beta_1 + (2 + \alpha_1)\beta_2}{2(\beta_1 + \beta_2)} \quad (105)$$

## State-space sections without limit cycles

- The RPSD expression for Structure I is then

$$\begin{aligned}
 \text{RPSD} &= 2 |G'_1(e^{j\omega})|^2 + 2 |G'_2(e^{j\omega})|^2 + 3 \\
 &= 2 \frac{1}{\lambda^2} |G''_1(e^{j\omega})|^2 + 2 \frac{1}{\lambda^2 \sigma^2} |G''_2(e^{j\omega})|^2 + 3 \quad (106)
 \end{aligned}$$

where  $G'_1(e^{j\omega})$  and  $G'_2(e^{j\omega})$  are the noise transfer functions for the scaled filter,  $\lambda$  is the scaling factor, and  $G''_1(e^{j\omega})$  and  $G''_2(e^{j\omega})$  are functions generated from  $G'_1(e^{j\omega})$  and  $G'_2(e^{j\omega})$ , when we remove the parameters  $\sigma$  and  $\lambda$  from them.

## State-space sections without limit cycles

- We now show that choosing  $\sigma$  according to equation (102) leads to the minimization of the output noise.
- From equation (106), we can infer that the output noise can be minimized when  $\sigma$  and  $\lambda$  are maximized, with the scaling coefficient given by

$$\lambda = \frac{1}{\max \{ \|F_1(z)\|_p, \|F_2(z)\|_p \}} \quad (107)$$

- However,  $F_1(z)$  is not a function of  $\sigma$ , and, as a consequence, the choice of  $\|F_2(z)\|_p = \|F_1(z)\|_p$  leads to a maximum value for  $\lambda$ .
- On the other hand, the maximum  $\sigma$ , without reducing  $\lambda$ , is

$$\sigma = \frac{\|F_1(e^{j\omega})\|_p}{\|F_2''(e^{j\omega})\|_p} \quad (108)$$

from which we can conclude that this choice for  $\sigma$  minimizes the roundoff noise at the filter output.

## State-space sections without limit cycles

- In order to design a cascade structure without limit cycles which realizes

$$H(z) = \prod_{i=1}^m H_i(z) = H_0 \prod_{i=1}^m \frac{z^2 + \gamma'_{1i}z + \gamma'_{2i}}{z^2 + \alpha_{1i}z + \alpha_{2i}} = \prod_{i=1}^m (d_i + H'_i(z)) \quad (109)$$

with  $H'_i(z)$  described in the form of the first term of the right side of equation (45), one must adopt the following procedure for Structure I:

- Step 1: Calculate  $\sigma_i$  and  $\lambda_i$  for each section using

$$\sigma_i = \frac{\left\| F_{1i}(z) \prod_{j=1}^{i-1} H_j(z) \right\|_p}{\left\| F''_{2i}(z) \prod_{j=1}^{i-1} H_j(z) \right\|_p} \quad (110)$$

$$\lambda_i = \frac{1}{\left\| F_{2i}(z) \prod_{j=1}^{i-1} H_j(z) \right\|_p} \quad (111)$$

## State-space sections without limit cycles

- Step 2: Determine  $\alpha$  and  $\zeta$  from equations (94).
- Step 3: Compute the coefficients of **A**, **B**, and **C** using equations (91)–(93).
- Step 4: Calculate the multiplier coefficients  $d_i$  by

$$d_i = \begin{cases} \frac{1}{\left\| \prod_{j=1}^i H_j(z) \right\|_p}, & \text{for } i = 1, 2, \dots, (m-1) \\ \frac{H_0}{\prod_{j=1}^{m-1} d_j}, & \text{for } i = m \end{cases} \quad (112)$$

in order to satisfy overflow constraints at the output of each section.

## State-space sections without limit cycles

- Step 5: Incorporate the scaling multipliers of sections  $2, 3, \dots, m$  into the output multipliers of sections  $1, 2, \dots, (m - 1)$ , generating

$$\left. \begin{aligned} c'_{1i} &= c_{1i} \frac{\lambda_{i+1}}{\lambda_i} \\ c'_{2i} &= c_{2i} \frac{\lambda_{i+1}}{\lambda_i} \\ d'_i &= d_i \frac{\lambda_{i+1}}{\lambda_i} \end{aligned} \right\} \quad (113)$$

## State-space sections without limit cycles

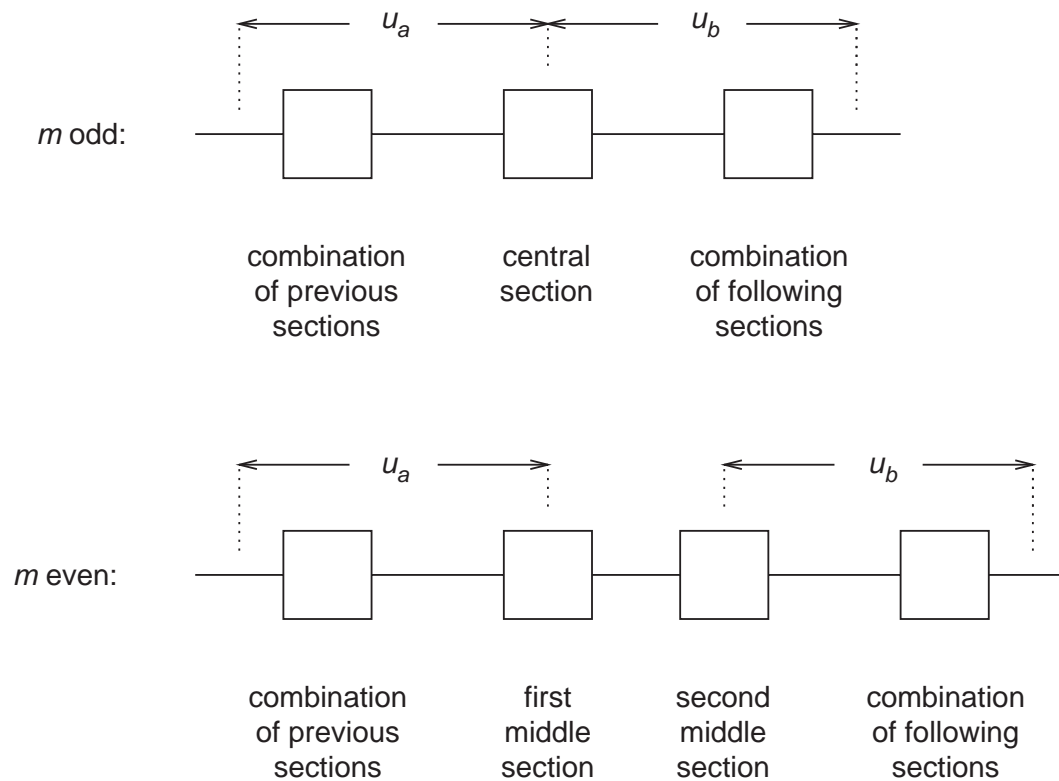


Figure 7: Ordering of state-space sections.



## State-space sections without limit cycles

- The cascade form design procedures employing the section-optimal and the limit-cycle-free state-space structures use the same strategy for pairing the poles and zeros as those employing direct-form sections.
- The section ordering depends on the definition of a parameter  $u_j$ , given by

$$u_j = \sum_{i=1}^2 \frac{\max \{|F_{ij}(e^{j\omega})|\}}{\min \{|F_{ij}(e^{j\omega})|\}} \quad (114)$$

where the maximum is computed for all  $\omega$ , while the minimum is calculated solely within the passband.

## State-space sections without limit cycles

- According to Figure 7, for an odd number of sections,  $m$ , the ordering consists of placing the section with highest value for  $u_j$  to be the central section. For an even number of sections, the two with largest  $u_j$  are placed in the central positions, these are called first and second middle sections. For odd  $m$ , the previous and following sections to the central block are chosen from the remaining sections so as to minimize the summation of  $u_a$  and  $u_b$  (see Figure 7), one referred to the combination of the central section and all previous ones, and the other referred to the combination of the central section and all the following ones. For even  $m$ , the sections before and after the central sections are chosen, among the remaining ones, in order to minimize the summation of  $u_a$  and  $u_b$  (see Figure 7), one referred to the combination of the first middle section and all previous sections, and the other referred to the combination of the second middle section and the sections following it. This approach is continuously employed until all second-order blocks have been ordered.

## State-space sections without limit cycles

- The output roundoff-noise PSD of the state-space structure without limit cycles in cascade form is expressed by

$$\Gamma_Y(e^{j\omega}) = \sigma_e^2 \sum_{i=1}^m \frac{1}{\lambda_{i+1}^2} (2G'_{1i}(e^{j\omega})G'_{1i}(e^{-j\omega}) + 2G'_{2i}(e^{j\omega})G'_{2i}(e^{-j\omega}) + 3) \quad (115)$$

where  $G'_{1i}(e^{j\omega})$  and  $G'_{2i}(e^{j\omega})$  are the noise frequency responses of the scaled sections and  $\lambda_{m+1} = 1$ .

### Example 13.3

- Repeat Example 13.1 using the cascade of optimal and limit-cycle-free state-space sections. Quantize the coefficients to 9 bits, including the sign bit, and verify the results.

### Example 13.3 - Solution

- In each case, all but the last section are scaled first to guarantee unit  $L_2$  norm at its output.
- After this initial scaling, for the optimal state-space structure, the transformation matrix  $\mathbf{T}$  is determined as given in equation (73) considering also the effect of the cumulative transfer function from previous blocks.
- Tables 7–10 list the coefficients of all the designed filters.
- Figure 8 depicts the magnitude responses obtained by the cascade of optimal state-space sections and of state-space sections without limit cycles. In all cases the coefficients were quantized to 9 bits including the sign bit.

### Example 13.3 - Solution

Table 7: Cascade structure using optimal state-space second-order sections.

Coefficient	Section 1	Section 2	Section 3
$a_{11}$	8.0271E—01	8.1339E—01	7.9823E—01
$a_{12}$	—5.9094E—01	—5.7910E—01	—6.0685E—01
$a_{21}$	5.7520E—01	5.7117E—01	5.8489E—01
$a_{22}$	8.0271E—01	8.1339E—01	7.9823E—01
$b_1$	8.0236E—02	6.4821E—03	2.9027E—02
$b_2$	1.5745E—01	—1.8603E—02	8.8313E—03
$c_1$	8.8747E—01	—8.8929E—01	2.5127E—02
$c_2$	4.5225E—01	3.0987E—01	8.2587E—02
$d$	8.8708E—02	1.2396E—01	1.3061E—02

### Example 13.3 - Solution

Table 8: Cascade structure using optimal state-space second-order sections quantized with 9 bits.

Coefficient	Section 1	Section 2	Section 3
$[a_{11}]_Q$	8.0078E-01	8.1250E-01	7.9688E-01
$[a_{12}]_Q$	-5.8984E-01	-5.7813E-01	-6.0547E-01
$[a_{21}]_Q$	5.7422E-01	5.7031E-01	5.8594E-01
$[a_{22}]_Q$	8.0078E-01	8.1250E-01	7.9688E-01
$[b_1]_Q$	8.2031E-02	7.8125E-03	2.7344E-02
$[b_2]_Q$	1.5625E-01	-1.9531E-02	7.8125E-03
$[c_1]_Q$	8.8672E-01	-8.9063E-01	2.3438E-02
$[c_2]_Q$	4.5313E-01	3.0859E-01	8.2031E-02
$[d]_Q$	8.9844E-02	1.2500E-01	1.1719E-02

### Example 13.3 - Solution

Table 9: Cascade structure without limit cycles using optimal state-space second-order sections.  $\lambda = 2.7202\text{E}-01$ .

Coefficient	Section 1	Section 2	Section 3
$a_{11}$	8.0272E-01	8.1339E-01	7.9822E-01
$a_{12}$	-5.8289E-01	-5.7823E-01	-5.8486E-01
$a_{21}$	5.8316E-01	5.7204E-01	6.0688E-01
$a_{22}$	8.0272E-01	8.1339E-01	7.9822E-01
$b_1$	1.9728E-01	1.8661E-01	2.0178E-01
$b_2$	-5.8316E-01	-5.7204E-01	-6.0688E-01
$c_1$	-9.2516E-03	-4.4228E-02	9.1891E-02
$c_2$	-2.8600E-02	1.6281E-02	-2.5557E-02
$d$	9.2516E-03	1.8323E-01	2.9191E-01

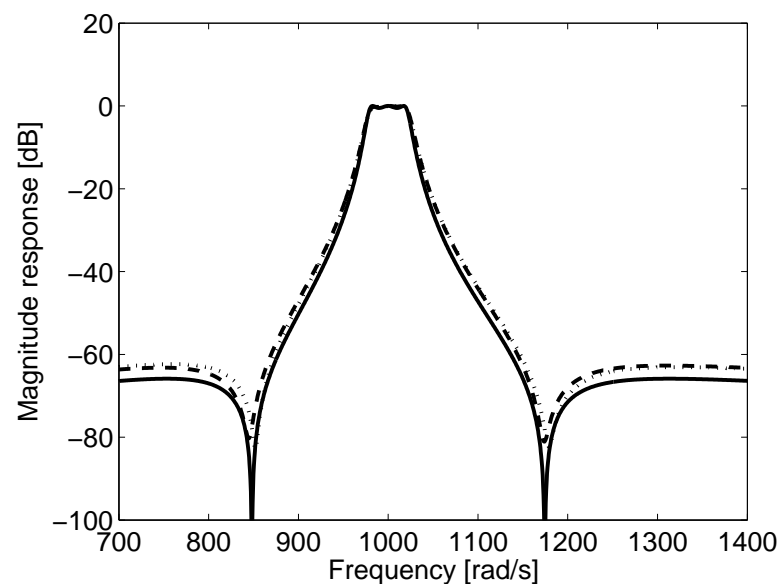


### Example 13.3 - Solution

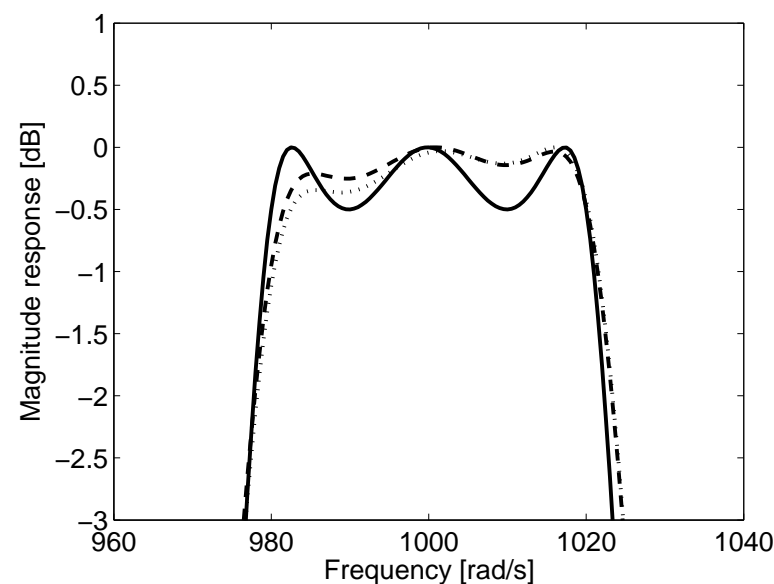
Table 10: Cascade structure without limit cycles using optimal state-space second-order sections quantized with 9 bits.  $[\lambda]_Q = 2.7344\text{E}-01$ .

Coefficient	Section 1	Section 2	Section 3
$[a_{11}]_Q$	8.0078E-01	8.1250E-01	7.9688E-01
$[a_{12}]_Q$	-5.8203E-01	-5.7812E-01	-5.8594E-01
$[a_{21}]_Q$	5.8203E-01	5.7031E-01	6.0547E-01
$[a_{22}]_Q$	8.0078E-01	8.1250E-01	7.9688E-01
$[b_1]_Q$	1.9922E-01	1.8750E-01	2.0313E-01
$[b_2]_Q$	-5.8203E-01	-5.7031E-01	-6.0547E-01
$[c_1]_Q$	-7.8125E-03	-4.2969E-02	9.3750E-02
$[c_2]_Q$	-2.7344E-02	1.5625E-02	-2.7344E-02
$[d]_Q$	7.8125E-03	1.8359E-01	2.9297E-01

## Example 13.3 - Solution



(a)



(b)

Figure 8: Coefficient-quantization effects in the cascade forms, using state-space second-order sections: (a) overall magnitude response; (b) passband detail. (Solid line – initial design; dashed line – cascade of optimal state-space sections (9 bits); dotted line – cascade of limit-cycle-free state-space sections (9 bits).)

## Lattice filters

- Consider a general IIR transfer function written in the form

$$H(z) = \frac{N_M(z)}{D_N(z)} = \frac{\sum_{i=0}^M b_{i,M} z^{-i}}{1 + \sum_{i=1}^N a_{i,N} z^{-i}} \quad (116)$$

- In the lattice construction, we concentrate first on the realization of the denominator polynomial through an order-reduction strategy. For that, we define the auxiliary  $N$ th-order polynomial, obtained by reversing the order of the coefficients of the denominator  $D_N(z)$ , as given by

$$zB_N(z) = D_N(z^{-1})z^{-N} = z^{-N} + \sum_{i=1}^N a_{i,N} z^{i-N} \quad (117)$$

## Lattice filters

- We can then calculate a reduced order polynomial as

$$\begin{aligned}
 (1 - a_{N,N}^2) D_{N-1}(z) &= D_N(z) - a_{N,N} z B_N(z) \\
 &= (1 - a_{N,N}^2) + \cdots + (a_{N-1,N} - a_{N,N} a_{1,N}) z^{-N+1}
 \end{aligned}
 \tag{118}$$

where we can also express  $D_{N-1}(z)$  as  $1 + \sum_{i=1}^{N-1} a_{i,N-1} z^{-i}$ .

- Note that the first and last coefficients of  $D_N(z)$  are 1 and  $a_{N,N}$ , whereas for the polynomial  $zB_N(z)$  they are  $a_{N,N}$  and 1, respectively.
- This strategy to achieve the order reduction guarantees a monic  $D_{N-1}(z)$ , that is,  $D_{N-1}(z)$  having the coefficient of  $z^0$  equal to 1.

## Lattice filters

- By induction, this order-reduction procedure can be performed repeatedly, thus yielding

$$zB_j(z) = D_j(z^{-1})z^{-j} \quad (119)$$

$$D_{j-1}(z) = \frac{1}{1 - a_{j,j}^2} (D_j(z) - a_{j,j}zB_j(z)) \quad (120)$$

for  $j = N, (N - 1), \dots, 1$ , with  $zB_0(z) = D_0(z) = 1$ .

- It can be shown that the above equations are equivalent to the following expression:

$$\begin{bmatrix} D_{j-1}(z) \\ B_j(z) \end{bmatrix} = \begin{bmatrix} 1 & -a_{j,j} \\ a_{j,j}z^{-1} & (1 - a_{j,j}^2)z^{-1} \end{bmatrix} \begin{bmatrix} D_j(z) \\ B_{j-1}(z) \end{bmatrix} \quad (121)$$

## Lattice filters

- The previous equation can be implemented, for example, by the two-port network  $TP_j$  shown in Figure 9.

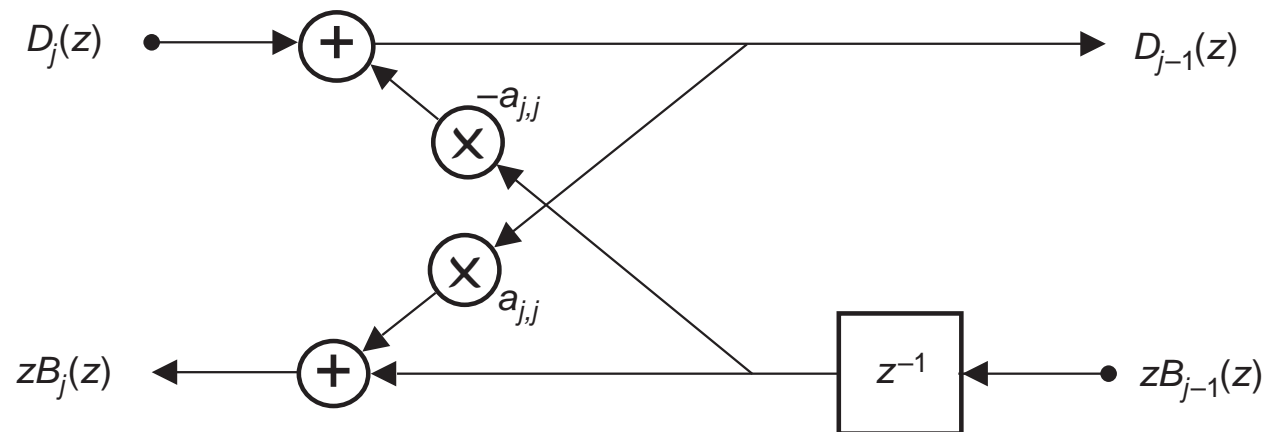


Figure 9: Two-multiplier network  $TP_j$  implementing equation (121).



## Lattice filters

- Since at the lower branches of the two-port networks in Figure 10 we have the signals  $\frac{zB_j(z)}{D_N(z)}$  available, then a convenient way to form the desired numerator is to apply weights to the polynomials  $zB_j(z)$ , such that

$$N_M(z) = \sum_{j=0}^M v_j zB_j(z) \quad (122)$$

where the tap coefficients  $v_j$  are calculated through the following order-reduction recursion

$$N_{j-1}(z) = N_j(z) - zv_j B_j(z) \quad (123)$$

for  $j = M, (M-1), \dots, 1$ , with  $v_M = b_{M,M}$  and  $v_0 = b_{0,0}$ .

- Then, a way of implementing the overall IIR transfer function  $H(z) = \frac{N(z)}{D(z)} = \frac{\sum_{j=0}^M v_j zB_j(z)}{D_N(z)}$  is to use the structure in Figure 11, which is called the IIR lattice realization.



