

# The Normal Lattice—A Cascade Digital Filter Structure

CASPER W. BARNES, SENIOR MEMBER, IEEE, AND SHU LEUNG

**Abstract**—A normal lattice is defined as a cascade of two-input two-output second-order normal form digital filter modules. An analysis of this structure is presented along with an iterative algorithm for constructing the system eigenvectors, a scaling algorithm, and an explicit expression for roundoff noise gain. These expressions and algorithms provide an explicit design procedure for realizing arbitrary rational transfer functions for single-input single-output filters.

Comparisons with other digital filter structures show that the normal lattice has very low roundoff noise and low coefficient sensitivity in realizations of narrow-band filters. This structure is also well adapted to the realization of complex filters.

## I. INTRODUCTION

In normal form realizations of digital filters, the basic building block for generating a conjugate pair of complex poles is the structure shown in Fig. 1. For the purposes of our discussion here

Manuscript received June 10, 1980; revised December 2, 1980 and July 2, 1981. This work was based on work supported by the National Science Foundation under Grant ENG 7818252 and by the Electronics Systems Group of Rockwell International Corporation under Contract ASWM-79470.

The authors are with the Electrical Engineering Department, University of California, Irvine, CA 92717.

93\$00.75 ©1982 IEEE

Fig. 1. Normal lattice of second-order module.



Fig. 2. Normal lattice structure.



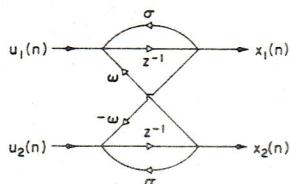


Fig. 1. The normal module.

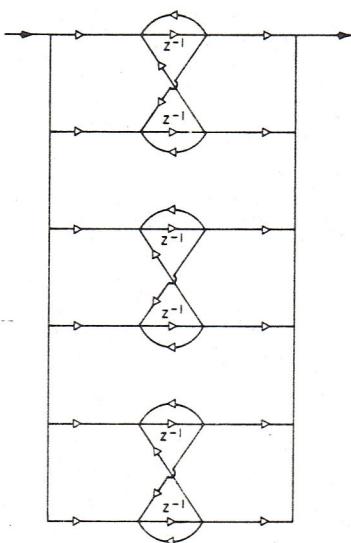


Fig. 2. Parallel structure of normal modules.

it will be convenient to regard this structure as two-input two-output module with state equation

$$\begin{pmatrix} x^1(n+1) \\ x^2(n+1) \end{pmatrix} = \begin{pmatrix} \sigma & \omega \\ -\omega & \sigma \end{pmatrix} \begin{pmatrix} x^1(n) \\ x^2(n) \end{pmatrix} + \begin{pmatrix} u_1(n) \\ u_2(n) \end{pmatrix} \quad (1)$$

and with matrix transfer function

$$H(z) = \frac{1}{(z-\sigma)^2 + \omega^2} \begin{pmatrix} z-\sigma & z\omega \\ -\omega & z-\sigma \end{pmatrix}. \quad (2)$$

As is well known, this structure will not support overflow limit cycles [1], has a uniform grid of possible pole locations for fixed-point realizations [2], and exhibits near optimal roundoff noise performance [3]–[6].

It has also been observed that this is a natural structure for realizing single complex poles in complex digital filters [7], with applications to complex demodulation. Thus if we interpret  $u_1(n)$  and  $u_2(n)$  as the real and imaginary parts of the input, and  $x^1(n)$  and  $x^2(n)$  as the real and imaginary parts of the state, then the complex state equation is

$$\begin{aligned} & [x^1(n+1) + ix^2(n+1)] \\ &= (\sigma - i\omega)[x^1(n) + ix^2(n)] + [u_1(n) + iu_2(n)] \quad (3) \end{aligned}$$

and the complex transfer function is

$$H_c(z) = \frac{1}{z - (\sigma - i\omega)}. \quad (4)$$

These modules can be combined in various ways to realize either real or complex digital filters. For example, a parallel structure of normal modules is shown in Fig. 2.

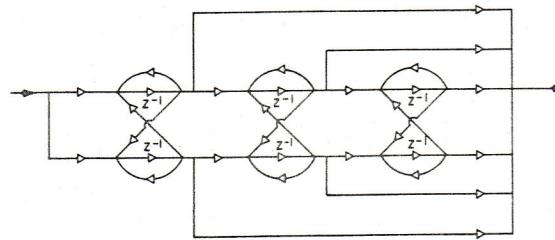


Fig. 3. Cascade structure of normal modules (normal lattice).

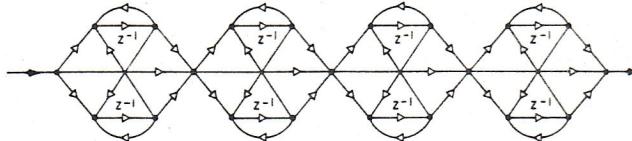


Fig. 4. Single-input single-output cascade of normal sections.