## Lattice filters

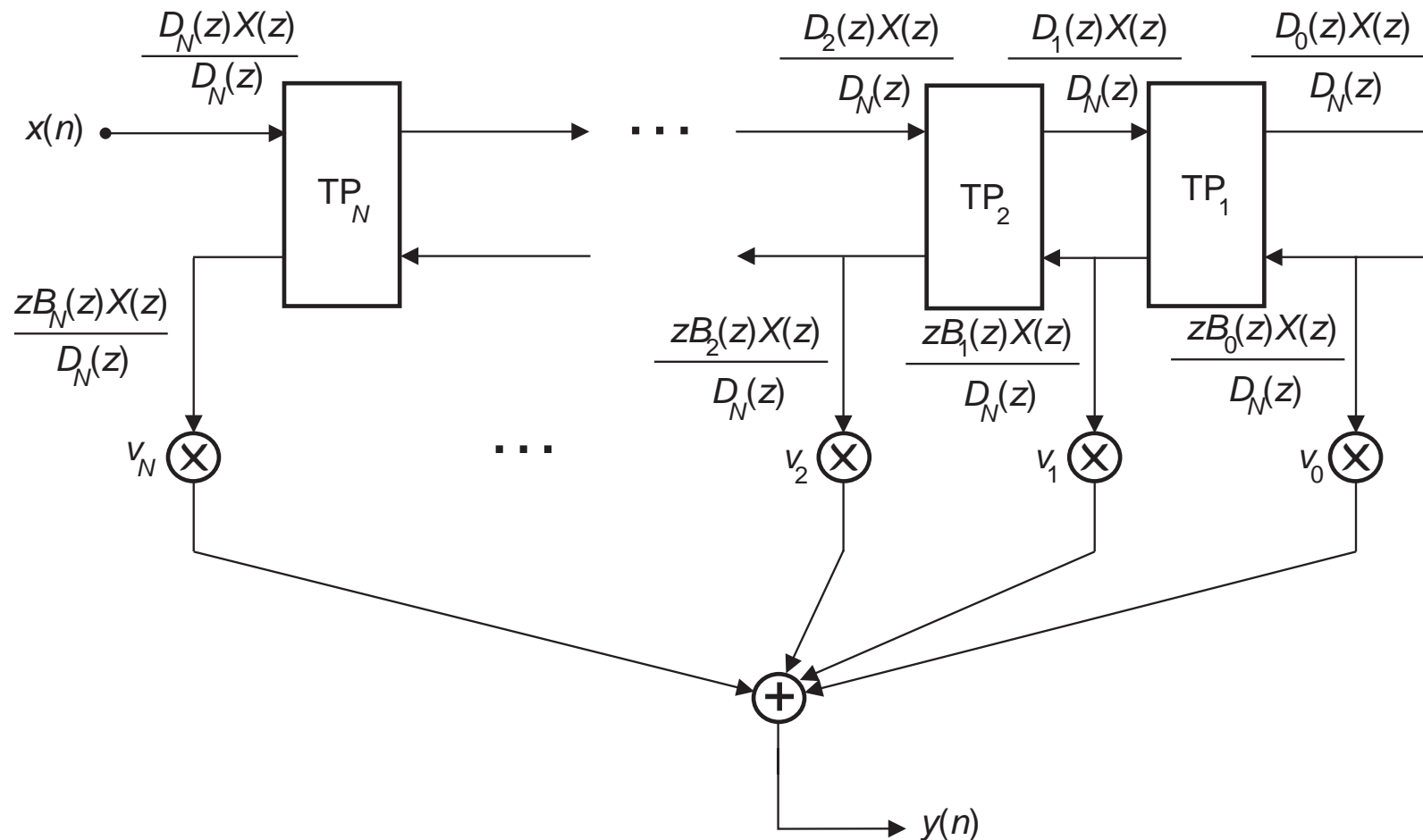


Figure 11: General IIR lattice digital filter structure.

## Lattice filters

- From the above, we have a simple procedure for obtaining the lattice network given the direct-form transfer function,  $H(z) = \frac{N_M(z)}{D_N(z)}$ :
  - (i) Obtain, recursively, the polynomials  $B_j(z)$  and  $D_j(z)$ , as well as the lattice coefficient  $\alpha_{j,j}$ , for  $j = N, (N - 1), \dots, 1$ , using equations (119) and (120).
  - (ii) Compute the coefficients  $v_j$ , for  $j = N, (N - 1), \dots, 1$ , using the recursion in equation (123).

## Lattice filters

- Conversely, if given the lattice realization we want to compute the direct-form transfer function, we can use the following procedure:
  - (i) Start with  $zB_0(z) = D_0(z)=1$ .
  - (ii) Compute recursively  $B_j(z)$  and  $D_j(z)$  for  $j = 1, 2, \dots, N$ , using the following relation:

$$\begin{bmatrix} D_j(z) \\ B_j(z) \end{bmatrix} = \begin{bmatrix} 1 & a_{j,j} \\ a_{j,j}z^{-1} & z^{-1} \end{bmatrix} \begin{bmatrix} D_{j-1}(z) \\ B_{j-1}(z) \end{bmatrix} \quad (124)$$

- (iii) Compute  $N_M(z)$  using equation (122).
- (iv) The direct-form transfer function is then  $H(z) = \frac{N_M(z)}{D_N(z)}$ .

## Lattice filters

- There are some important properties related to the lattice realization which should be mentioned.
- If  $D_N(z)$  has all the roots inside the unit circle the lattice structure will have all coefficients  $a_{j,j}$  with magnitude less than one. Otherwise,  $H(z) = \frac{N(z)}{D(z)}$  represents an unstable system. This straightforward stability condition makes the lattice realizations useful for implementing time-varying filters.
- In addition, the polynomials  $zB_j(z)$ , for  $j = 0, 1, \dots, M$ , form an orthogonal set. This property justifies the choice of these polynomials to form the desired numerator polynomial  $N_M(z)$ , as described in equation (123).

## Lattice filters

- Since, in Figure 10, the two-port system consisting of section  $TP_j$  and all the sections to its right, which relates the output signal  $\frac{zB_j(z)X(z)}{D_N(z)}$  to the input signal  $\frac{D_j(z)X(z)}{D_N(z)}$  is linear, its transfer function remains unchanged if we multiply its input signal by  $\lambda_j$  and divide its output by the same amount
- Therefore,  $\frac{zB_N(z)}{D_N(z)}$  will not change if we multiply the signal entering the upper-left branch of section  $TP_j$  by  $\lambda_j$  and divide the signal leaving the lower-left branch by  $\lambda_j$ .
- This is equivalent to scaling the section  $TP_j$  by  $\lambda_j$ .

## Lattice filters

- If we do this for every branch  $j$ , the signals entering and leaving at the left of section  $N$  remain unchanged, the signals entering and leaving at the left of section  $(N - 1)$  will be scaled by  $\lambda_N$ , the signals entering and leaving at the left of section  $(N - 2)$  will be multiplied by  $\lambda_N \lambda_{N-1}$ , and so on, leading to the scaled signals  $\frac{\overline{D}_j(z)X(z)}{D_N(z)}$  and  $\frac{z\overline{B}_j(z)X(z)}{D_N(z)}$  at the left of section  $TP_j$ , such that

$$\overline{D}_j(z) = \left( \prod_{i=N}^{j+1} \lambda_i \right) D_j(z) \quad (125)$$

$$\overline{B}_j(z) = \left( \prod_{i=N}^{j+1} \lambda_i \right) B_j(z) \quad (126)$$

for  $j = (N - 1), (N - 2), \dots, 1$ , with  $\overline{D}_N(z) = D_N(z)$  and  $\overline{B}_N(z) = B_N(z)$ .

## Lattice filters

- Therefore, in order to maintain the transfer function of the scaled lattice realization unchanged, we must make

$$\bar{v}_j = \frac{v_j}{\left( \prod_{i=N}^{j+1} \lambda_i \right)} \quad (127)$$

for  $j = (N - 1), (N - 2), \dots, 1$ , with  $\bar{v}_N = v_N$ .

- Based on the above property, we can derive a more economical two-port network using a single multiplier, as shown in Figure 12, where the plus-or-minus signs indicate that two different realizations are possible.

## Lattice filters

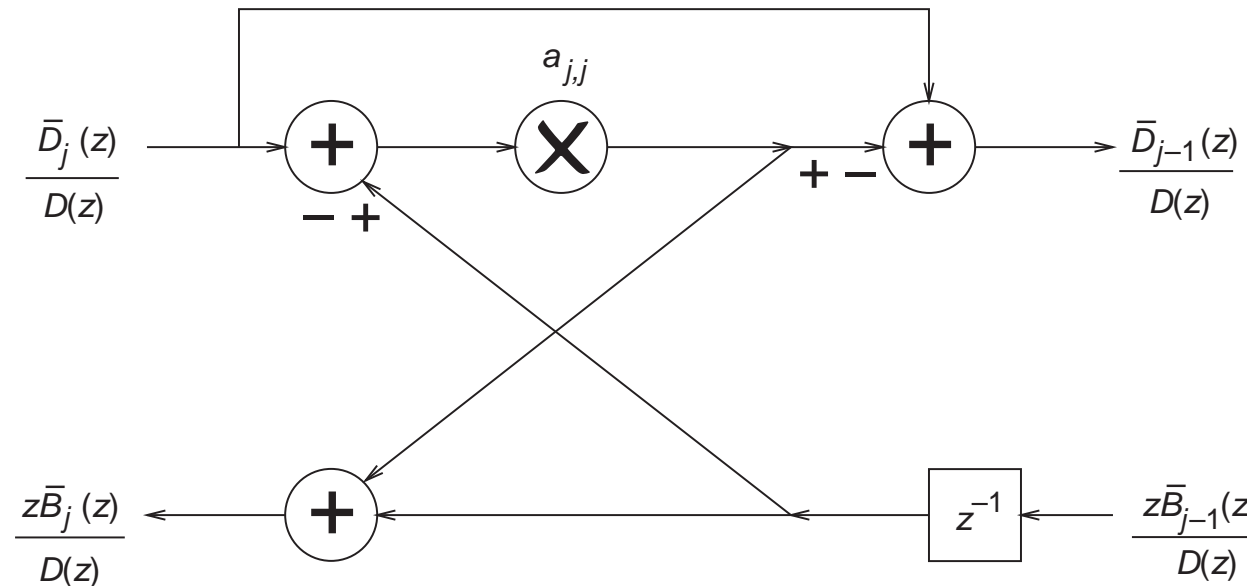


Figure 12: The one-multiplier network for equation (121).

- The choice of these signs can vary from section to section, aiming at the reduction of the quantization noise at the filter output.
- This network is equivalent to the one in Figure 9 scaled using  $\lambda_j = 1 \pm a_{j,j}$ , and therefore the coefficients  $\bar{v}_j$  should be computed using equation (127).



## Lattice filters

- Another important realization for the two-port network results when the scaling parameters  $\lambda_i$  are chosen such that all the internal nodes of the lattice network have a transfer function from the input of unit  $L_2$  norm. The appropriate scaling can be derived by first noting that at the left of section  $TP_j$ , the norms of the corresponding transfer functions are given by

$$\left\| \frac{\overline{D}_j(z)}{\overline{D}_N(z)} \right\|_2 = \left\| \frac{z\overline{B}_j(z)}{\overline{D}_N(z)} \right\|_2 \quad (128)$$

since, from equation (119),  $z\overline{B}_j(z) = D_j(z^{-1})z^{-j}$ .

- From the above equations, if we want unit  $L_2$  norm at the internal nodes of the lattice network, we must have

$$\left\| \frac{z\overline{B}_0(z)}{\overline{D}_N(z)} \right\|_2 = \dots = \left\| \frac{z\overline{B}_{N-1}(z)}{\overline{D}_N(z)} \right\|_2 = \left\| \frac{z\overline{B}_N(z)}{\overline{D}_N(z)} \right\|_2 = \left\| \frac{\overline{D}_N(z)}{\overline{D}_N(z)} \right\|_2 = 1 \quad (129)$$

## Lattice filters

- Then, using  $\lambda_j$  from equations (124) to (126), it can be derived that

$$\lambda_j = \frac{\left\| \frac{zB_j(z)}{D_N(z)} \right\|_2}{\left\| \frac{zB_{j-1}(z)}{D_N(z)} \right\|_2} = \sqrt{1 - a_{j,j}^2} \quad (130)$$

- It is easy to show that section  $TP_j$  of the normalized lattice can be implemented as depicted in Figure 13. The most important feature of the normalized lattice realization is that, since all its internal nodes have transfer function with unit  $L_2$  norm, it presents an automatic scaling in the  $L_2$ -norm sense. This explains the low roundoff noise generated by the normalized lattice realization as compared with the other forms of the lattice realization. Note that the coefficients  $\bar{v}_j$  have to be computed using equation (127).

## Lattice filters

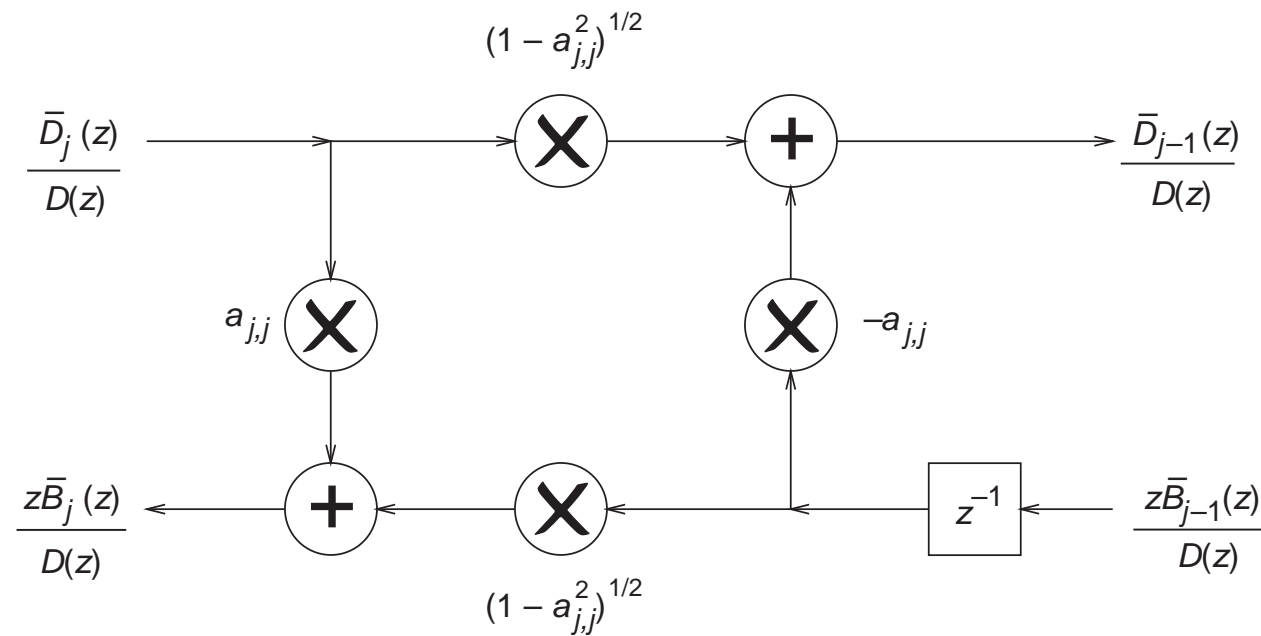


Figure 13: The normalized network for equation (121).

### Example 13.4

- Repeat Example 13.1 using the one-multiplier, two-multiplier, and normalized lattice forms. Quantize the coefficients of the normalized lattice using 9 bits, including the sign bit, and verify the results.

### Example 13.4 - Solution

- The two-multiplier IIR lattice can be determined from the direct form using the MATLAB command `tf2latc`. For the one-multiplier, we use  $\lambda_j = (1 + a_{j,j})$  in equation (127) to determine the feedforward coefficients, whereas for the normalized lattice, we use  $\lambda_j = \sqrt{1 - a_{j,j}^2}$ . The resulting coefficients in each case are seen in Tables 11–14.

## Example 13.4 - Solution

Table 11: Coefficients of two-multiplier lattice.

Section $j$	$a_{j,j}$	$v_j$
0		$-2.1521\text{E}-06$
1	$8.0938\text{E}-01$	$-1.1879\text{E}-06$
2	$-9.9982\text{E}-01$	$9.3821\text{E}-06$
3	$8.0903\text{E}-01$	$3.4010\text{E}-06$
4	$-9.9970\text{E}-01$	$8.8721\text{E}-05$
5	$8.0884\text{E}-01$	$-2.3326\text{E}-04$
6	$-9.6906\text{E}-01$	$-1.4362\text{E}-04$

## Example 13.4 - Solution

Table 12: Coefficients of one-multiplier lattice.

Section $j$	$a_{j,j}$	$\bar{v}_j$
0		$-2.1371E+02$
1	$8.0938E-01$	$-2.1342E+02$
2	$-9.9982E-01$	$3.0663E-01$
3	$8.0903E-01$	$2.0108E-01$
4	$-9.9970E-01$	$1.5850E-03$
5	$8.0884E-01$	$-7.5376E-03$
6	$-9.6905E-01$	$-1.4362E-04$

## Example 13.4 - Solution

Table 13: Coefficients of normalized lattice.

Section $j$	$a_{j,j}$	$\bar{v}_j$
0		$-9.1614\text{E}-02$
1	$8.0938\text{E}-01$	$-2.9697\text{E}-02$
2	$-9.9982\text{E}-01$	$4.4737\text{E}-03$
3	$8.0903\text{E}-01$	$9.5319\text{E}-04$
4	$-9.9970\text{E}-01$	$6.1121\text{E}-04$
5	$8.0884\text{E}-01$	$-9.4494\text{E}-04$
6	$-9.6905\text{E}-01$	$-1.4362\text{E}-04$



## Example 13.4 - Solution

Table 14: Coefficients of normalized lattice quantized with 9 bits.

Section $j$	$a_{j,j}$	$\bar{v}_j$
0		$-8.9844\text{E} - 02$
1	$8.0938\text{E} - 01$	$-3.1250\text{E} - 02$
2	$-9.9982\text{E} - 01$	$3.9063\text{E} - 03$
3	$8.0903\text{E} - 01$	$0.0000\text{E} + 00$
4	$-9.9970\text{E} - 01$	$0.0000\text{E} + 00$
5	$8.0884\text{E} - 01$	$0.0000\text{E} + 00$
6	$-9.6905\text{E} - 01$	$0.0000\text{E} + 00$

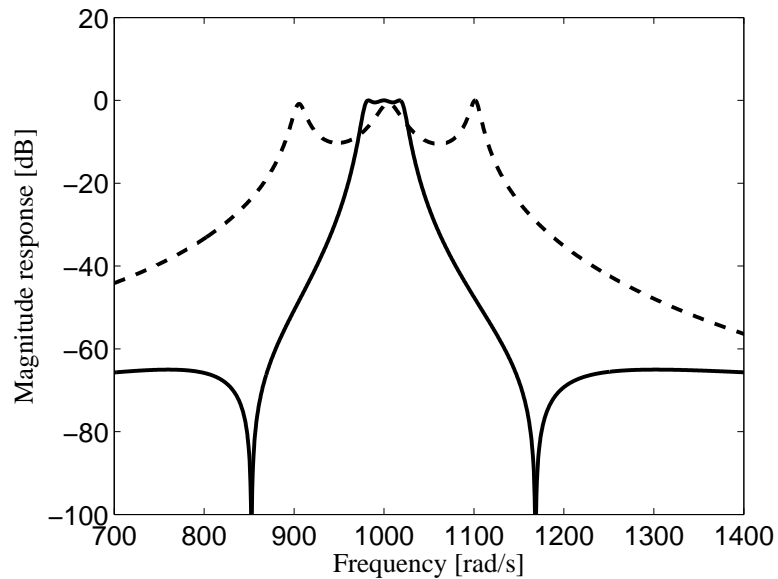
### Example 13.4 - Solution

- In the quantization procedure, we must guarantee that the absolute value of all feedback coefficients  $a_{j,j}$  remain below 1 to guarantee stability of the resulting filter.
- From Tables 11–14, one observes that the three lattice forms have serious quantization issues due to the wide range covered by their coefficients.
- It must be added that the normalized lattice performs much better than the two- and one-multiplier lattices with respect to quantization effects.
- It is also worth mentioning that in the two-multiplier structure, the feedforward coefficients assume very small values, forcing the use of more than 9 bits for their representation. This normally happens when designing a filter having poles very close to the unit circle, as is the case in this example.

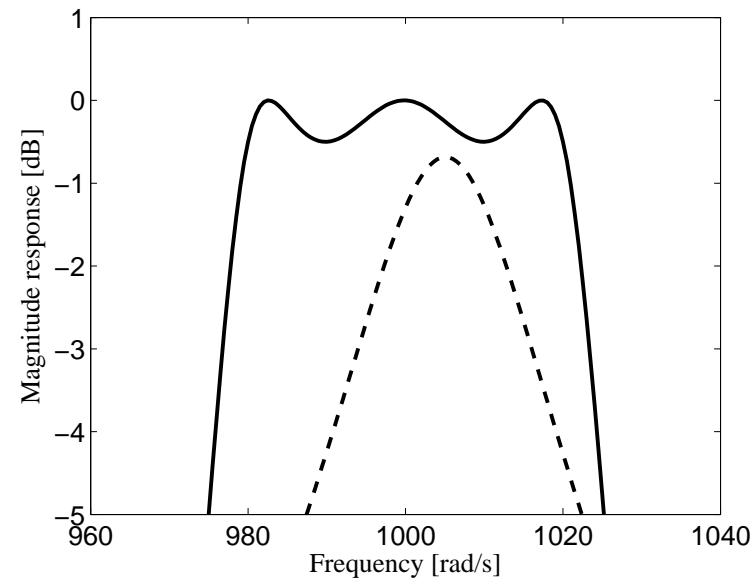
## Example 13.4 - Solution

- Figure 14 depicts the magnitude responses obtained by the original and quantized normalized lattices. Note that the magnitude responses of the normalized lattice structure are significantly different to the ideal one, especially when compared to the results shown in previous examples.

## Example 13.4 - Solution



(a)



(b)

Figure 14: Coefficient-quantization effects in the normalized lattice form: (a) overall magnitude response; (b) passband detail. (Solid line – initial design; dashed line – normalized lattice (9 bits).)

## Doubly complementary filters

- In this section the class of doubly complementary filter is discussed as it plays an important role in alias-free 2-band filter banks and some audio applications.
- **Theorem:** Two transfer functions  $H_0(z)$  and  $H_1(z)$  are referred to as doubly complementary if their frequency responses are allpass complementary, that is,

$$|H_0(e^{j\omega}) + H_1(e^{j\omega})|^2 = 1 \quad (131)$$

and also power complementary, such that

$$|H_0(e^{j\omega})|^2 + |H_1(e^{j\omega})|^2 = 1 \quad (132)$$

for all  $\omega$ .

## Doubly complementary filters

- For doubly complementary filters, we can write that

$$H_0(z) + H_1(z) = F_0(z) \quad (133)$$

$$H_0(z) - H_1(z) = F_1(z) \quad (134)$$

where  $F_0(z)$  and  $F_1(z)$  are stable allpass transfer functions, and then

$$H_0(z) = \frac{1}{2} (F_0(z) + F_1(z)) \quad (135)$$

$$H_1(z) = \frac{1}{2} (F_0(z) - F_1(z)) \quad (136)$$

whose implementation can be as shown in Figure 15.

## Doubly complementary filters

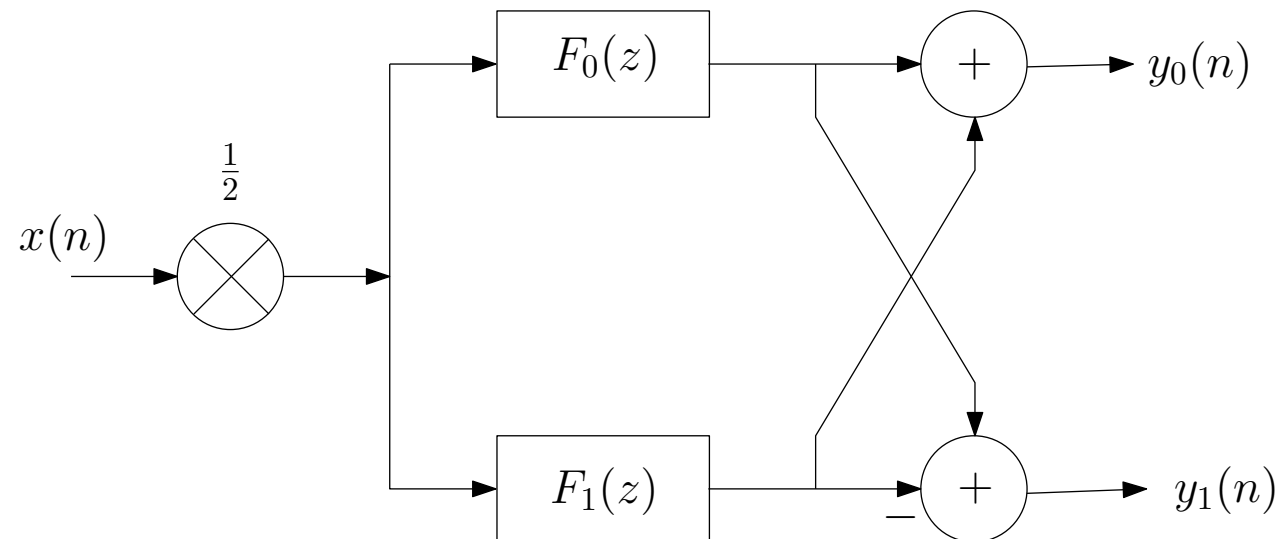


Figure 15: Doubly complementary filters.

## Doubly complementary filters

- **Proof:** The doubly complementary frequency responses can be described in polar form as follows:

$$H_0(e^{j\omega}) = r_0(\omega)e^{j\phi_0(\omega)} \quad (137)$$

$$H_1(e^{j\omega}) = r_1(\omega)e^{j\phi_1(\omega)} \quad (138)$$

- Using these expressions, the left-hand side  $L$  of equation (131) can be written as

$$\begin{aligned}
 L &= |r_0(\omega)e^{j\phi_0(\omega)} + r_1(\omega)e^{j\phi_1(\omega)}|^2 \\
 &= \left(r_0(\omega)e^{j\phi_0(\omega)} + r_1(\omega)e^{j\phi_1(\omega)}\right) \left(r_0(\omega)e^{-j\phi_0(\omega)} + r_1(\omega)e^{-j\phi_1(\omega)}\right) \\
 &= r_0^2(\omega) + r_1^2(\omega) + r_0(\omega)r_1(\omega)e^{j(\phi_0(\omega)-\phi_1(\omega))} + r_0(\omega)r_1(\omega)e^{-j(\phi_0(\omega)-\phi_1(\omega))} \\
 &= r_0^2(\omega) + r_1^2(\omega) + 2r_0(\omega)r_1(\omega)\cos(\phi_0(\omega) - \phi_1(\omega)) \quad (139)
 \end{aligned}$$



## Doubly complementary filters

- Since equation (132) is equivalent to  $r_0^2(\omega) + r_1^2(\omega) = 1$ , then for allpass complementary  $H_0(e^{j\omega})$  and  $H_1(e^{j\omega})$  we must have that

$$2r_0(\omega)r_1(\omega)\cos(\phi_0(\omega) - \phi_1(\omega)) = 0 \quad (140)$$

- By following the same procedure as the derivation of equation (139), it is possible to show that

$$\begin{aligned} |H_0(e^{j\omega}) - H_1(e^{j\omega})|^2 &= r_0^2(\omega) + r_1^2(\omega) - 2r_0(\omega)r_1(\omega)\cos(\phi_0(\omega) - \phi_1(\omega)) \\ &= 1 \end{aligned} \quad (141)$$

## Doubly complementary filters

- By applying the expressions of equations (135) and (136) in equation (132) and using the polar representation, it is straightforward to show that

$$|F_0(e^{j\omega})|^2 + |F_1(e^{j\omega})|^2 = 2 \quad (142)$$

Also applying equations (135) and (136) along with (140) in equation (131) it follows that

$$|F_0(e^{j\omega})|^2 = r_0^2(\omega) + r_1^2(\omega) = 1 \quad (143)$$

and then

$$|F_1(e^{j\omega})|^2 = r_0^2(\omega) + r_1^2(\omega) = 1 \quad (144)$$

Therefore,  $F_0(z)$  and  $F_1(z)$  are both allpass filters.

## Doubly complementary filters

- Allpass transfer functions have the following general form

$$F_i(z) = \frac{\sum_{l=0}^{N_i} a_{N_i-l,i} z^{-l}}{\sum_{l=0}^{N_i} a_{l,i} z^{-l}} = \frac{D_i(z^{-1})}{z^{-N_i} D_i(z)} = z^{N_i} \frac{D_i(z^{-1})}{D_i(z)} \quad (145)$$

for  $i = 0, 1$  and  $D_i(z) = a_{0,i} z^{N_i} + a_{1,i} z^{N_i-1} + \dots + a_{N_i,i}$ .

- The phase responses of the allpass filters are given by

$$\theta_i(\omega) = -N_i \omega + 2 \arctan \left( \frac{\sum_{l=0}^{N_i} a_{l,i} \sin(l\omega)}{\sum_{l=0}^{N_i} a_{l,i} \cos(l\omega)} \right) \quad (146)$$

## Doubly complementary filters

- Given that  $F_0(e^{j\omega})$  and  $F_1(e^{j\omega})$  are allpass frequency responses they can be expressed in polar form as

$$F_0(e^{j\omega}) = e^{j\theta_0(\omega)} \quad (147)$$

$$F_1(e^{j\omega}) = e^{j\theta_1(\omega)} \quad (148)$$

in such way that

$$|H_0(e^{j\omega})| = \frac{1}{2} \left| e^{j(\theta_0(\omega) - \theta_1(\omega))} + 1 \right| \quad (149)$$

$$|H_1(e^{j\omega})| = \frac{1}{2} \left| e^{j(\theta_0(\omega) - \theta_1(\omega))} - 1 \right| \quad (150)$$

## Doubly complementary filters

- Assuming that at frequency  $\omega = 0$  both allpass filters have zero phase, as a result  $|H_0(1)| = 1$  and  $|H_1(1)| = 0$ , which are typical features of lowpass and highpass filters, respectively. On the other hand, for  $\omega = \pi$

$$|H_0(e^{j\pi})| = \frac{1}{2} |e^{j(N_0 - N_1)\pi} + 1| \quad (151)$$

$$|H_1(e^{j\pi})| = \frac{1}{2} |e^{j(N_0 - N_1)\pi} - 1| \quad (152)$$

so that if the difference  $(N_0 - N_1)$  is odd, then  $|H_0(e^{j\pi})| = 0$  and  $|H_1(e^{j\pi})| = 1$ , again a typical property of lowpass and highpass filters, respectively.

## Doubly complementary filters

- Let us consider a simple, and yet useful, choice for the allpass transfer functions, that is

$$F_0(z) = z^{-N_0} \quad (153)$$

$$F_1(z) = z^{-1} F'_1(z) \quad (154)$$

where  $F'_1(z)$  is a standard allpass transfer function of order  $N_0$  of the form in equation (145).

- With an odd  $(N_0 - N_1) = 1$ , it is possible to generate doubly complementary transfer functions with lowpass and high shapes given by

$$H_0(z) = z^{-N_0} + z^{-1} F'_1(z) \quad (155)$$

$$H_1(z) = z^{-N_0} - z^{-1} F'_1(z) \quad (156)$$

## Doubly complementary filters

- The difference in the phase response of the allpass filters are given by

$$\begin{aligned}
 \theta_0(\omega) - \theta_1(\omega) &= -N_0\omega + \theta_1(\omega) \\
 &= (-N_0 - 1)\omega + \angle F'_1(e^{j\omega}) \\
 &= (-N_0 - 1)\omega + N_0\omega - 2\arctan \left( \frac{\sum_{l=0}^{N_0} a_{l,0} \sin(l\omega)}{\sum_{l=0}^{N_0} a_{l,0} \cos(l\omega)} \right) \\
 &= -\omega - 2\arctan \left( \frac{\sum_{l=0}^{N_0} a_{l,0} \sin(l\omega)}{\sum_{l=0}^{N_0} a_{l,0} \cos(l\omega)} \right) \quad (157)
 \end{aligned}$$

### Example 13.5

- Design doubly complementary filters  $H_0(z)$  and  $H_1(z)$  satisfying the specification for the lowpass filter below:

$$\left. \begin{aligned} A_r &= 40 \text{ dB} \\ \Omega_p &= 0.5\pi \text{ rad/s} \\ \Omega_r &= 0.6\pi \text{ rad/s} \end{aligned} \right\} \quad (158)$$



## Example 13.5 - Solution

- In this solution we employ the simple choice for the first allpass filter  $F_0(z) = z^{-N_0}$ . In this case, we start the solution by first designing an allpass filter  $F_1(z)$  whose phase response follows as close as possible the following specifications

$$\theta_1(\omega) = \begin{cases} -N_0\omega, & \text{for } 0 \leq \omega \leq \Omega_p \\ -N_0\omega + \pi, & \text{for } \Omega_r \leq \omega \leq \pi \end{cases} \quad (159)$$

considering the sampling frequency  $\Omega_s = 2\pi$ .

- With this strategy the phase difference  $\theta_0(\omega) - \theta_1(\omega) = -N_0\omega - \theta_1(\omega)$  in equation (157) will be approximately zero at low frequencies and approximately  $\pi$  after the frequency  $\frac{\pi}{2}$ , thus enforcing the doubly complementary property.

### Example 13.5 - Solution

- There are several ways to design allpass filters satisfying prescribed phase specifications such as those based on the  $L_p$ -norm minimization criteria described in Section 6.5. Other specialized methods are described in the associated literature. We employed them to design an allpass filter  $F_1(z)$  whose coefficients are shown in Table 15. A sixth-order allpass filter sufficed to generate a stopband attenuation of about 40 dB.

## Example 13.5 - Solution

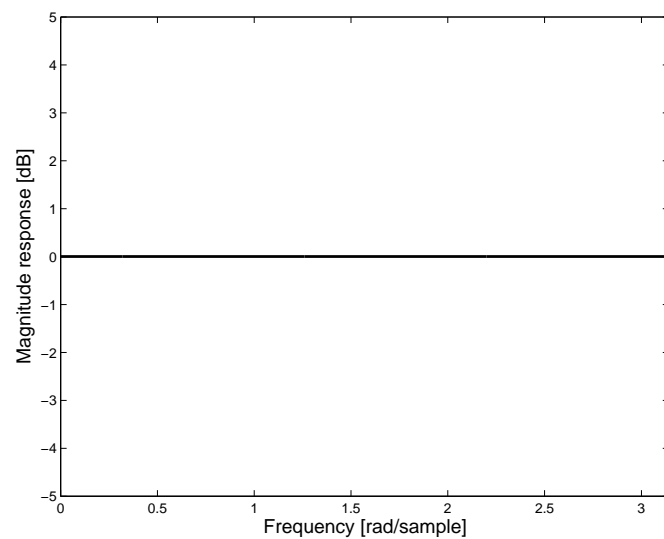
Table 15:  $F_1(z)$  Allpass filter coefficients.

Coefficient $a_{j,1}$	
$a_{0,1} =$	1.0000
$a_{1,1} =$	0.0000
$a_{2,1} =$	0.4780
$a_{3,1} =$	0.0000
$a_{4,1} =$	-0.0941
$a_{5,1} =$	0.0000
$a_{6,1} =$	0.0283

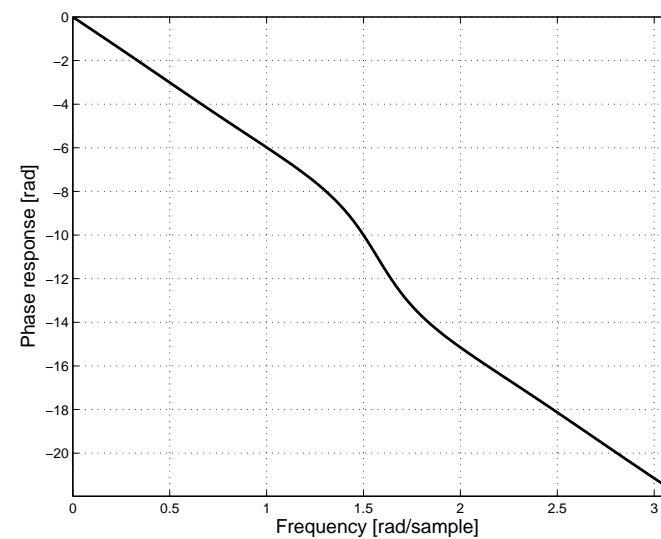
### Example 13.5 - Solution

- As can be observed, the even-order coefficients of the allpass filter are zero, a property originated from the fact that the allpass filter has symmetric response around the frequency  $\frac{\pi}{2}$  as any half-band filter.
- Figure 16 depicts the magnitude and phase responses of the allpass filter  $F'_1(z)$ , where from the phase response it is possible to observe the differences in the phase delays at the low and high frequency ranges.
- Figure 17 depicts the magnitude and phase responses of the doubly complementary filters  $H_0(z)$  and  $H_1(z)$ , respectively, generated according to equations (155) and (156).

## Example 13.5 - Solution



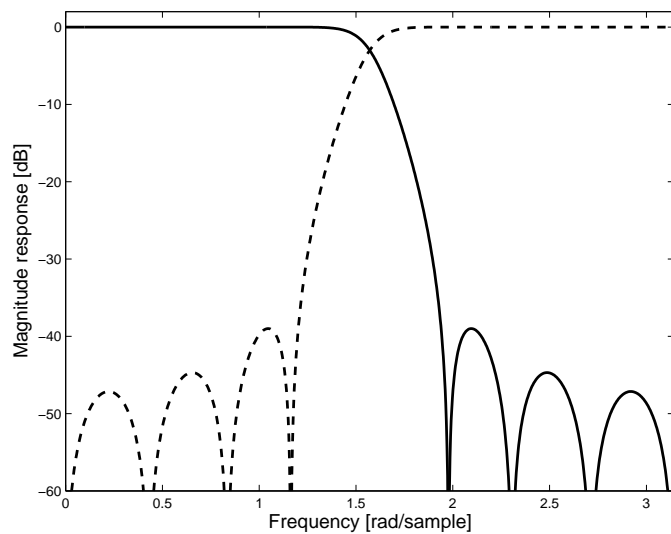
(a)



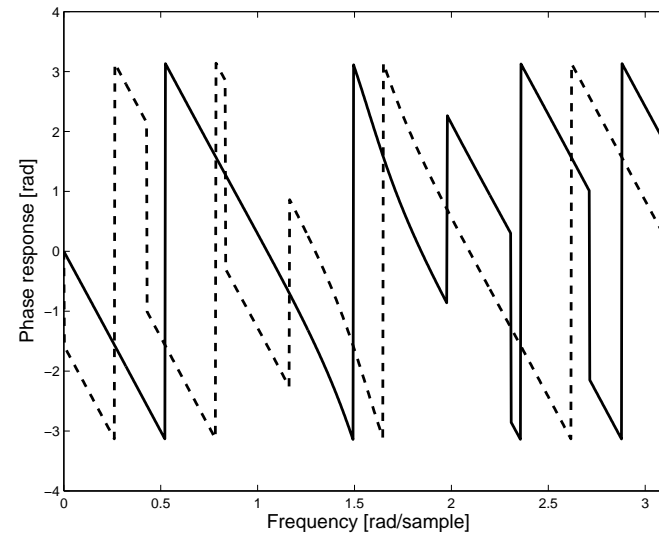
(b)

Figure 16: Allpass filter of order  $N = 6$ : (a) magnitude response of the allpass filter; (b) unwrapped phase response.

## Example 13.5 - Solution



(a)



(b)

Figure 17: Doubly complementary filter of order  $N = 7$ : (a) magnitude responses of the lowpass and highpass filters; (b) phase responses.

## QMF filter bank implementation

- We consider the case where  $H_0(z)$  and  $H_1(z)$  satisfy the doubly complementary conditions, and they are chosen as the lowpass and highpass analysis filters from a filter bank, respectively.
- The synthesis filters are selected according to the QMF conditions, namely  $G_0(z) = H_1(-z)$  and  $G_1(z) = -H_0(-z)$ , and the overall transfer function of the 2-band QMF filter bank is given by

$$H(z) = \frac{1}{2}(H_0(z)H_1(-z) - H_1(z)H_0(-z)) \quad (160)$$

## QMF filter bank implementation

- If  $H_0(z)$  and  $H_1(z)$  are chosen according to equations (135) and (136), then

$$\begin{aligned}
 H(z) &= \frac{1}{2} \left[ \frac{1}{4} (F_0(z) + F_1(z)) (F_0(-z) - F_1(-z)) \right. \\
 &\quad \left. - \frac{1}{4} (F_0(z) - F_1(z)) (F_0(-z) + F_1(-z)) \right] \\
 &= \frac{1}{2} \left[ \frac{1}{2} (F_0(z)F_1(-z) - F_0(-z)F_1(z)) \right] \quad (161)
 \end{aligned}$$



## QMF filter bank implementation

- The overall transfer function of a 2-band QMF filter bank whose analysis filters are  $H_0(z)$  and  $H_1(z)$  is given by

$$H(z) = -\frac{1}{2}z^{-1}\hat{F}_0(z^2)\hat{F}_1(z^2) \quad (162)$$

where  $F_0(z) = \hat{F}_0(z^2)$  and  $F_1(z) = z^{-1}\hat{F}_1(z^2)$ , since  $F_0(z)$  and  $F_1(z)$  are half band filters, given that the specifications of  $H_0(z)$  and  $H_1(z)$  are the symmetric of each other around  $\frac{\pi}{2}$ .

- It is possible to observe that the filter bank transfer function is alias free and has no magnitude distortion as  $H(z)$  consists of a product of allpass functions.

## Wave filters

- In classical analog filter design, it is widely known that doubly terminated LC lossless filters have zero sensitivity of the transfer function, with respect to the lossless L's and C's components, at frequencies where the maximal power is transferred to the load.
- Filter transfer functions that are equiripple in the passband, such as Chebyshev and elliptic filters, have several frequencies of maximal power transfer.
- Since the ripple values are usually kept small within the passband, the sensitivities of the transfer function to variations in the filter components remain small over the frequency range consisting of the entire passband. This is the reason why several methods have been proposed to generate realizations that attempt to emulate the internal operations of the doubly terminated lossless filters.

## Wave filters

- In digital-filter design, the first attempt to derive a realization starting from an analog prototype consisted of applying the bilinear transformation to the continuous-time transfer function, establishing a direct correspondence between the elements of the analog prototype and of the resulting digital filter.
- However, the direct simulation of the internal quantities, such as voltages and currents, of the analog prototype in the digital domain leads to delay-free loops, as will be seen below. These loops can not be computed sequentially, since not all node values in the loop are initially known.
- An alternative approach results from the fact that any analog  $n$ -port network can be characterized using the concepts of incident and reflected waves' quantities known from distributed parameter theory.

## Wave filters

- Through the application of the wave characterization, digital filter realizations without delay-free loops can be obtained from passive and active analog filters, using the bilinear transformation, as originally proposed by Fettweis.
- The realizations obtained using this procedure are known as wave digital filters.
- The name wave digital filter derives from the fact that wave quantities are used to represent the internal analog signals in the digital-domain simulation. The possible wave quantities are voltage, current, and power quantities. The choice between voltage or current waves is irrelevant, whereas power waves lead to more complicated digital realizations. Traditionally, voltage wave quantities are the most widely used, and, therefore, we base our presentation on that approach.

## Wave filters

- Another great advantage of the wave digital filters, when imitating doubly terminated lossless filters, is their inherent stability under linear conditions (infinite-precision arithmetic), as well as in the nonlinear case, where the signals are subjected to quantization. Also, if the states of a wave digital filter structure, imitating a passive analog network, are quantized using magnitude truncation and saturation arithmetic, no zero-input or overflow limit cycles can be sustained.
- The wave digital filters are also adequate to simulate certain analog systems, such as power systems, due to the topological equivalence with their analog counterparts.
- The transformation of a transfer function  $T(s)$  representing a continuous-time system into a discrete-time transfer function  $H(z)$  may be performed using the bilinear transformation in the following form

$$H(z) = T(s) \Big|_{s = \frac{2}{T} \frac{z-1}{z+1}} \quad (163)$$

## Wave filters

- Given the doubly terminated LC network depicted in Figure 18, if we use voltage and current variables to simulate the analog components, we have that

$$\left. \begin{aligned} I_1 &= \frac{V_i - V_2}{R_1} \\ I_2 &= \frac{V_2}{Z_C} \\ I_3 &= \frac{V_2}{Z_L} \\ I_4 &= \frac{V_2}{R_2} \\ I_1 &= I_2 + I_3 + I_4 \end{aligned} \right\} \quad (164)$$

## Wave filters

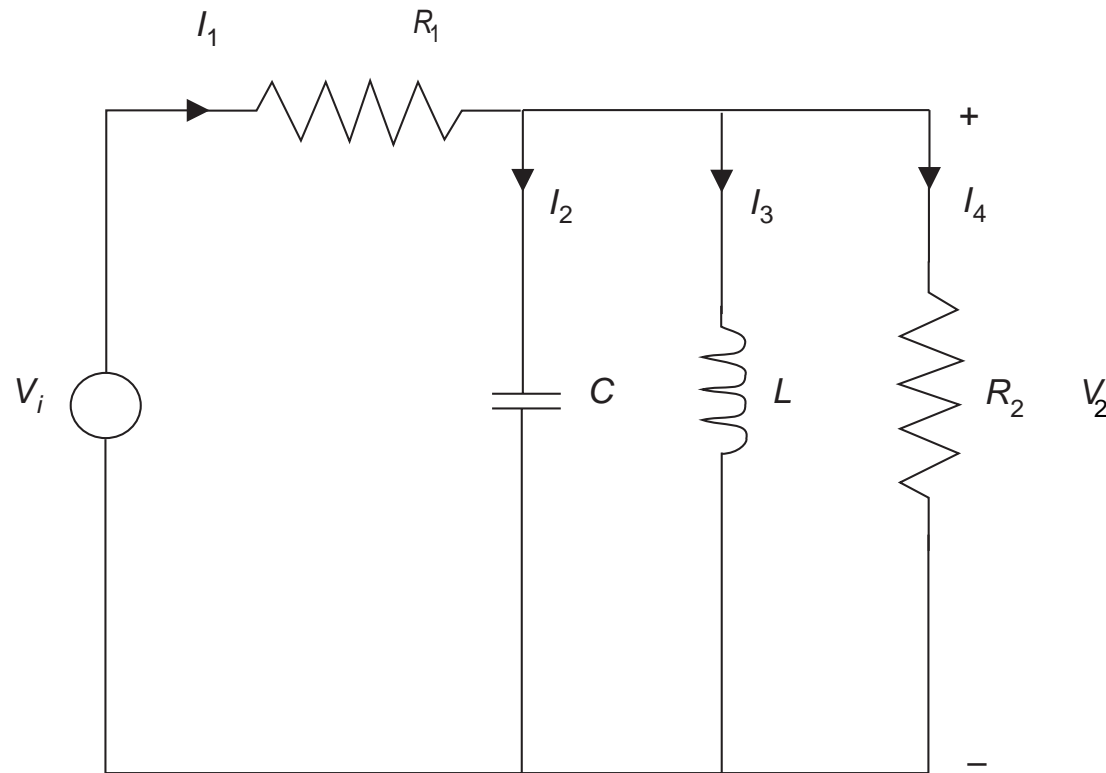


Figure 18: Doubly terminated LC network.