In this report we shall be concerned with single-input single-output filters realized as cascades of normal modules, as shown in Fig. 3. The distinctive feature of this cascade structure that distinguishes it from usual cascade realizations [8] is that internally it is a two-input two-output cascade, rather than a single-input single-output cascade. In the ensuing discussion, we shall refer to this structure as a normal lattice. It should be noted that scaled second-order normal filter sections require two input multipliers and two output multipliers and, therefore, naturally lend themselves to a two-input two-output cascade structure.

An alternative cascade structure of normal sections, proposed by Mullis and Roberts [9] and Jackson, Lindgren, and Kim [3], is the single-input single-output cascade shown in Fig. 4. We see that the normal lattice requires one fewer multiplier per stage than the single-input single-output cascade. An additional property of the normal lattice that simplifies design is that the optimal ordering of filter sections can be found by examining only pole pair orderings (rather than both pole pair and zero pair orderings, as required by single-input single-output cascades [8]) since the zeros of the filter are determined by the state variable tap weights feeding the output, as shown in Fig. 3.

In the following we shall provide a state variable description of the normal lattice, derive the design equations for realizing an arbitrary rational transfer function with simple poles, present an algorithm for scaling the structure, and derive an expression for the roundoff noise gain.

Finally, the roundoff noise performance of the normal lattice is compared with the noise performances of other filter structures. Over the range of parameters considered, the normal lattice exhibited the lowest noise gain. In addition, a specific narrow band filter design is used to illustrate the relatively low coefficient sensitivity of the normal lattice.

## II. STATE VARIABLE DESCRIPTION OF THE NORMAL LATTICE

We consider a general cascade of  $N$  normal modules, as shown in Fig. 5. The  $2 \times 2$  matrix transfer function for each normal module is

$$H_n = (zI - A_n)^{-1} \quad (5)$$

where

$$A_n = \begin{pmatrix} \sigma_n & \omega_n \\ -\omega_n & \sigma_n \end{pmatrix} \quad (6)$$

for  $n = 1, 2, \dots, N$ . The normal modules are separated by  $2 \times 2$  scaling matrices of the form

$$\mathbf{K}_n = \begin{pmatrix} K_n^{11} & K_n^{12} \\ K_n^{21} & K_n^{22} \end{pmatrix} \quad (7)$$

$n = 2, 3, \dots, N$ . For this general discussion we allow the scaling matrices to take the most general  $2 \times 2$  form. As we shall later show, complete scaling can be accomplished with diagonal scaling matrices, as shown in Fig. 3. For complex filters the scaling multipliers would have the form:  $K_n^{11} = K_n^{22}$  and  $K_n^{21} = -K_n^{12}$ .

The outputs of the normal modules (which are the outputs of the delay elements) form a convenient set of state variables with which to describe the behavior of the complete system in the standard form<sup>1</sup>

$$\begin{aligned} \mathbf{x}(n+1) &= \mathbf{A}\mathbf{x}(n) + \mathbf{b}u(n) \\ y(n) &= \mathbf{c}^\dagger \mathbf{x}(n) \end{aligned} \quad (8)$$

where

$$\mathbf{A} = \begin{pmatrix} A_1 & 0 & \cdots & 0 & 0 \\ K_2 & A_2 & \cdots & 0 & 0 \\ 0 & K_3 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & A_{N-1} & 0 \\ 0 & 0 & \cdots & K_N & A_N \end{pmatrix} \quad (9)$$

$$\mathbf{b} = \text{col}(\mathbf{b}_1, \mathbf{0}, \mathbf{0}, \dots, \mathbf{0}) \quad (10) \quad \text{and}$$

$$\mathbf{c} = \text{col}(\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N) \quad (11)$$

and

$$\mathbf{x}(n) = \text{col}(\mathbf{x}_1(n), \mathbf{x}_2(n), \dots, \mathbf{x}_N(n)) \quad (12) \quad \text{where}$$

where  $A_n$  and  $K_n$  are given by (6) and (7), respectively, and

$$\mathbf{b}_1 = \text{col}(\mathbf{b}_1^1, \mathbf{b}_1^2) \quad (13)$$

$$\mathbf{c}_k = \text{col}(\mathbf{c}_k^1, \mathbf{c}_k^2) \quad (14)$$

$$\mathbf{x}_k(n) = \text{col}(\mathbf{x}_k^1(n), \mathbf{x}_k^2(n)). \quad (15)$$