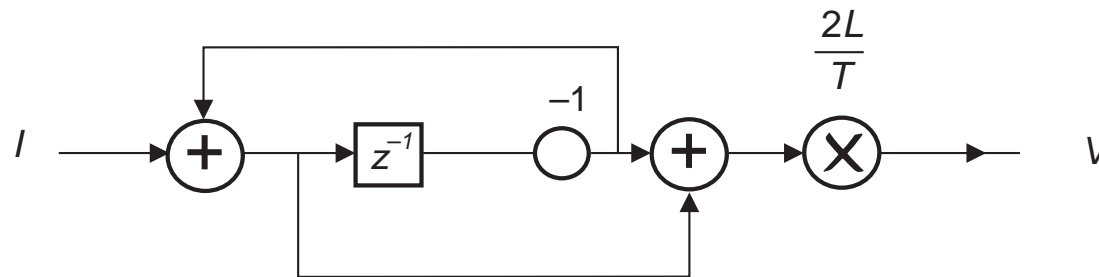
## Wave filters

- The possible representations for an inductor in the  $z$  plane will be in one of the forms shown in Figure 19, that is

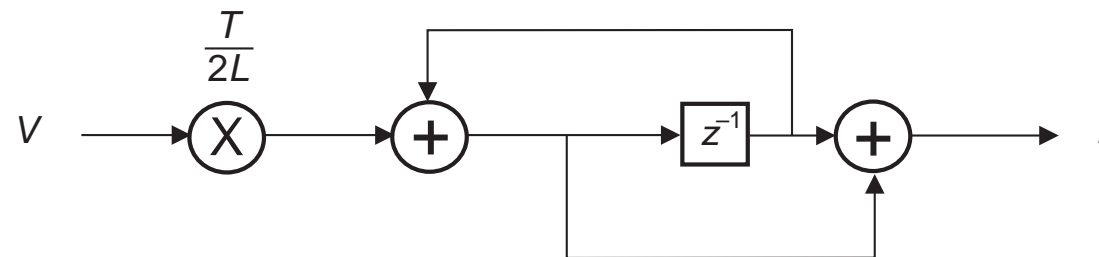
$$\left. \begin{aligned} V &= sLI = \frac{2L}{T} \frac{z-1}{z+1} I \\ I &= \frac{V}{sL} = \frac{T}{2L} \frac{z+1}{z-1} V \end{aligned} \right\} \quad (165)$$



## Wave filters



(a)



(b)

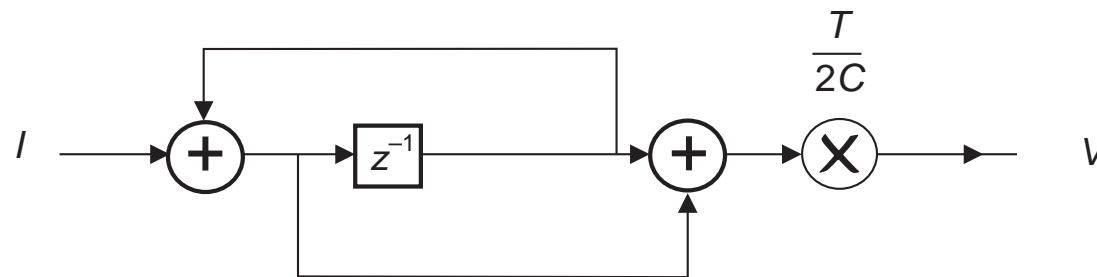
Figure 19: Two possible inductor realizations.

## Wave filters

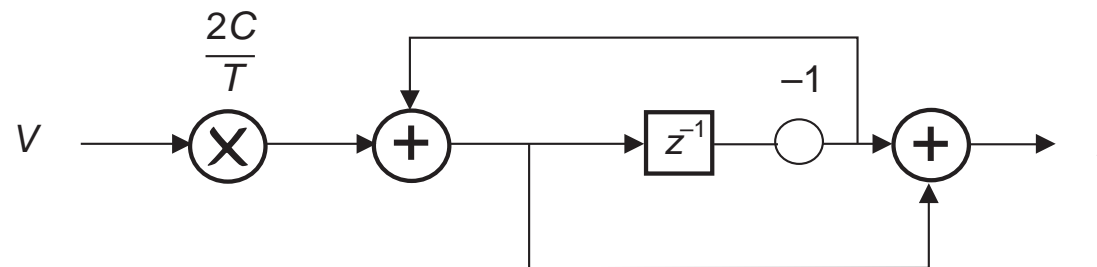
- For a capacitor, the resulting possible representations are depicted in Figure 20, such that

$$\left. \begin{aligned} V &= \frac{I}{sC} = \frac{T}{2C} \frac{z+1}{z-1} I \\ I &= sCV = \frac{2C}{T} \frac{z-1}{z+1} V \end{aligned} \right\} \quad (166)$$

## Wave filters



(a)



(b)

Figure 20: Two possible capacitor realizations.

## Wave filters

- The sources and the loads are represented in Figure 21. Therefore, using Figures 19–21, the digital simulation of the doubly terminated LC network of Figure 18 leads to the digital network shown in Figure 22, where we notice the existence of delayless loops.
- In the next subsection, we show how these loops can be avoided using the concept of wave digital filters.

## Wave filters

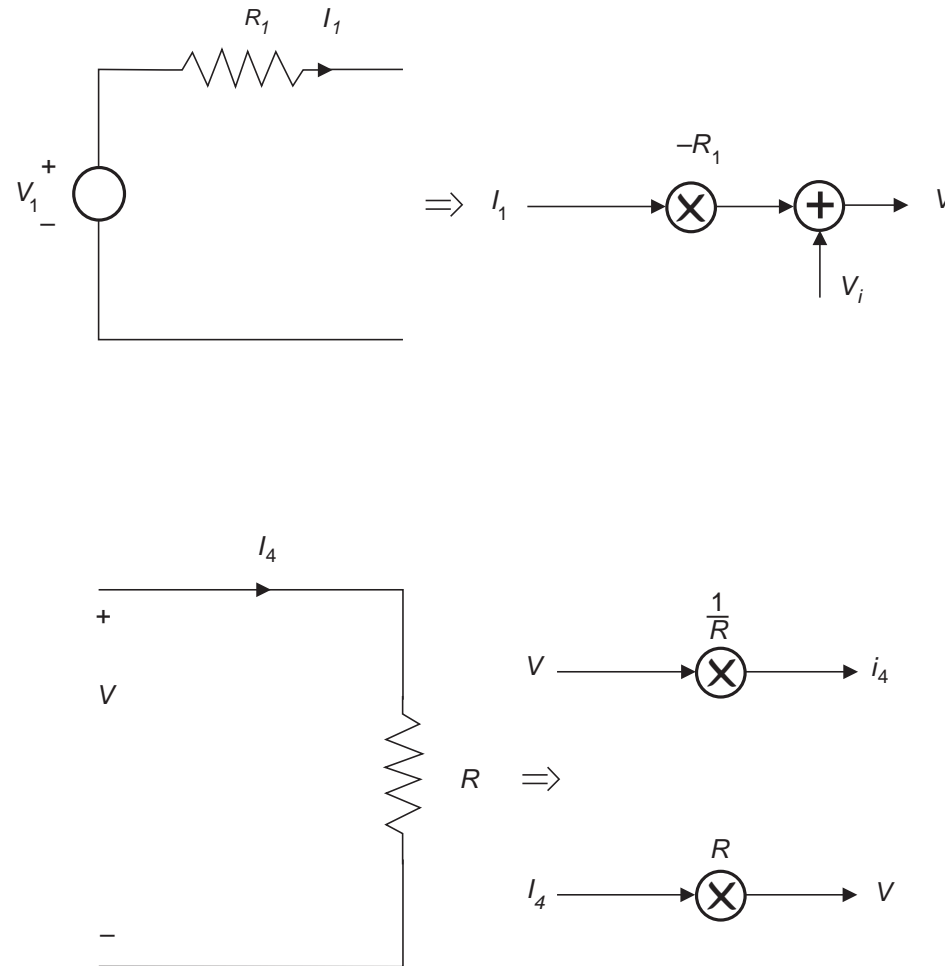


Figure 21: Termination realizations.

## Wave filters

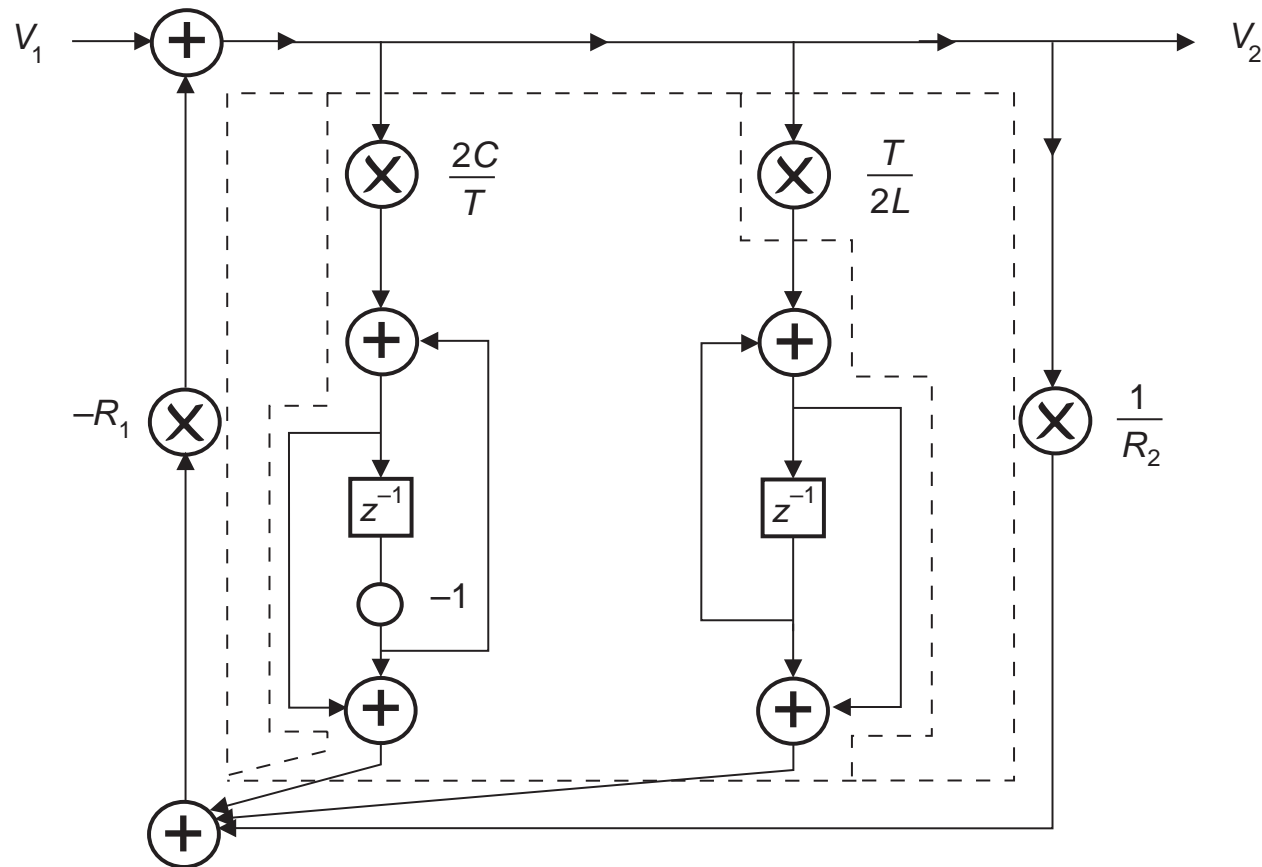


Figure 22: Digital network with delay-free loops.

## Wave elements

- As discussed in the previous subsection, the direct simulation of the branch elements of the analog network introduces delayless loops, generating a non-computable digital network. This problem can be circumvented by simulating the analog network using the wave equations that represent multiport networks instead of representing the voltages and currents in a straightforward manner.
- As shown in Figure 23, an analog one-port network can be described in terms of a wave characterization as a function of the variables

$$\left. \begin{aligned} a &= v + Ri \\ b &= v - Ri \end{aligned} \right\} \quad (167)$$

where  $a$  and  $b$  are the incident and reflected voltage wave quantities, respectively, and  $R$  is the port resistance assigned to the one-port network.

## Wave elements

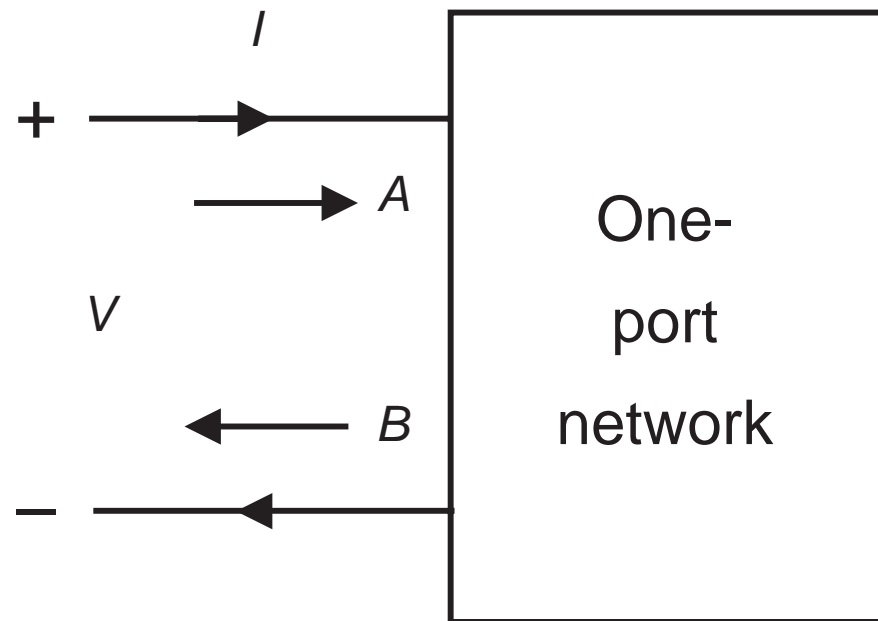


Figure 23: Convention for the incident and reflected waves.



## Wave elements

- In the frequency domain the wave quantities are  $A$  and  $B$ , such that

$$\left. \begin{aligned} A &= V + RI \\ B &= V - RI \end{aligned} \right\} \quad (168)$$

Notice how the voltage waves consist of linear combinations of the voltage and current of the one-port network.

- The value of  $R$  is a positive parameter called port resistance. A proper choice of  $R$  leads to simple multiport network realizations. In the following, we examine how to represent several analog elements using the incident and reflected waves.

## One-port elements

- For a capacitor, the following equations apply

$$\left. \begin{aligned} V &= \frac{1}{sC} I \\ B &= A \frac{V - RI}{V + RI} \end{aligned} \right\} \quad (169)$$

thus

$$B = A \frac{\frac{1}{sC} - R}{\frac{1}{sC} + R} \quad (170)$$

- By applying the bilinear transformation, we find

$$B = \frac{\frac{T}{2C}(z+1) - R(z-1)}{\frac{T}{2C}(z+1) + R(z-1)} A \quad (171)$$

## One-port elements

- The value of  $R$  that leads to a significant simplification in the implementation of  $B$  as a function of  $A$ , is

$$R = \frac{T}{2C} \quad (172)$$

and then

$$B = z^{-1} A \quad (173)$$

- The realization of  $B$  as a function of  $A$  is done as shown in Figure 24.

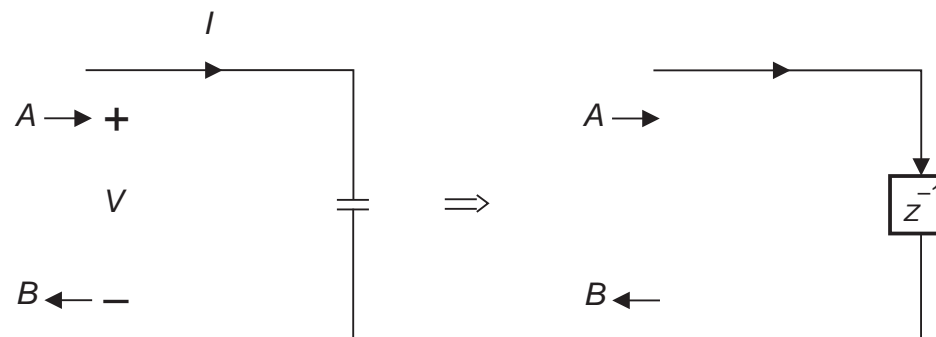


Figure 24: Wave realization for a capacitor.

## One-port elements

- Following a similar reasoning, the digital representation of several other one-port elements can be derived, along with the respective wave equations.

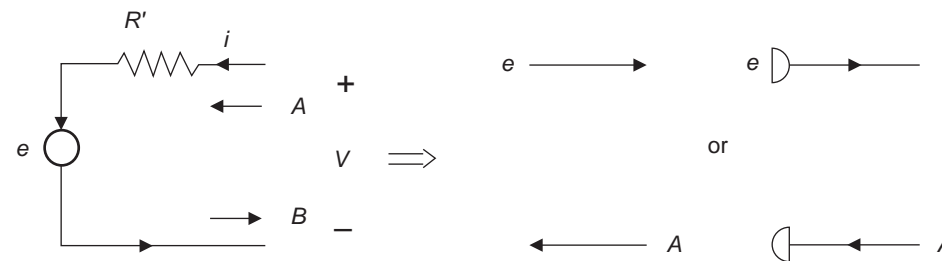


Figure 25: Wave realization for a series connection of a voltage source with resistor:  $e = V - R'I = V - RI = B$ , for  $R' = R$ .

## One-port elements

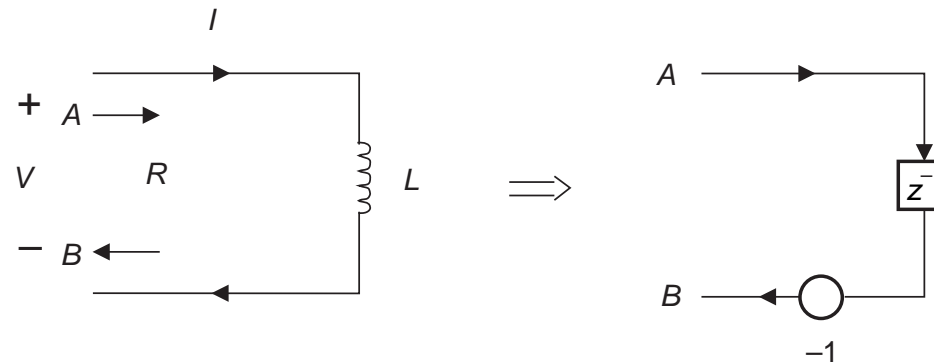


Figure 26: Wave realization for a inductor:  $R = 2L/T$ .

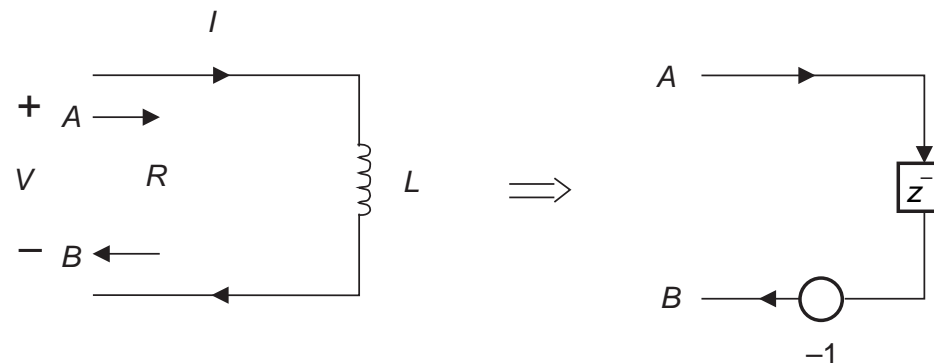


Figure 27: Wave realization for a resistor:  $B = 0$ ,  $A = 2RI$ .

## One-port elements

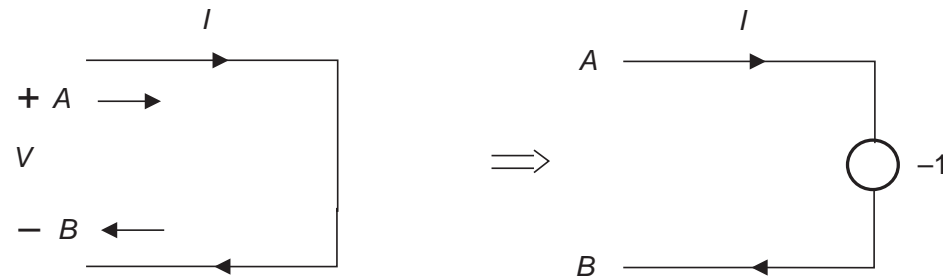


Figure 28: Wave realization for a short circuit:  $A = RI$ ,  $B = -RI$ .

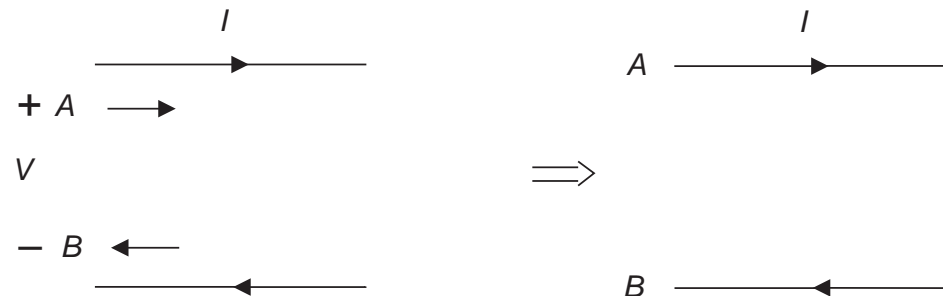


Figure 29: Wave realization for an open circuit:  $A = V$ ,  $B = V$ .

## One-port elements

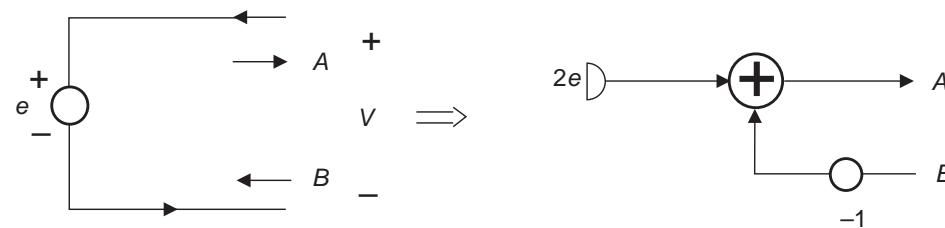


Figure 30: Wave realization for a voltage source:  $A = 2e - B$ .

## Voltage generalized immittance converter

- The voltage generalized immittance converter (VGIC), depicted in Figure 31, is a two-port network characterized by

$$\left. \begin{aligned} V_1(s) &= r(s)V_2(s) \\ I_1(s) &= -I_2(s) \end{aligned} \right\} \quad (174)$$

where  $r(s)$  is the so-called conversion function, and the pairs  $(V_1, I_1)$  and  $(V_2, I_2)$  are the VGIC voltages and currents at ports 1 and 2, respectively.

- The VGICs are not employed in the design of analog circuits due to difficulties in implementation when using conventional active devices such as transistors and operational amplifiers. However, there is no difficulty in utilizing VGICs in the design of digital filters.



## Voltage generalized immittance converter

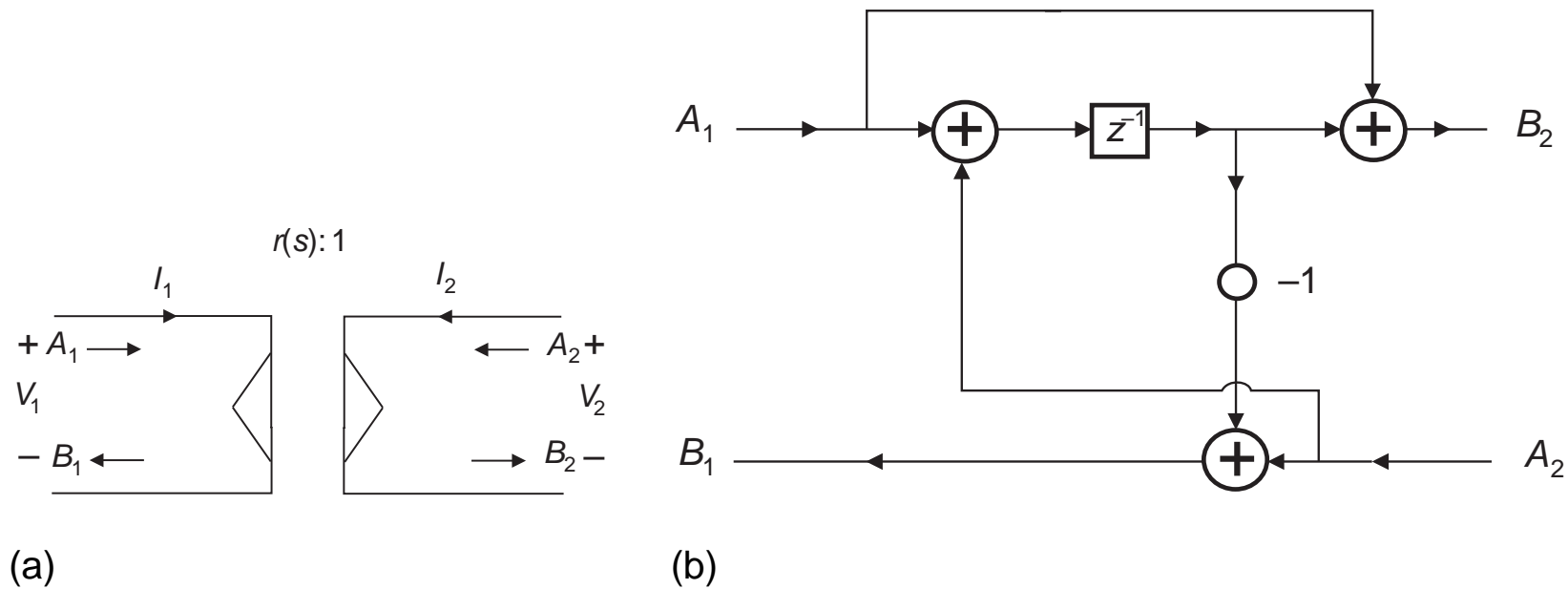


Figure 31: VGIC: (a) analog symbol; (b) digital realization:  $r(s) = T_s/2$ .

## Voltage generalized immittance converter

- The VGIC of Figure 31 may be described in terms of wave equations by

$$\left. \begin{aligned}
 A_1 &= V_1 + \frac{I_1}{G_1} \\
 A_2 &= V_2 + \frac{I_2}{G_2} \\
 B_1 &= V_1 - \frac{I_1}{G_1} \\
 B_2 &= V_2 - \frac{I_2}{G_2} \\
 V_1(s) &= r(s)V_2(s) \\
 I_1(s) &= -I_2(s)
 \end{aligned} \right\} \quad (175)$$

where  $A_i$  and  $B_i$  are the incident and reflected waves of each port, respectively, and  $G_i$  represents the conductance of port  $i$ , for  $i = 1, 2$ .

## Voltage generalized immittance converter

- After some algebraic manipulation, we can calculate the values of  $B_1$  and  $B_2$ , as functions of  $A_1$ ,  $A_2$ ,  $G_1$ ,  $G_2$ , and  $r(s)$ , as

$$\left. \begin{aligned} B_1 &= \frac{r(s)G_1 - G_2}{r(s)G_1 + G_2}A_1 + \frac{2r(s)G_2}{r(s)G_1 + G_2}A_2 \\ B_2 &= \frac{2G_1}{r(s)G_1 + G_2}A_1 + \frac{G_2 - r(s)G_1}{r(s)G_1 + G_2}A_2 \end{aligned} \right\} \quad (176)$$

- By applying the bilinear transformation, and by choosing  $G_2 = G_1$  and  $r(s) = \frac{T}{2}s$ , which lead to a simple digital realization, as seen in Figure 31b, the following relations result:

$$\left. \begin{aligned} B_1 &= -z^{-1}A_1 + (1 - z^{-1})A_2 \\ B_2 &= (1 + z^{-1})A_1 + z^{-1}A_2 \end{aligned} \right\} \quad (177)$$

## Current generalized immittance converter

- The current generalized immittance converter CGIC is described by

$$\left. \begin{aligned} V_1 &= V_2 \\ I_1 &= -h(s)I_2 \end{aligned} \right\} \quad (178)$$

Choosing  $G_1 = \frac{2G_2}{T}$  and  $h(s) = s$ , a simple realization for the CGIC results, as illustrated in Figure 32.

## Current generalized immittance converter

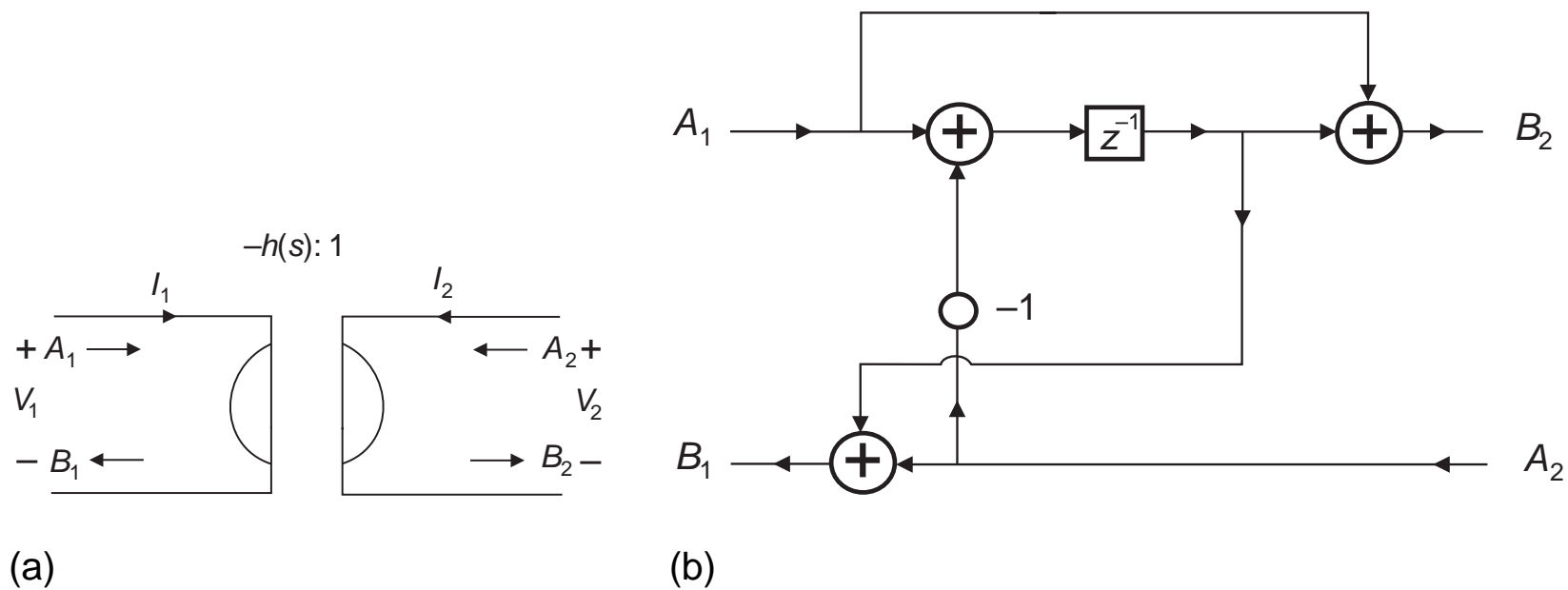


Figure 32: CGIC: (a) analog symbol; (b) digital realization:  $h(s) = s$ .

## Transformer

- A transformer with a turn ratio of  $n : 1$ , and with port resistances  $R_1$  and  $R_2$ , with  $\frac{R_2}{R_1} = \frac{1}{n^2}$ , has a digital representation as shown in Figure 33.

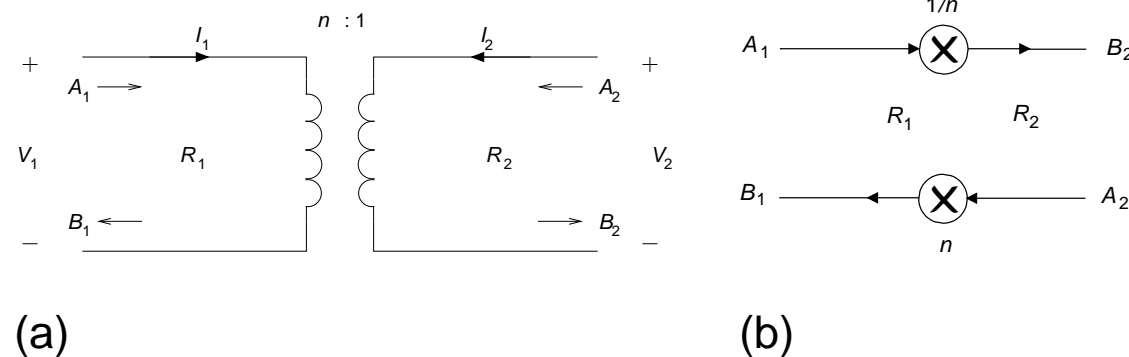


Figure 33: Transformer digital representation: (a) analog symbol; (b) digital realization:

$$R_2/R_1 = 1/n^2.$$

## Gyrator

- A gyrator is a lossless two-port element described by

$$\left. \begin{aligned} V_1 &= -RI_2 \\ V_2 &= RI_1 \end{aligned} \right\} \quad (179)$$

- It can be easily shown in this case that  $B_2 = A_1$  and  $B_1 = -A_2$ , with  $R_1 = R_2 = R$ . The digital realization of a gyrator is depicted in Figure 34.

## Gyrator

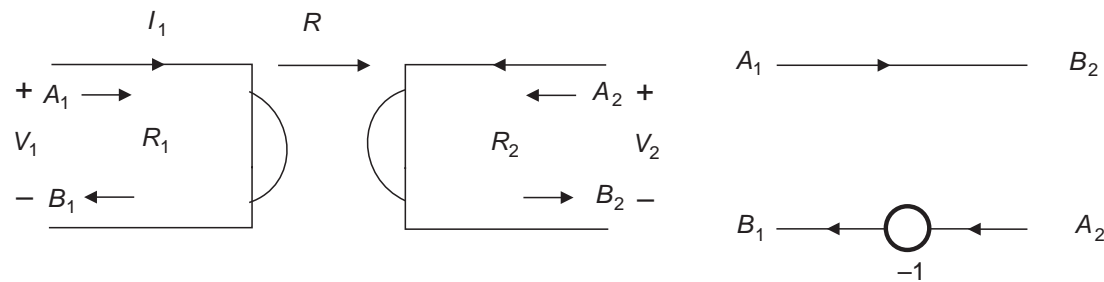


Figure 34: Gyrator digital representation.



## Wave elements

- We are now equipped with the digital representations of the main analog elements which serve as the basic building blocks for the realization of wave digital filters.
- However, to achieve our goal, we still have to learn how to interconnect these building blocks.
- To avoid delay-free loops, this must be done in the same way as these blocks are interconnected in the reference analog filter.
- Since the port resistances of the various elements are different, there is also a need to derive the so-called adaptors to allow the interconnection.
- Such adaptors guarantee that the current and voltage Kirchhoff laws are satisfied at all series and parallel interconnections of ports with different port resistances.

## Two-port adaptors

- Consider the parallel interconnection of two elements with port resistances given by  $R_1$  and  $R_2$ , respectively, as shown in Figure 35. The wave equations in this case are given by

$$\left. \begin{aligned} A_1 &= V_1 + R_1 I_1 \\ A_2 &= V_2 + R_2 I_2 \\ B_1 &= V_1 - R_1 I_1 \\ B_2 &= V_2 - R_2 I_2 \end{aligned} \right\} \quad (180)$$

## Two-port adaptors

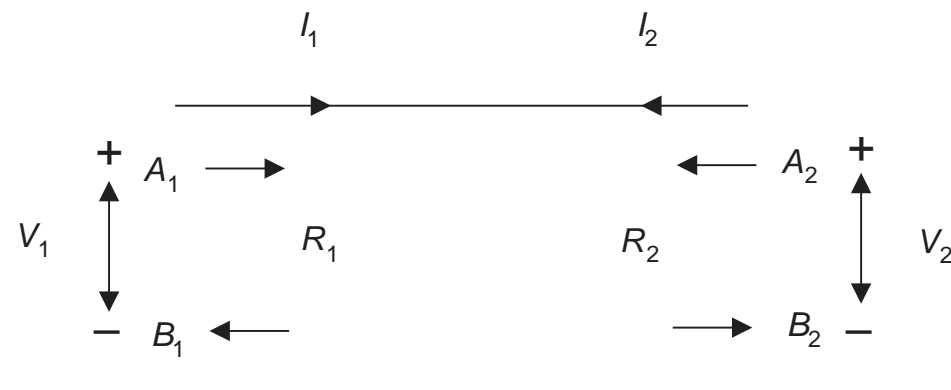


Figure 35: Two-port adaptor.

## Two-port adaptors

- Since  $V_1 = V_2$  and  $I_1 = -I_2$ , we have that

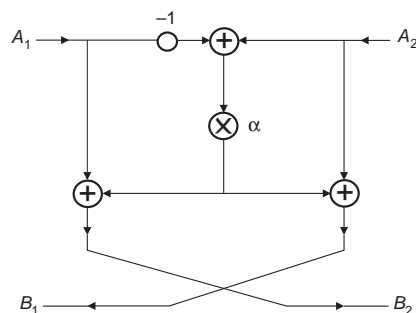
$$\left. \begin{aligned} A_1 &= V_1 + R_1 I_1 \\ A_2 &= V_1 - R_2 I_1 \\ B_1 &= V_1 - R_1 I_1 \\ B_2 &= V_1 + R_2 I_1 \end{aligned} \right\} \quad (181)$$

- Eliminating  $V_1$  and  $I_1$  from the above equations, we obtain

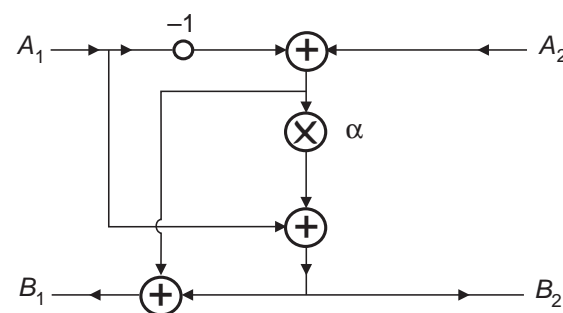
$$\left. \begin{aligned} B_1 &= A_2 + \alpha (A_2 - A_1) \\ B_2 &= A_1 + \alpha (A_2 - A_1) \end{aligned} \right\} \quad (182)$$

where  $\alpha = (R_1 - R_2)/(R_1 + R_2)$ . A realization for this general two-port adaptor is depicted in Figure 36a.

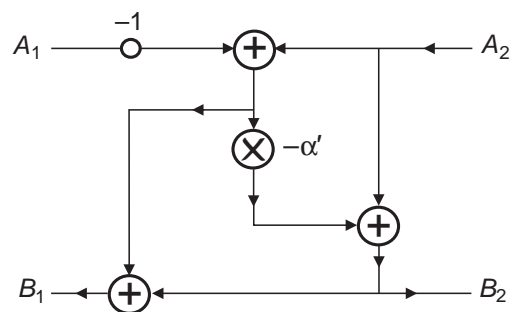
## Two-port adaptors



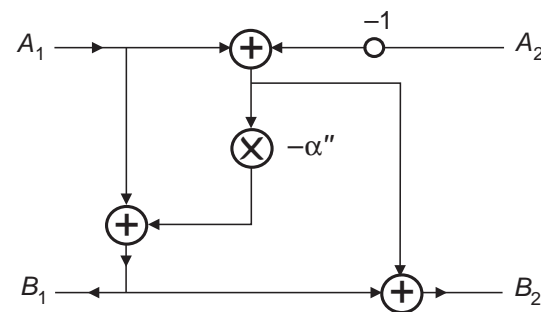
(a)



(b)



(c)



(d)

Figure 36: Possible digital realizations of the general two-port adaptors based on:  
 (a) equation (182); (b) equation (183); (c) equation (184); (d) equation (185).

## Two-port adaptors

- Expressing  $B_1$  as a function of  $B_2$  in equation (182), we get

$$\left. \begin{aligned} B_1 &= B_2 - A_1 + A_2 \\ B_2 &= A_1 + \alpha (A_2 - A_1) \end{aligned} \right\} \quad (183)$$

leading to a modified version of the two-port adaptor, as shown in Figure 36b.

- Other alternative forms of two-port adaptors are generated by expressing the equations for  $B_1$  and  $B_2$  in different ways, such as

$$\left. \begin{aligned} B_1 &= B_2 - A_1 + A_2 \\ B_2 &= A_2 - \alpha' (A_2 - A_1) \end{aligned} \right\} \quad (184)$$

where  $\alpha' = 2R_2/(R_1 + R_2)$ , generating the structure seen in Figure 36c.

## Two-port adaptors

- Or

$$\left. \begin{aligned} B_1 &= A_1 - \alpha'' (A_1 - A_2) \\ B_2 &= B_1 + A_1 - A_2 \end{aligned} \right\} \quad (185)$$

where  $\alpha'' = 2R_1 / (R_1 + R_2)$ , leading to the structure of Figure 36d.

- It is worth observing that the incident and reflected waves in any port could be expressed in the time domain, that is, through the instantaneous signal values ( $a_i(k)$  and  $b_i(k)$ ), or in the frequency domain ( $A_i(z)$  and  $B_i(z)$ ), corresponding to the steady-state description of the wave signals.

## **$n$ -port parallel adaptor**

- In cases where we need to interconnect  $n$  elements in parallel, with port resistances given by  $R_1, R_2, \dots, R_n$ , it is necessary to use an  $n$ -port parallel adaptor. The symbol to represent the  $n$ -port parallel adaptor is shown in Figure 37. Figure 38 illustrates a three-port parallel adaptor.



## $n$ -port parallel adaptor

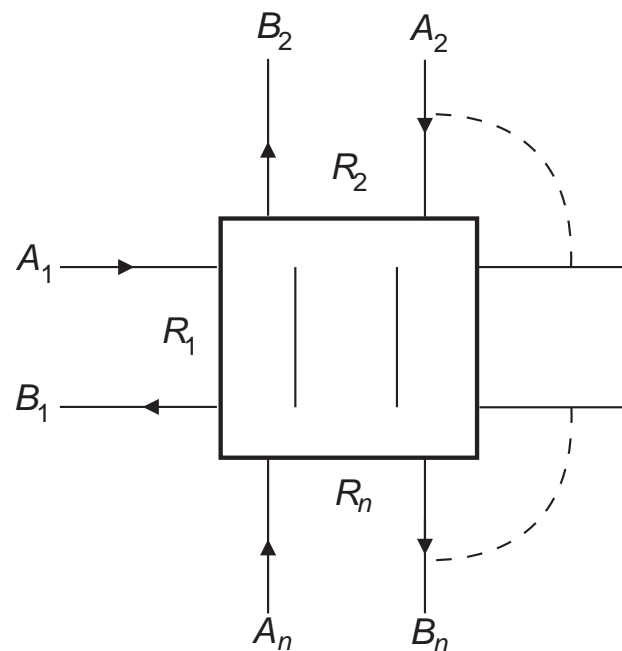


Figure 37: Symbol of the  $n$ -port parallel adaptor.