If this system has only simple poles, then the eigenvectors and reciprocal eigenvectors (left-handed eigenvectors) of the system  $\mathbf{A}$ -matrix form a bi-orthonormal pair of bases for the state space. In the following, we shall use these eigenvector bases to formulate the design equations, and to derive a general expression for the roundoff noise. It is notationally convenient to group the eigenvectors and reciprocal eigenvectors in pairs, as follows:

$$\mathbf{A}\phi_n^k = \lambda_n^k \phi_n^k, \quad n = 1, 2, \dots, N; k = 1, 2 \quad (16)$$

and

$$\psi_n^{k\dagger} \mathbf{A} = \lambda_n^k \psi_n^{k\dagger}, \quad n = 1, 2, \dots, N; k = 1, 2 \quad (17)$$

where

$$\lambda_n^1 = \sigma_n + i\omega_n, \quad n = 1, 2, \dots, N \quad (18)$$

$$\lambda_n^2 = \sigma_n - i\omega_n, \quad n = 1, 2, \dots, N \quad (19)$$

and where the eigenvectors and reciprocal eigenvectors are nor-

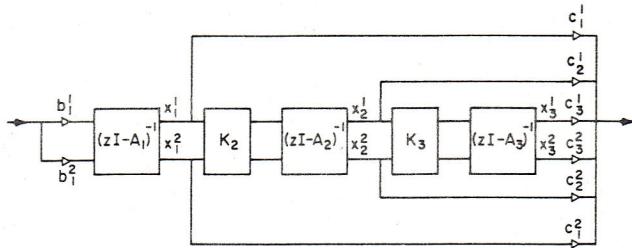


Fig. 5. Block diagram for normal lattice.

malized so that

$$\psi_n^{k\dagger} \phi_m^l = \delta_{nm} \delta_{kl}. \quad (20)$$

In the preceding equations  $\phi$  and  $\psi$  are column matrices of length  $2N$ . In Appendix A we describe a recursive procedure for constructing the eigenvectors and reciprocal eigenvectors for a normal cascade of any length; Fortran code is provided for the algorithm.

We can now expand the system matrix parameters on these bases as follows:

$$\mathbf{A} = \sum_{n=1}^N \sum_{k=1}^2 \lambda_n^k \phi_n^k \psi_n^{k\dagger} \quad (21)$$

$$\mathbf{b} = \sum_{n=1}^N \sum_{k=1}^2 \beta_n^k \phi_n^k \quad (22)$$

$$\mathbf{c}^\dagger = \sum_{n=1}^N \sum_{k=1}^2 \gamma_n^k \psi_n^{k\dagger} \quad (23)$$

$$\beta_n^k = \psi_n^{k\dagger} \mathbf{b} \quad (24)$$

and

$$\gamma_n^k = \mathbf{c}^\dagger \phi_n^k. \quad (25)$$

The transfer function of the normal lattice is then given by

$$\begin{aligned} H(z) &= \mathbf{c}^\dagger (zI - \mathbf{A})^{-1} \mathbf{b} \\ &= \sum_{n=1}^N \sum_{k=1}^2 \frac{\alpha_n^k}{z - \lambda_n^k} \end{aligned} \quad (26)$$

where

$$\alpha_n^k = \beta_n^k \gamma_n^k. \quad (27)$$

### Scaling

Scaling of the normal lattice can be accomplished by the following sequential procedure (refer to Fig. 5). We first choose  $b_1^1$  and  $b_1^2$  to scale the dynamic range of the state variable pair  $x_1^1$  and  $x_1^2$ . Then we choose  $K_2$  to scale the dynamic range of  $x_2^1$  and  $x_2^2$ , and so forth, proceeding sequentially to the last normal module in the lattice. The equations that describe the scaling relations are cumbersome, but the numerical process can be accomplished easily on a digital computer.

An algorithm for the  $L_2$ -scaling of a special class of normal lattices is described in Appendix B, and Fortran code for the algorithm is provided.

<sup>1</sup>Here and throughout the superscript dagger denotes matrix conjugate transpose.

### Design Procedure

A normal lattice realization of a rational transfer function can be obtained as follows.

- 1) Expand the transfer function into partial fractions as in (26).
- 2) Choose  $\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_N$ , as in (6), to obtain the required poles; i.e.,  $\lambda_n = \sigma_n + i\omega_n$ .
- 3) Choose  $\mathbf{b}_1, \mathbf{K}_2, \dots, \mathbf{K}_N$  to scale the dynamic ranges of the state variables (as shown in Appendix B).  $\mathbf{A}$  and  $\mathbf{b}$ , as given by (9) and (10), are now determined.
- 4) Calculate the reciprocal eigenvectors of  $\mathbf{A}$  (using the recursive procedure described in Appendix A).
- 5) Calculate  $\{\beta_n^k\}$  and  $\{\gamma_n^k\}$ , from (24) and (27), respectively.
- 6) Calculate  $\mathbf{c}^\dagger$  from (23).