## $n$ -port parallel adaptor

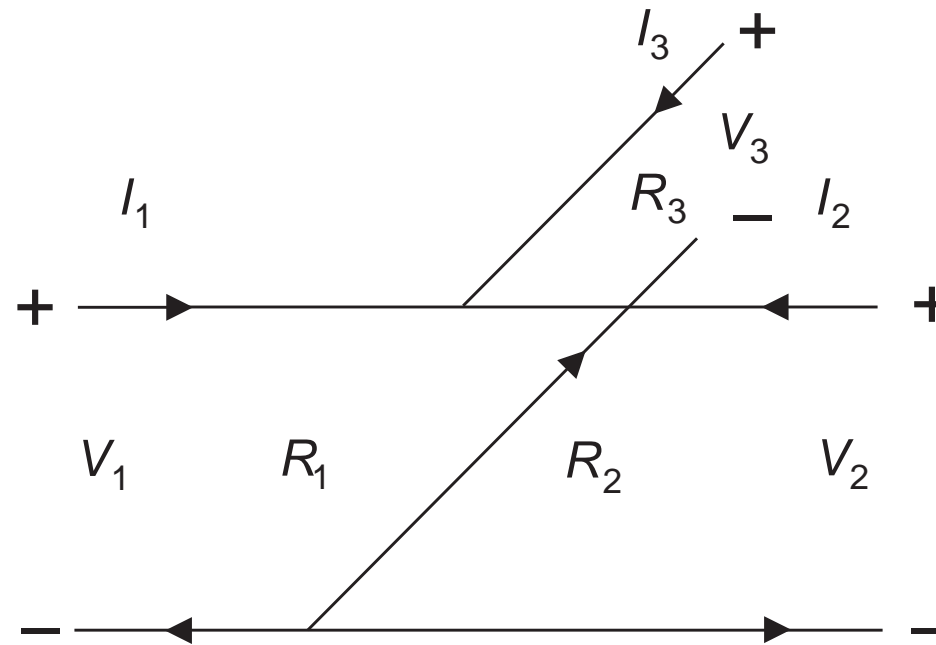


Figure 38: The three-port parallel adaptor.

## n-port parallel adaptor

- The wave equation on each port of a parallel adaptor is given by

$$\left. \begin{aligned} A_k &= V_k + R_k I_k \\ B_k &= V_k - R_k I_k \end{aligned} \right\} \quad (186)$$

for  $k = 1, 2, \dots, n$ .

- As all ports are connected in parallel, we then have that

$$\left. \begin{aligned} V_1 &= V_2 = \dots = V_n \\ I_1 + I_2 + \dots + I_n &= 0 \end{aligned} \right\} \quad (187)$$

## **n-port parallel adaptor**

- After some algebraic manipulation to eliminate  $V_k$  and  $I_k$ , we have that

$$B_k = (A_0 - A_k) \quad (188)$$

where

$$A_0 = \sum_{k=1}^n \alpha_k A_k \quad (189)$$

with

$$\alpha_k = \frac{2G_k}{G_1 + G_2 + \cdots + G_n} \quad (190)$$

and

$$G_k = \frac{1}{R_k} \quad (191)$$

## n-port parallel adaptor

- From equation 190, we note that

$$\alpha_1 + \alpha_2 + \cdots + \alpha_n = 2 \quad (192)$$

hence, one  $\alpha_i$  can be determined from the others.

- If we calculate  $\alpha_n$  as a function of the remaining  $\alpha_i$ , we can express  $A_0$  as

$$\begin{aligned}
 A_0 &= \sum_{k=1}^{n-1} \alpha_k A_k + \alpha_n A_n \\
 &= \sum_{k=1}^{n-1} \alpha_k A_k + [2 - (\alpha_1 + \alpha_2 + \cdots + \alpha_{n-1})] A_n \\
 &= 2A_n + \sum_{k=1}^{n-1} \alpha_k (A_k - A_n)
 \end{aligned} \quad (193)$$

where only  $(n - 1)$  multipliers are required to calculate  $A_0$ .

## **n-port parallel adaptor**

- In this case, port  $n$  is called the dependent port. It is also worth observing that if we have several port resistances  $R_k$  with the same value, the number of multiplications can be further reduced. If, however,  $\sum_{k=1}^{n-1} \alpha_k \approx 2$ , the error in computing  $\alpha_n$  may be very high due to quantization effects. In this case, it is better to choose another port  $k$ , with  $\alpha_k$  as large as possible, to be the dependent one.
- In practice, the three-port adaptors are the most widely used in wave digital filters. A possible implementation for a three-port parallel adaptor is shown in Figure 39a, which corresponds to the direct realization of equation (188), with  $\bar{A}_0$  calculated using equation (193).

## n-port parallel adaptor

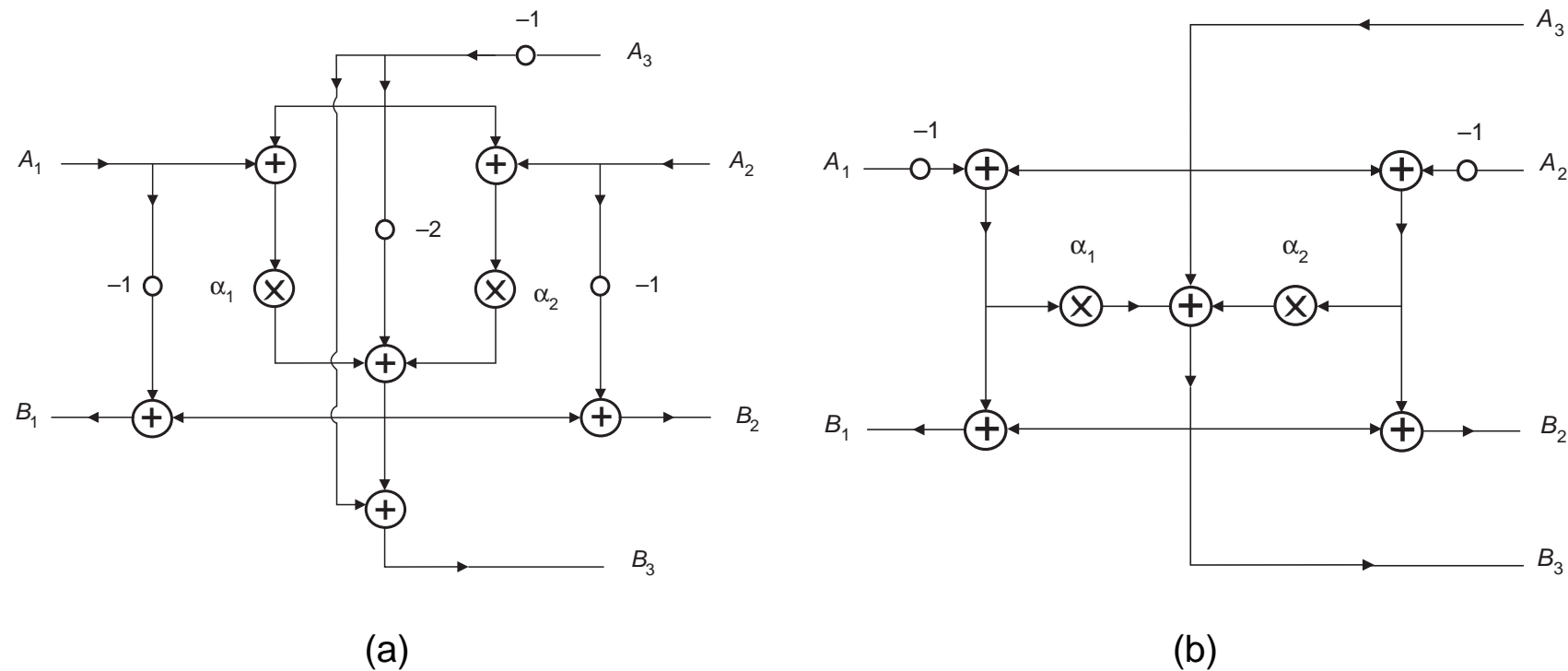


Figure 39: Possible digital realizations of the three-port parallel adaptor based on: (a) equation (193); (b) equation (194).

## n-port parallel adaptor

- Substituting equation (193) into equation (188), with  $k = n$ , then

$$\left. \begin{aligned} B_n &= (A_0 - A_n) = A_n + \sum_{k=1}^{n-1} \alpha_k (A_k - A_n) \\ B_k &= (A_0 - A_k) = B_n + A_n - A_k \end{aligned} \right\} \quad (194)$$

for  $k = 1, 2, \dots, (n - 1)$ , and we end up with an alternative realization for the three-port parallel adaptor, as seen in Figure 39b.

- Analyzing equation (194), we observe that the reflection wave  $B_i$  is directly dependent on the incident wave  $A_i$  at the same port. Hence, if two arbitrary adaptors are directly interconnected, a delay-free loop will appear between the two adaptors, as shown in Figure 40.



## n-port parallel adaptor

- A solution to this problem is to choose one of the  $\alpha$  in the adaptor to be equal to 1. For example,  $\alpha_n = 1$ . In this case, according to equations (188), (190), and (193), the equations describing the parallel adaptor become

$$\left. \begin{aligned} G_n &= G_1 + G_2 + \cdots + G_{n-1} \\ B_n &= \sum_{k=1}^{n-1} \alpha_k A_k \end{aligned} \right\} \quad (195)$$

and the expression for  $B_n$  becomes independent of  $A_n$ , thus eliminating the delay-free loops at port  $n$ .

- In this case, equation (192) becomes

$$\alpha_1 + \alpha_2 + \cdots + \alpha_{n-1} = 1 \quad (196)$$

which still allows one of the  $\alpha_i$ , for  $i = 1, 2, \dots, (n - 1)$ , to be expressed as a function of the others, thus eliminating one multiplication.

## **n-port parallel adaptor**

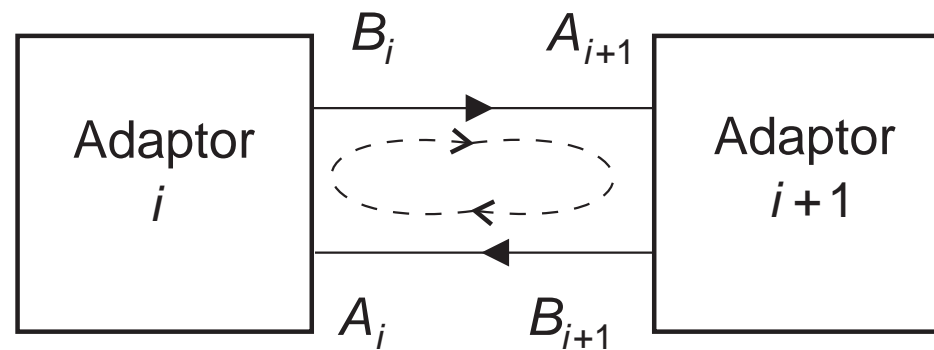


Figure 40: Adaptor interconnection.

## n-port parallel adaptor

- It is worth observing that choosing an  $\alpha = 1$  does not imply any loss of generality in the filter design. In fact, at the ports corresponding to these coefficients, the port resistances can be chosen arbitrarily, since they are used only for interconnection, and therefore they do not depend on any component value of the analog prototype.
- For example, the resistance of the ports common to two interconnected adaptors must be the same but can have arbitrary values. For instance, in the case of a three-port parallel adaptor, the describing equations would be

$$\left. \begin{aligned}
 \alpha_2 &= 1 - \alpha_1 \\
 B_3 &= \alpha_1 A_1 + (1 - \alpha_1) A_2 \\
 B_2 &= (A_0 - A_2) = \alpha_1 A_1 + (1 - \alpha_1) A_2 + A_3 - A_2 = \alpha_1 (A_1 - A_2) + A_3 \\
 B_1 &= \alpha_1 A_1 + (1 - \alpha_1) A_2 + A_3 - A_1 = (1 - \alpha_1) (A_2 - A_1) + A_3
 \end{aligned} \right\} \quad (197)$$

the realization of which is depicted in Figure 41.

## **$n$ -port parallel adaptor**

- Note that the port with no possibility of delay-free loops is marked with  $(\vdash)$ , and is known as the reflection-free port.
- A parallel adaptor, as illustrated in Figure 42, can be interpreted as a parallel connection of  $n$  ports with  $(n - 2)$  auxiliary ports, which are introduced to provide separation among several external ports. The same figure also shows the symbolic representation of the  $n$ -port parallel adaptor consisting of several three-port adaptors.

## n-port parallel adaptor

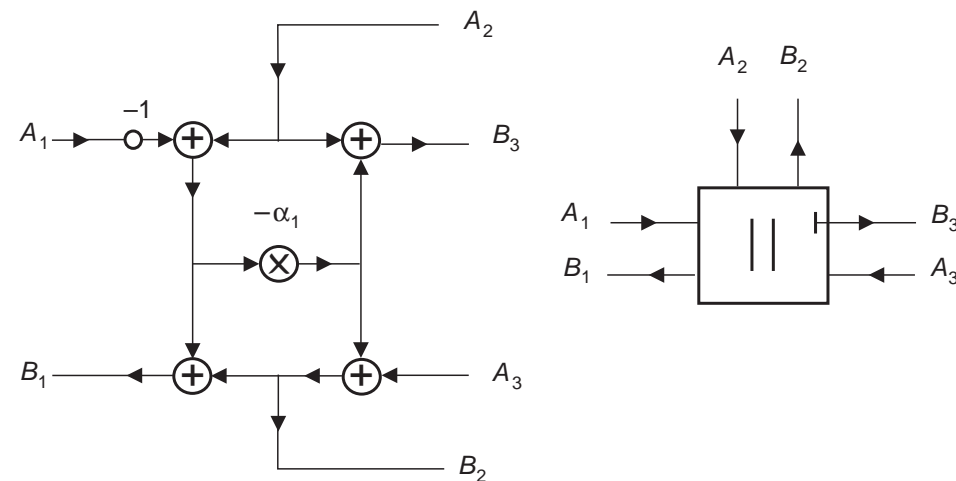
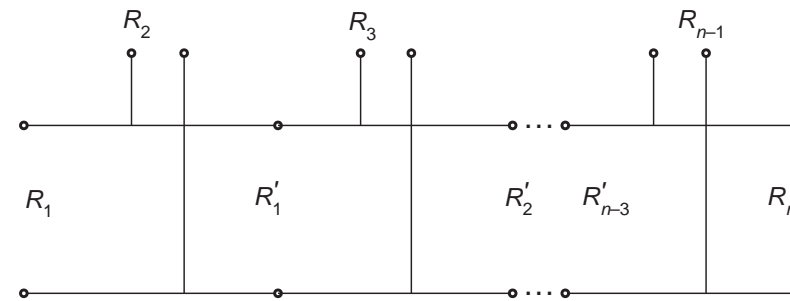
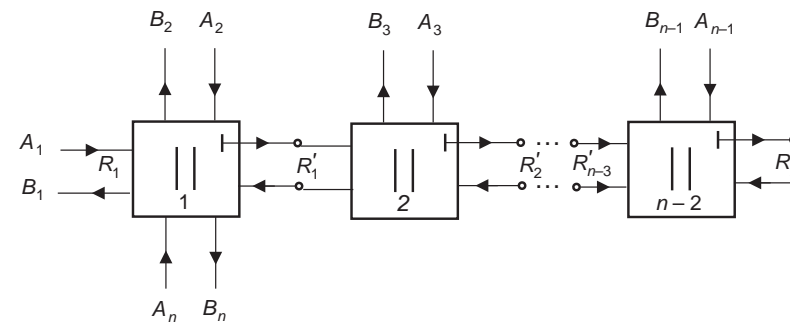


Figure 41: Reflection-free parallel adaptor at port 3.

## n-port parallel adaptor



(a)



(b)

Figure 42: The  $n$ -port parallel adaptor: (a) equivalent connection; (b) interpretation as several three-port parallel adaptors.

## **n-port series adaptor**

- In the situation where we need to interconnect  $n$  elements in series with distinct port resistances  $R_1, R_2, \dots, R_n$ , we need to use an  $n$ -port series adaptor, whose symbol is shown in Figure 43.
- In this case, the wave equations for each port are

$$\left. \begin{aligned} A_k &= V_k + R_k I_k \\ B_k &= V_k - R_k I_k \end{aligned} \right\} \quad (198)$$

for  $k = 1, 2, \dots, n$ .

## **n-port series adaptor**

- We then must have that

$$\left. \begin{array}{l} V_1 + V_2 + \cdots + V_n = 0 \\ I_1 = I_2 = \cdots = I_n = I \end{array} \right\} \quad (199)$$

- Figure 44 depicts a possible three-port series adaptor.



## $n$ -port series adaptor

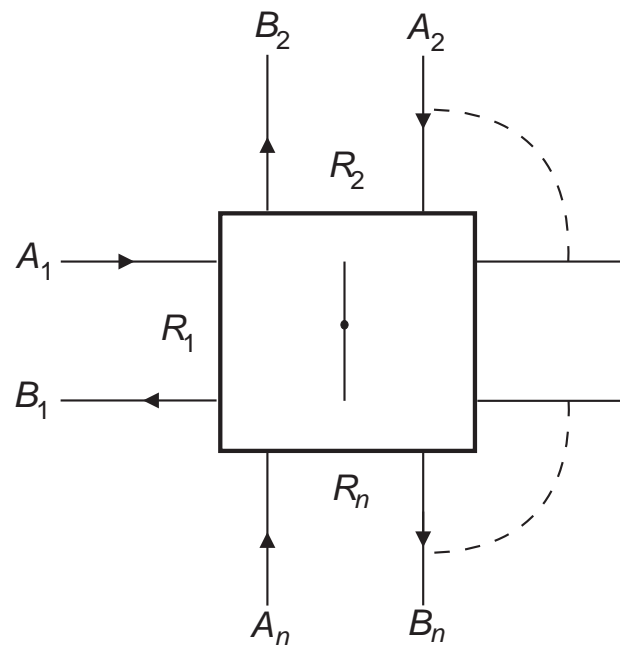


Figure 43: Symbol of the  $n$ -port series adaptor.

## **n-port series adaptor**

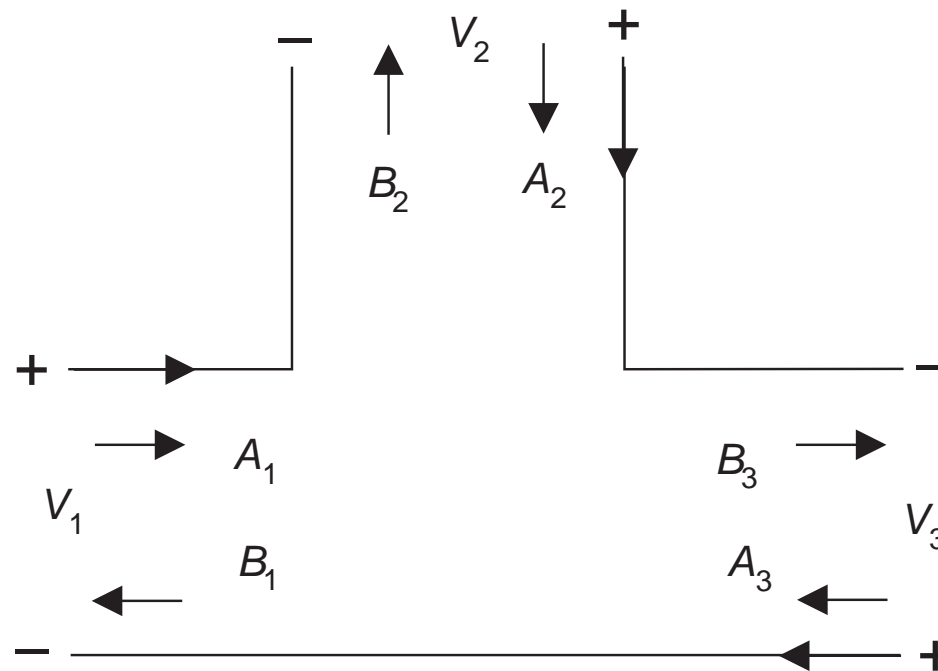


Figure 44: The three-port series adaptor.

## n-port series adaptor

- Since

$$\sum_{i=1}^n A_i = \sum_{i=1}^n V_i + \sum_{i=1}^n R_i I_i = I \sum_{i=1}^n R_i \quad (200)$$

it follows that

$$B_k = A_k - 2R_k I_k = A_k - \frac{2R_k}{\sum_{i=1}^n R_i} \sum_{i=1}^n A_i = A_k - \beta_k A_0 \quad (201)$$

where

$$\left. \begin{aligned} \beta_k &= \frac{2R_k}{\sum_{i=1}^n R_i}; & A_0 &= \sum_{i=1}^n A_i \end{aligned} \right\} \quad (202)$$

## n-port series adaptor

- From equation (202), we observe that

$$\sum_{k=1}^n \beta_k = 2 \quad (203)$$

- Therefore, it is possible to eliminate a multiplication, as previously done for the parallel adaptor, by expressing  $\beta_n$  as a function of the remaining  $\beta$ , and then

$$B_n = A_n + (-2 + \beta_1 + \beta_2 + \cdots + \beta_{n-1})A_0 \quad (204)$$

## n-port series adaptor

- Since from equation (11.165)

$$\sum_{k=1}^{n-1} B_k = \sum_{k=1}^{n-1} A_k - A_0 \sum_{k=1}^{n-1} \beta_k \quad (205)$$

then

$$B_n = A_n - 2A_0 + \sum_{k=1}^{n-1} A_k - \sum_{k=1}^{n-1} B_k = -A_0 - \sum_{k=1}^{n-1} B_k \quad (206)$$

where port  $n$  is the so-called dependent port.

- The three-port adaptor realized from the equations above is shown in Figure 45, where

$$\left. \begin{aligned} B_1 &= A_1 - \beta_1 A_0 \\ B_2 &= A_2 - \beta_2 A_0 \\ B_3 &= -(A_0 + B_1 + B_2) \end{aligned} \right\} \quad (207)$$

where  $A_0 = A_1 + A_2 + A_3$ .

## **n-port series adaptor**

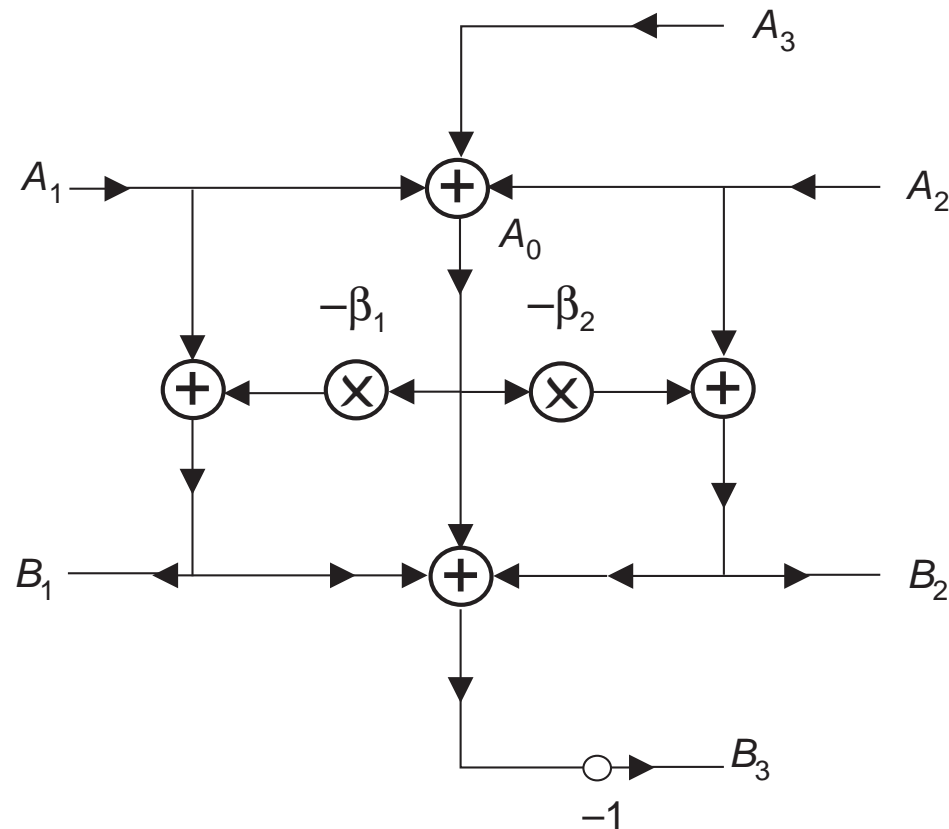


Figure 45: Digital realization of the three-port series adaptor.

## n-port series adaptor

- For the series adaptors, we avoid the delay-free loops by choosing one of the  $\beta$  equal to 1. For example, if we choose  $\beta_n = 1$ , we have that

$$\left. \begin{aligned} R_n &= R_1 + R_2 + \cdots + R_{n-1} \\ B_n &= (A_n - \beta_n A_0) = -(A_1 + A_2 + \cdots + A_{n-1}) \end{aligned} \right\} \quad (208)$$

- Equation (203) can now be replaced by

$$\beta_1 + \beta_2 + \cdots + \beta_{n-1} = 1 \quad (209)$$

which allows one of the  $\beta_i$  to be calculated from the remaining ones.

## n-port series adaptor

- A three-port series adaptor having  $\beta_3 = 1$  and  $\beta_2$  described as a function of  $\beta_1$  is described by

$$\beta_2 = 1 - \beta_1$$

$$B_1 = A_1 - \beta_1 A_0$$

$$B_2 = A_2 - (1 - \beta_1) A_0 = A_2 - (A_1 + A_2 + A_3) + \beta_1 A_0 = \beta_1 A_0 - (A_1 + A_3)$$

$$B_3 = (A_3 - A_0) = -(A_1 + A_2)$$

(210)

and its implementation is depicted in Figure 46. Note again that the port which avoids delay-free loops is marked with  $(\vdash)$ .



## n-port series adaptor

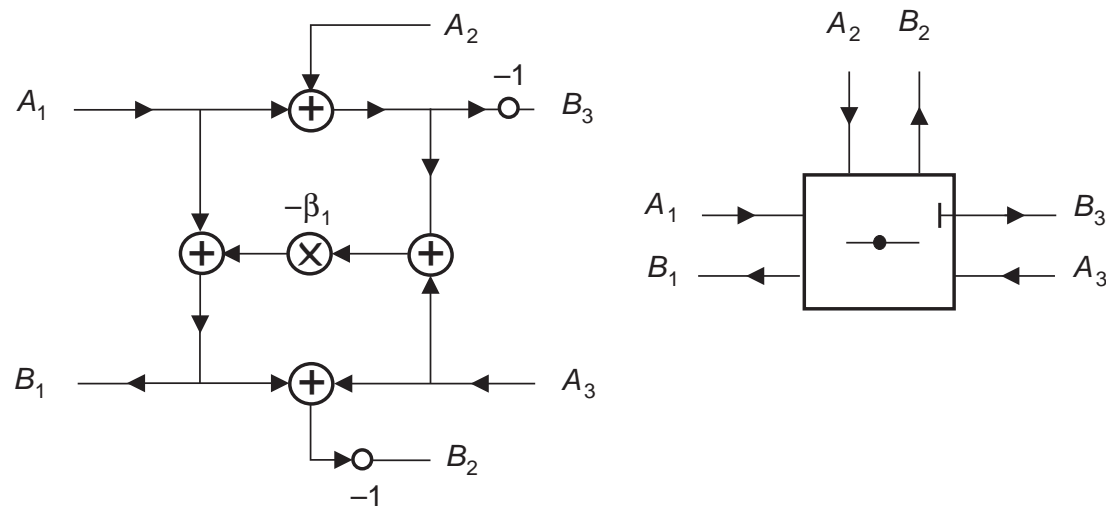
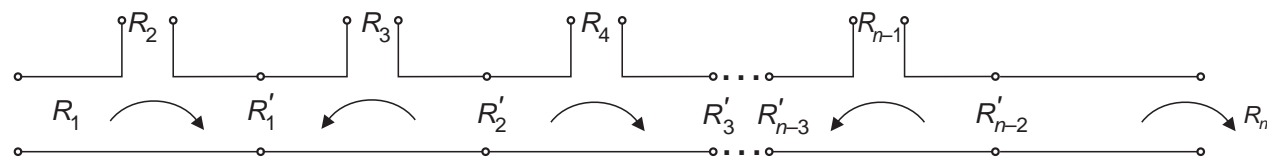


Figure 46: Reflection-free series adaptor at port 3.

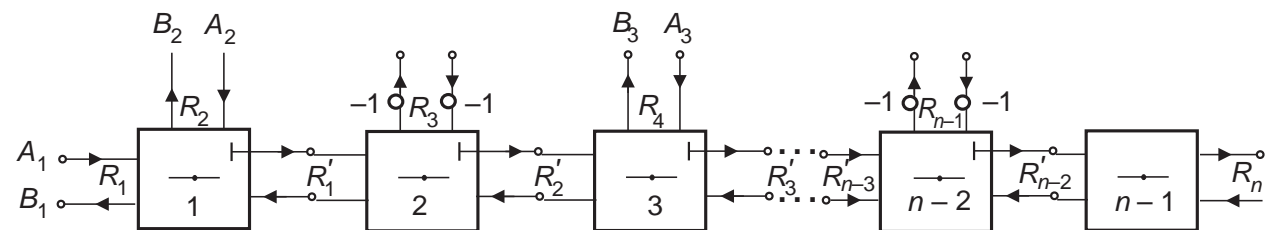
## **$n$ -port series adaptor**

- Consider now a series connection with inverted odd-port orientations, starting from port 3 and with  $(n - 2)$  auxiliary ports used for separation, as depicted in Figure 47a. We can easily show that such a series connection can be implemented through several elementary series adaptors, as shown in Figure 47b.

## n-port series adaptor



(a)



(b)

Figure 47: The  $n$ -port series adaptor: (a) equivalent connection; (b) interpretation as several three-port series adaptors.

## Lattice wave digital filters

- In some designs, it is preferable to implement a wave digital filter from a lattice analog network rather than from a ladder analog network. This is because, when implemented digitally, the lattice structures have low sensitivity in the filter passband and high sensitivity in the stopband.
- There are two explanations for the sensitivity properties of the lattice realization. First, any changes in the coefficients belonging to the adaptors of the lattice structure do not destroy its symmetry, whereas in the symmetric ladder this symmetry can be lost. The second reason applies to the design of filters with zeros on the unit circle.
- In the lattice structures, quantization usually moves these zeros away from the unit circle, whereas in the ladder structures the zeros always move around the unit circle.

## Lattice wave digital filters

- Fortunately, all symmetric ladder networks can be transformed into lattice networks by applying the so-called Bartlett bisection theorem. This subsection deals with the implementation of lattice wave digital filters.
- Given the symmetric analog lattice network of Figure 48, where  $Z_1$  and  $Z_2$  are the lattice impedances, it can be shown that the incident and reflected waves are related by

$$\left. \begin{aligned} B_1 &= \frac{S_1}{2} (A_1 - A_2) + \frac{S_2}{2} (A_1 + A_2) \\ B_2 &= -\frac{S_1}{2} (A_1 - A_2) + \frac{S_2}{2} (A_1 + A_2) \end{aligned} \right\} \quad (211)$$

where

$$\left. S_1 = \frac{Z_1 - R}{Z_1 + R}; S_2 = \frac{Z_2 - R}{Z_2 + R} \right\} \quad (212)$$

with  $S_1$  and  $S_2$  being the reflectances of the impedances  $Z_1$  and  $Z_2$ , respectively.

## Lattice wave digital filters

- That is,  $S_1$  and  $S_2$  correspond to the ratio between the reflected and incident waves at ports  $Z_1$  and  $Z_2$  with port resistance  $R$ , since

$$\frac{B}{A} = \frac{V - RI}{V + RI} = \frac{Z - R}{Z + R} \quad (213)$$

- The lattice realization then consists solely of impedance realizations, as illustrated in Figure 49. Note that, in this figure, since the network is terminated by a resistance, then  $A_2 = 0$ .

## Lattice wave digital filters

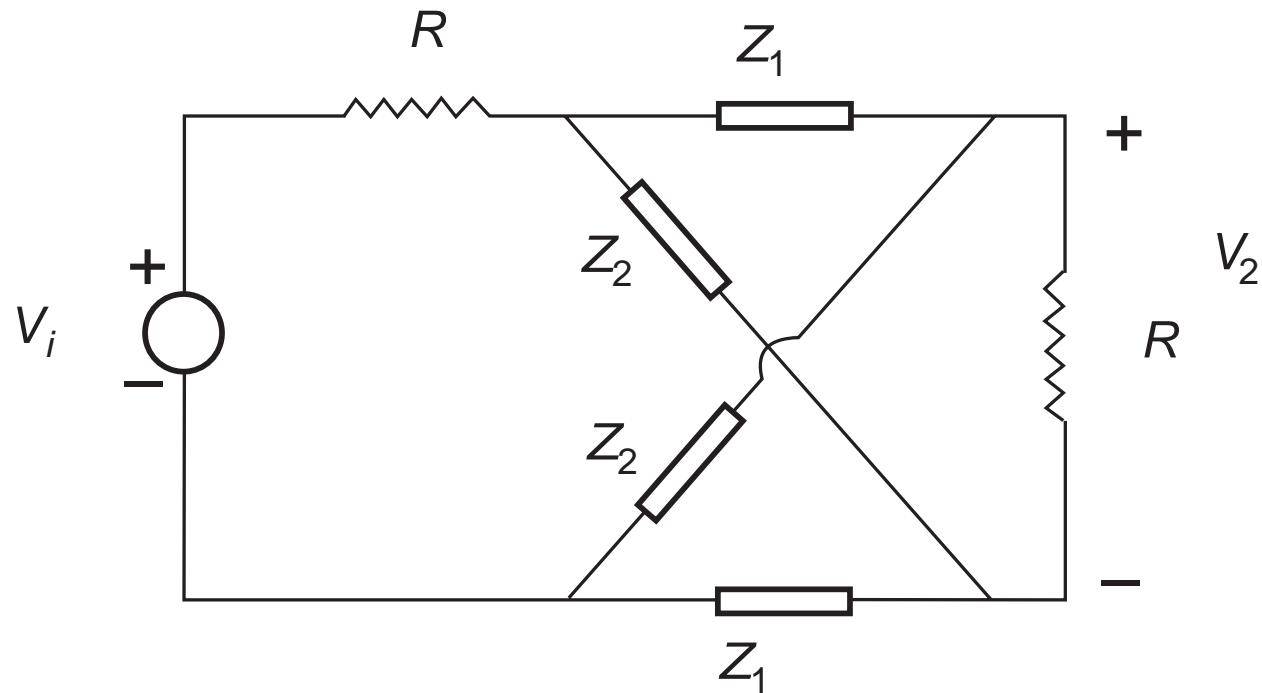


Figure 48: Analog lattice network.

## Lattice wave digital filters

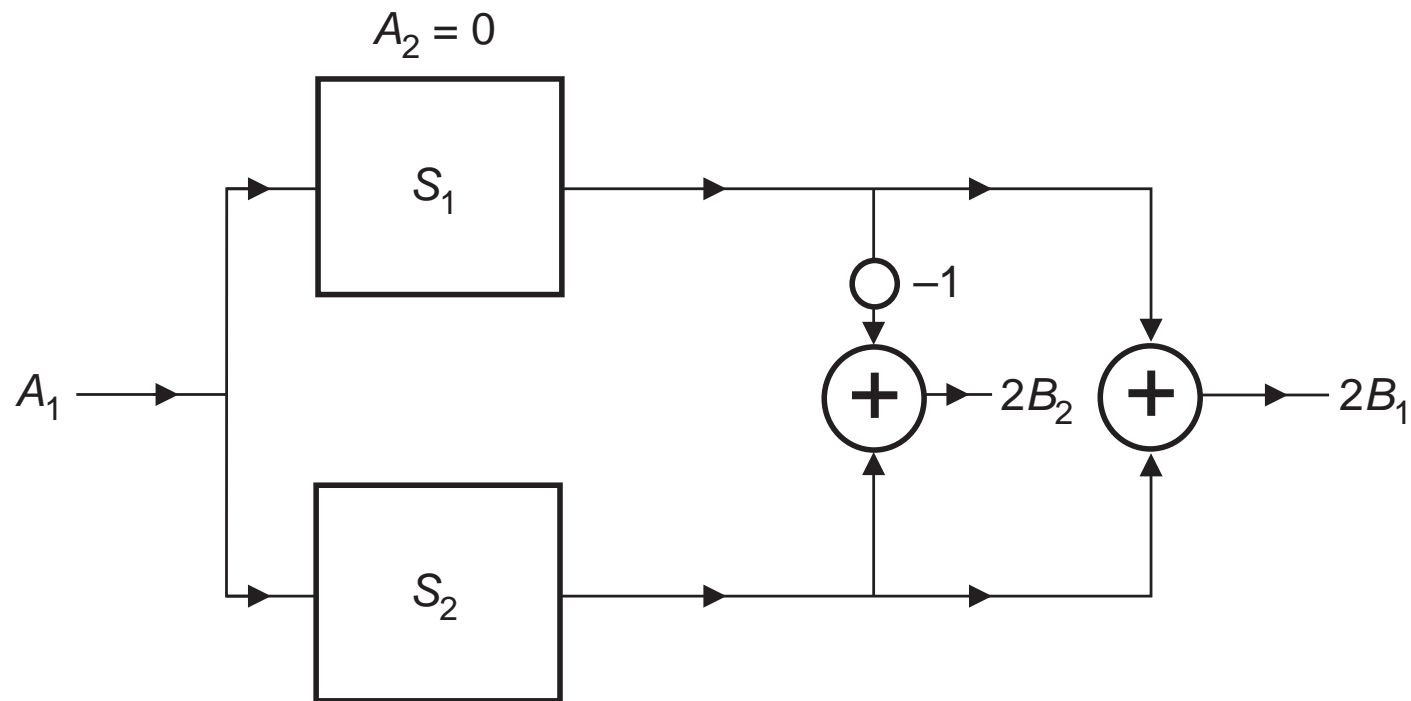


Figure 49: Wave digital lattice representation.



### Example 13.6

- Realize the lowpass filter represented in Figure 50, using a wave lattice network.

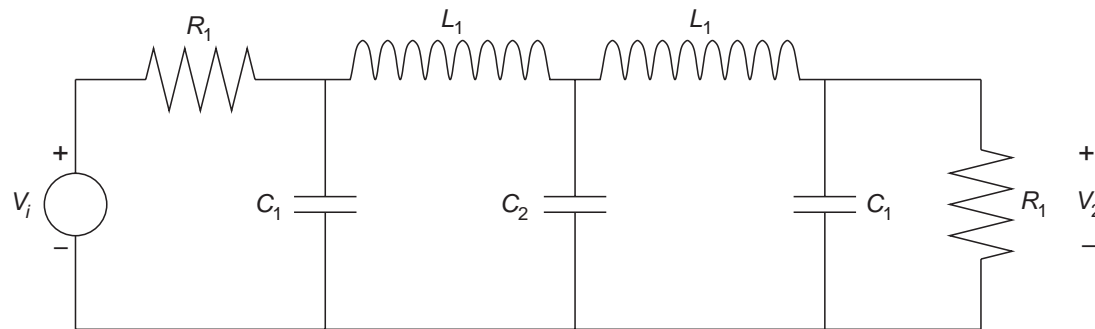


Figure 50: Lowpass RLC network.

## Example 13.6 - Solution

- The circuit in Figure 50 is a symmetric ladder network, as is made clear when it is redrawn as in Figure 51.

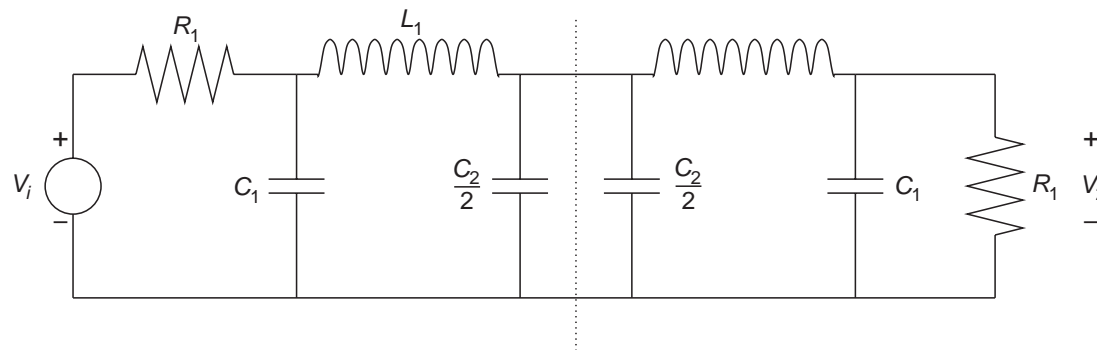


Figure 51: Symmetric ladder network.

## Example 13.6 - Solution

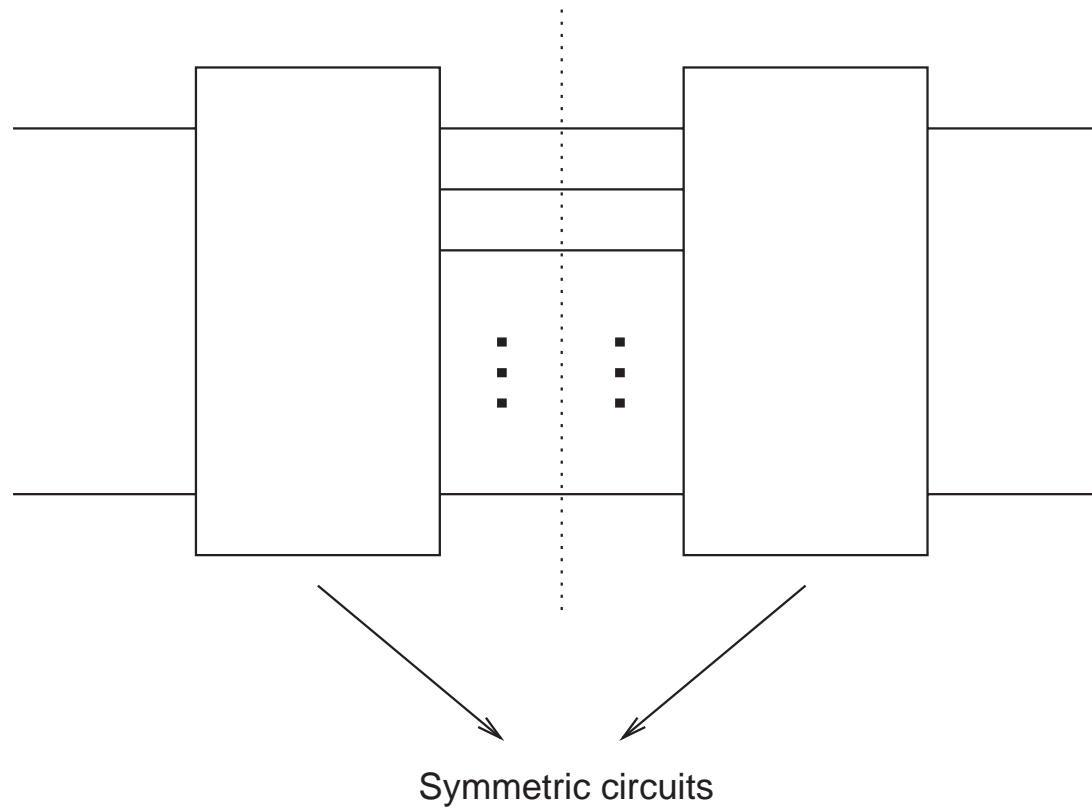


Figure 52: Generic symmetric ladder network.

## Example 13.6 - Solution

- Bartlett's bisection theorem states that, when you have a two-port network composed of two equal half networks connected by any number of wires, as illustrated in Figure 52, then it is equivalent to a lattice network as in Figure 48.
- Figure 51 is a good example of a symmetric ladder network. The impedance  $Z_1$  is equal to the input impedance of any half-network when the connection wires to the other half are short circuited, and  $Z_2$  is equal to the input impedance of any half-network when the connection wires to the other half are open.
- This is illustrated in Figure 53 below, where the determinations of the impedances  $Z_1$  and  $Z_2$  of the equivalent lattice are shown.
- Figure 54 shows the computation of  $Z_1$  and  $Z_2$  for this example.

### Example 13.6 - Solution

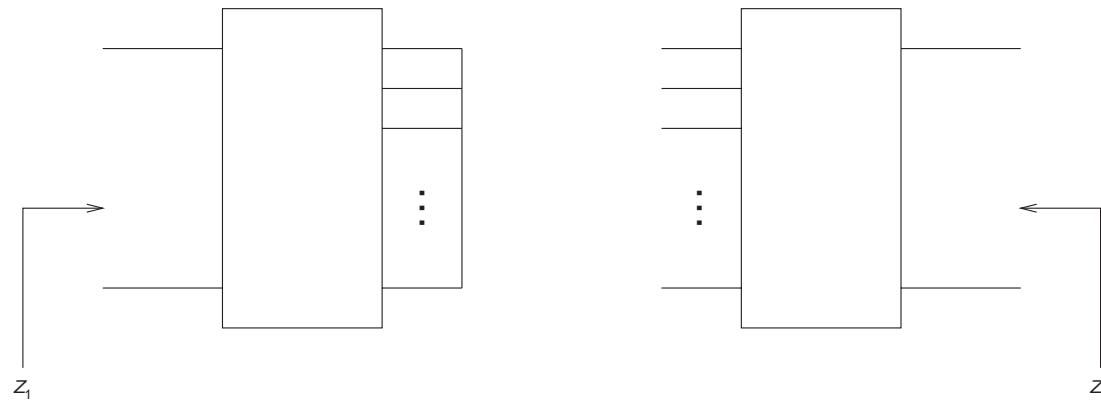


Figure 53: Computation of the lattice's impedances  $Z_1$  and  $Z_2$  for the generic case.

## Example 13.6 - Solution

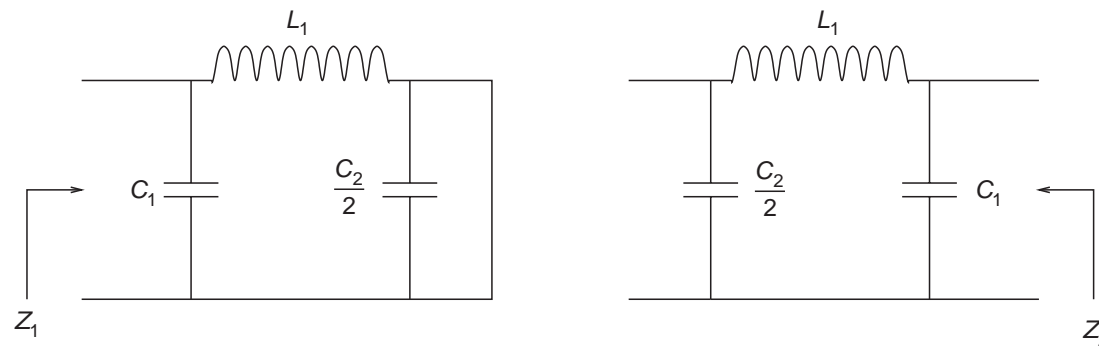


Figure 54: Computation of the lattice's impedances  $Z_1$  and  $Z_2$  for Example 13.4.

### Example 13.6 - Solution

- From the resulting lattice network, the final wave filter is then as represented in Figure 55, where

$$\alpha_1 = \frac{2G_1}{G_1 + \frac{2C_1}{T} + \frac{T}{2L_1}} \quad (214)$$

$$\alpha_2 = \frac{2\frac{2C_1}{T}}{G_1 + \frac{2C_1}{T} + \frac{T}{2L_1}} \quad (215)$$

with

$$\alpha_3 = \frac{2G_1}{G_1 + \frac{2C_1}{T} + G_3} + \frac{2G_1}{2G_1 + \frac{4C_1}{T}} = \frac{G_1}{G_1 + \frac{2C_1}{T}} \quad (216)$$

$$G_3 = G_1 + \frac{2C_1}{T} \quad (217)$$

$$\beta_1 = \frac{2R_3}{R_3 + \frac{T}{C_2} + \frac{2L_1}{T}} \quad (218)$$

$$\beta_2 = \frac{\frac{2T}{C_2}}{R_3 + \frac{T}{C_2} + \frac{2L_1}{T}} \quad (219)$$



## Example 13.6 - Solution

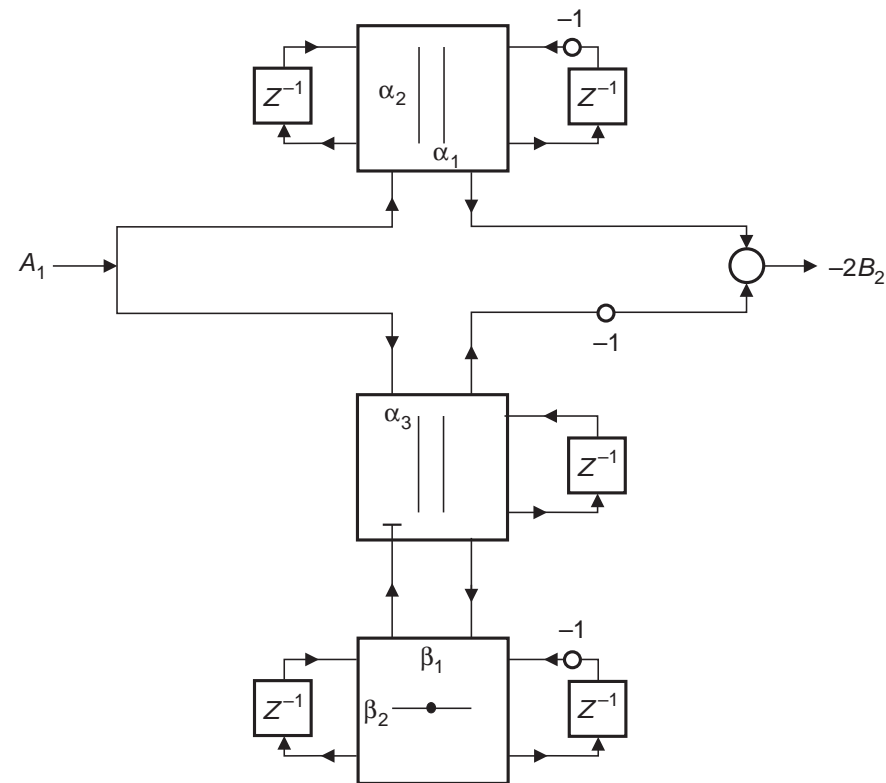


Figure 55: Resulting digital lattice network.

### Example 13.6 - Solution

- One should note that, since we want a network that generates the output voltage at the load, then  $B_2$  is the only variable of interest to us, and therefore  $B_1$  does not need to be computed.

### Example 13.6

- Realize the ladder filter represented in Figure 56 using a wave network.

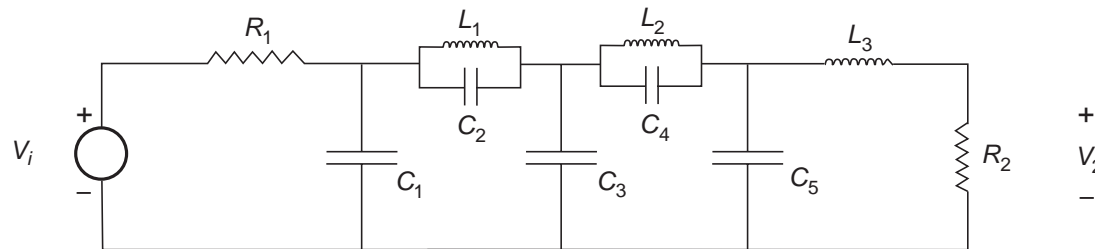


Figure 56: Ladder RLC network.

## Example 13.6 - Solution

- The connections among the elements can be interpreted as illustrated in Figure 57.

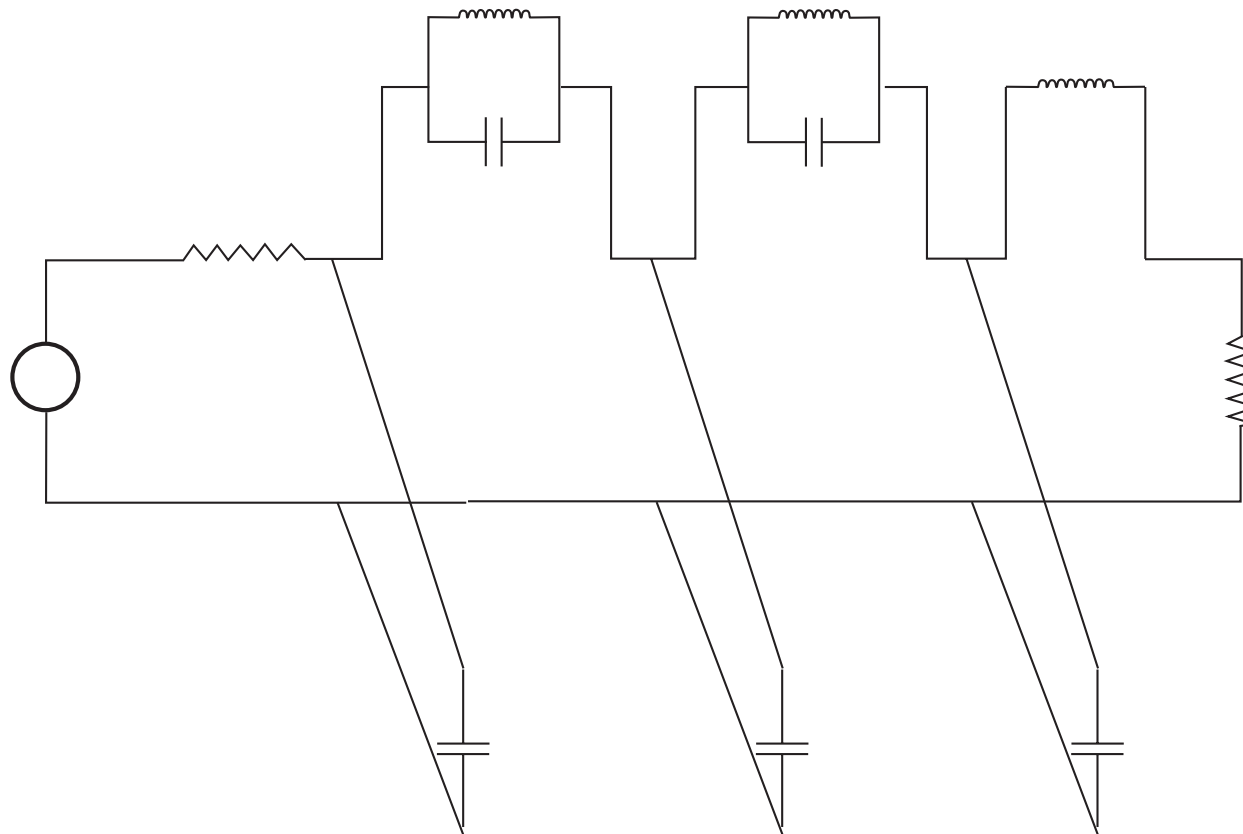


Figure 57: Component connections.

## Example 13.6 - Solution

- The resulting wave filter should be represented as in Figure 58, where the choice of the reflection-free ports is arbitrary.

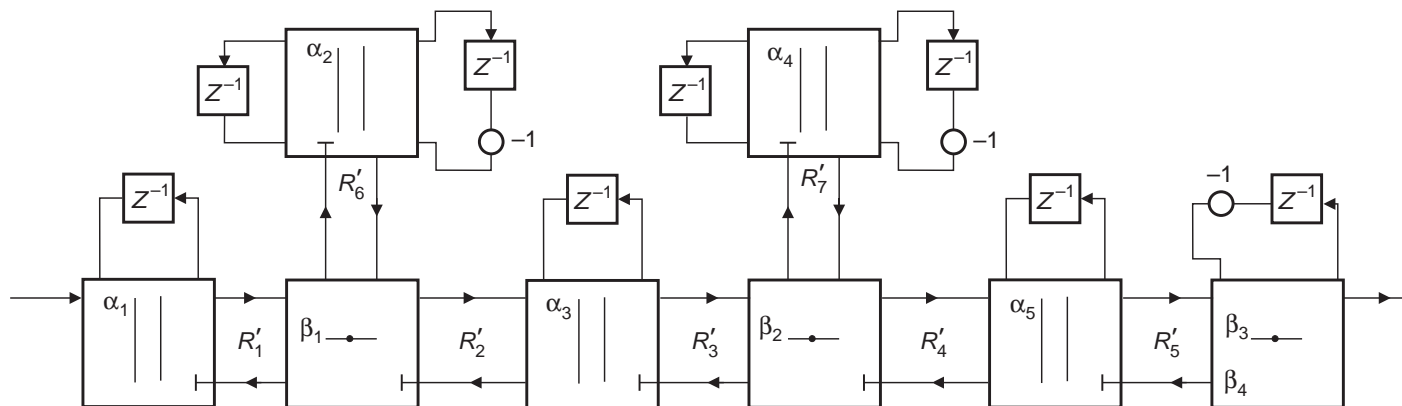


Figure 58: Resulting digital wave network.

## Example 13.6 - Solution

- The equations below describe how to calculate the multiplier coefficient of each adaptor, and Figure 59 depicts the resulting wave digital filter realization.

$$G'_1 = G_1 + \frac{2C_1}{T} \quad (220)$$

$$\alpha_1 = \frac{2G_1}{G_1 + \frac{2C_1}{T} + G'_1} = \frac{2G_1}{2G_1 + \frac{4C_1}{T}} = \frac{G_1}{G_1 + \frac{2C_1}{T}} \quad (221)$$

$$G'_6 = \frac{2C_2}{T} + \frac{T}{2L_1} \quad (222)$$

$$\alpha_2 = \frac{2\frac{2C_2}{T}}{\frac{2C_2}{T} + \frac{T}{2L_1} + G'_6} = \frac{\frac{2C_2}{T}}{\frac{2C_2}{T} + \frac{T}{2L_1}} \quad (223)$$

$$R'_2 = R'_1 + R'_6 \quad (224)$$

### Example 13.6 - Solution

$$\beta_1 = \frac{R'_1}{R'_1 + R'_6} = \frac{1}{1 + \frac{G_1 + \frac{2C_1}{T}}{\frac{2C_2}{T} + \frac{T}{2L_1}}} \quad (225)$$

$$G'_3 = G'_2 + \frac{2C_3}{T} \quad (226)$$

$$\alpha_3 = \frac{G'_2}{G'_2 + \frac{2C_3}{T}} \quad (227)$$

$$G'_7 = \frac{2C_4}{T} + \frac{T}{2L_2} \quad (228)$$

$$\alpha_4 = \frac{\frac{2C_4}{T}}{\frac{2C_4}{T} + \frac{T}{2L_2}} \quad (229)$$

$$R'_4 = R'_3 + R'_7 \quad (230)$$

### Example 13.6 - Solution

$$\beta_2 = \frac{R'_3}{R'_3 + R'_7} \quad (231)$$

$$G'_5 = G'_4 + \frac{2C_5}{T} \quad (232)$$

$$\alpha_5 = \frac{2G'_4}{G'_4 + \frac{2C_5}{T}} \quad (233)$$

$$\beta_3 = \frac{2R'_5}{R'_5 + \frac{2L_3}{T} + R_2} \quad (234)$$

$$\beta_4 = \frac{2R_2}{R'_5 + \frac{2L_3}{T} + R_2} \quad (235)$$



## Example 13.6 - Solution

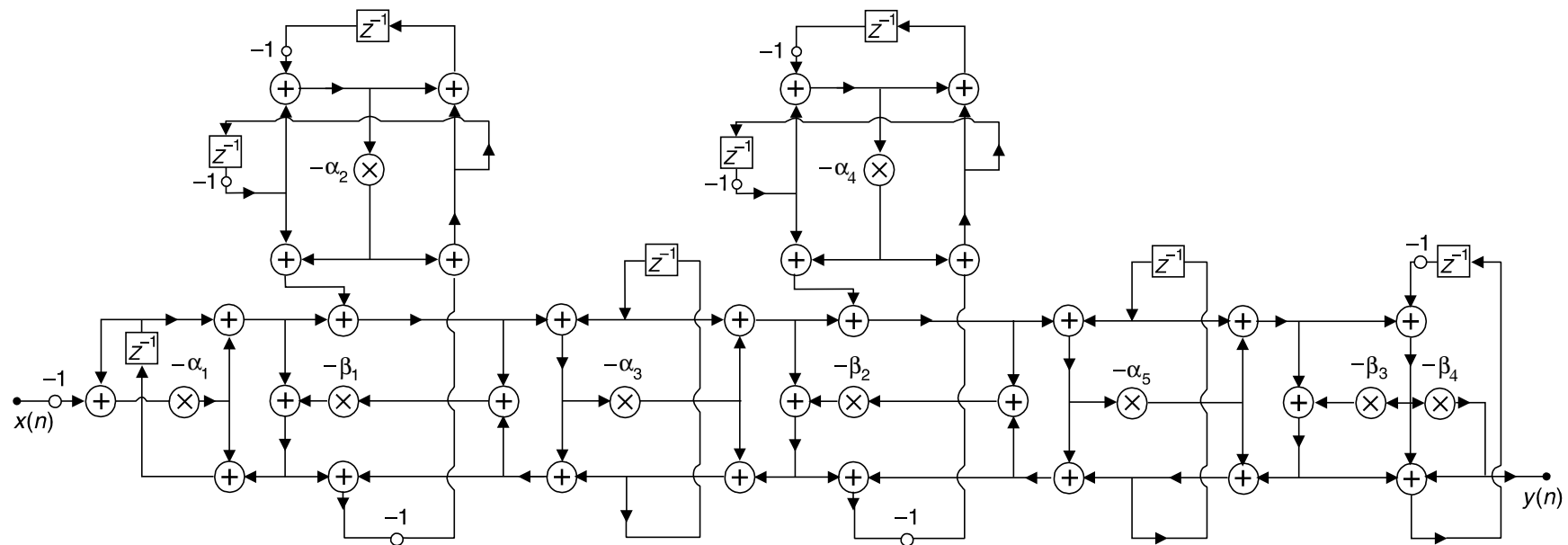


Figure 59: Resulting wave filter realization.

## Do-it-yourself: Efficient IIR structures

- **Experiment 13.1:** Consider the elliptic bandpass filter designed in Example 13.1 satisfying the specifications:

$$\left. \begin{aligned}
 A_p &= 0.5 \text{ dB} \\
 A_r &= \text{dB} \\
 \Omega_{r_1} &= 850 \text{ rad/s} \\
 \Omega_{p_1} &= 980 \text{ rad/s} \\
 \Omega_{p_2} &= 1020 \text{ rad/s} \\
 \Omega_{r_2} &= 1150 \text{ rad/s} \\
 \Omega_s &= 10\,000 \text{ rad/s}
 \end{aligned} \right\} \quad (236)$$

## Experiment 13.1

- The resulting cascade realization is described below.

Table 16: Cascade structure using direct-form second-order sections. Gain constant:

$$h_0 = 1.4362 \text{ E} - 04.$$

Coefficient	Section 1	Section 2	Section 3
$\gamma_0$	1.0000	1.0000	1.0000
$\gamma_1$	0.0000	−1.4848	−1.7198
$\gamma_2$	−1.0000	1.0000	1.0000
$m_1$	−1.6054	−1.5965	−1.6268
$m_2$	0.9843	0.9921	0.9924

## Experiment 13.1

- Given the transfer function of each section in the form

$$H_i(z) = \frac{N_i(z)}{D_i(z)} = \frac{\gamma_{0i} z^2 + \gamma_{1i} z + \gamma_{2i}}{z^2 + m_{1i} z + m_{2i}} \quad (237)$$

the corresponding peaking factor,  $P_i$ , as defined in equation (16), can be determined in MATLAB as

```
N_i = [gamma0i gamma1i gamma2i]; D_i = [1 m1i m2i];
npoints = 1000;
Hi = freqz(N_i,D_i,npoints);
Hi_infty = max(abs(Hi));
Hi_2 = sqrt(sum(abs(Hi).^2)/npoints);
Pi = Hi_infty/Hi_2;
```

## Experiment 13.1

- Using the coefficient values provided in Table 16, one gets

$$\left. \begin{array}{l} P_1 = 11.2719 \\ P_2 = 13.4088 \\ P_3 = 12.5049 \end{array} \right\} \quad (238)$$

- If we scale the filter with the  $L_2$  norm, then, in order to minimize the maximum value of the output-noise PSD, one should change the section order to get a decreasing  $P_i$  sequence, as seen in Figure 60.

## Experiment 13.1

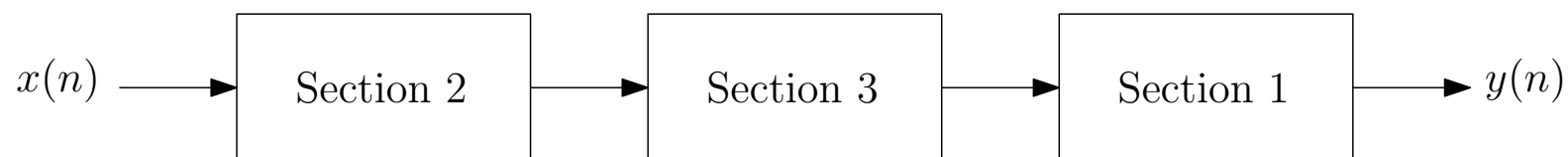


Figure 60: New section ordering in cascade filter to minimize peak value of output-noise PSD.

## Experiment 13.1

- Scaling the old Section 2 using  $L_2$  norm, we get

$$\lambda_2 = \frac{1}{\|F_2(z)\|_2} = \frac{1}{\|\frac{h_0}{D_2(z)}\|_2} \quad (239)$$

such that

```
D_2 = [1 m12 m22];
```

```
F_2 = freqz(h_0,D_2,npoints);
```

```
lambda_2 = 1/sqrt(sum(abs(F_2).^2)/npoints);
```

## Experiment 13.1

- Scaling the old Section 3, taking into account that it now comes after the old Section 2, we get

$$\lambda_3 = \frac{1}{\|H_2(z)F_3(z)\|_2} = \frac{1}{\left\| \frac{h_0 N_2(z)}{D_2(z)D_3(z)} \right\|_2} \quad (240)$$

such that

```
N_2 = [gamma02 gamma12 gamma22];
D_3 = [1 m13 m23];
D_2D_3 = conv(D_2,D_3);
H_2F_3 = freqz(h_0N_2,D_2D_3,npoints);
lambda_3 = 1/sqrt(sum(abs(H_2F_3).^2)/npoints);
```



## Experiment 13.1

- Finally, scaling the old Section 1, taking into account that it comes after the old Sections 2 and 3, we get

$$\lambda_1 = \frac{1}{\| H_2(z)H_3(z)F_1(z) \|_2} = \frac{1}{\| \frac{h_0 N_2(z)N_3(z)}{D_2(z)D_3(z)D_1(z)} \|_2} \quad (241)$$

such that

```
N_3 = [gamma03 gamma13 gamma23];
```

```
D_1 = [1 m11 m21];
```

```
N_2N_3 = conv(N_2,N_3);
```

```
D_2D_3D_1 = conv(D_2D_3,D_1);
```

```
H_2H_3F_1 = freqz(h_0N_2N_3,D_2D_3D_1,npoints);
```

```
lambda_1 = 1/sqrt(sum(abs(H_2H_3F_1).^2)/npoints);
```

yielding

$$\left. \begin{aligned} \lambda_2 = \text{lambda}_2 &= 522.2077 \\ \lambda_3 = \text{lambda}_3 &= 83.8126 \\ \lambda_1 = \text{lambda}_1 &= 12.1895 \end{aligned} \right\} \quad (242)$$