### III. NOISE GAIN

The matrix transfer function from the state variable nodes to output is given by

$$\begin{aligned} F(z) &= z^{-1} \mathbf{c}^\dagger (z\mathbf{I} - \mathbf{A})^{-1} \\ &= z^{-1} \sum_{n=1}^N \sum_{k=1}^2 \frac{\gamma_n^k}{z - \lambda_n^k} \psi_n^\dagger. \end{aligned} \quad (28)$$

If we model the roundoff errors at each summing node as uncorrelated white noise, each with variance  $\sigma_0^2$ , then the variance of the roundoff error at the filter output is given by

$$\sigma_y^2 = (1 + \mu) \sigma_0^2 \quad (29)$$

where

$$\begin{aligned} \mu &= \frac{1}{2\pi i} \oint \frac{dz}{z} F(z) F^\dagger (1/z^*) \\ &= \sum_{n,m=1}^N \sum_{k,l=1}^2 \frac{\gamma_n^k \gamma_m^{l*}}{1 - \lambda_n^k \lambda_m^l} \psi_n^\dagger \psi_m^l. \end{aligned} \quad (30)$$

Equation (30) is suitable for numerical evaluation of the unit noise gain  $\mu$ , since values for all of the required quantities can be computed from the transfer function.

To determine the minimum noise gain, (30) must be evaluated for all possible pole-pair sequences. For each new pole sequence, the filter must be rescaled, and the reciprocal eigenvectors recomputed.

### IV. NUMERICAL COMPARISONS

To illustrate some of the properties of the normal lattice structure, and compare its performance with other structures, we

$\mathbf{A} =$

0.97536457	0.016554894	0.0	0.0	0.0	0.0	0.0	0.0	0.0
-0.016554894	0.97536457	0.0	0.0	0.0	0.0	0.0	0.0	0.0
-0.0035905249	0.0	0.98241036	0.043118061	0.0	0.0	0.0	0.0	0.0
0.0	0.053246011	-0.043118061	0.98241036	0.0	0.0	0.0	0.0	0.0
0.0	0.0	-0.026864306	0.0	0.99000867	0.057908279	0.0	0.0	0.0
0.0	0.0	0.0	0.046561198	-0.057908279	0.99000867	0.0	0.0	0.0
0.0	0.0	0.0	0.0	-0.016569622	0.0	0.99555249	0.06397247	0.99555249
0.0	0.0	0.0	0.0	0.0	0.023718661	-0.06397247	0.99555249	

(31)

consider the design of an eighth-order elliptic low-pass filter with a 0.1-dB passband ripple, and -72-dB stopband gain. In Fig. 6 we plot the filter roundoff noise gain as a function of passband

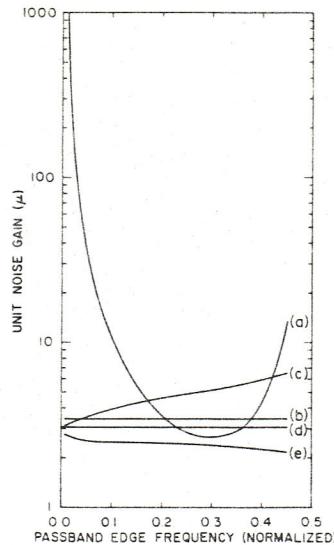


Fig. 6. Comparison of roundoff noise gains as a function of passband. Edge for five different realizations of the same eighth-order elliptic low-pass filter transfer function. Performance specifications: passband ripple = 0.1 dB; stopband ripple = -72 dB. (a) Direct form parallel realization. (b) Normal form parallel realization. (c) Normalized orthogonal polynomial structure. (d) Block-optimal parallel realization. (e) Normal lattice.

edge frequency for five different filter realizations—the direct form parallel [8], normal form parallel [4], the normalized orthogonal polynomial structure<sup>2</sup> described by Gray and Markel [10]–[12], block optimal parallel [9], and the normal lattice structure described in this paper. In all cases the filter was scaled for an  $L_2$  norm of unity for the transfer functions from the input to each internal summing node. The normal lattice was found to have the lowest noise gain over the range of filter bandwidths from 0.005 to 0.45.

In addition we find that for very narrow band filters the multiplier coefficients have reasonable values (i.e., neither large nor small compared to unity), and that the filter response is relatively insensitive to coefficient quantization. To illustrate these properties we consider an eighth-order elliptic low-pass filter designed to meet the following performance specifications:

$$\begin{aligned} \text{passband edge frequency} &= 0.01 \\ \text{stopband edge frequency} &= 0.013 \\ \text{passband ripple} &= 0.1 \text{ dB} \\ \text{stopband gain} &= -72 \text{ dB}. \end{aligned}$$

Following the design procedure given in Section II, we obtain the following matrix parameters for a normal lattice:

<sup>2</sup>The orthogonal polynomial filter noise gains were computed using a program written by T. Miyawaki.

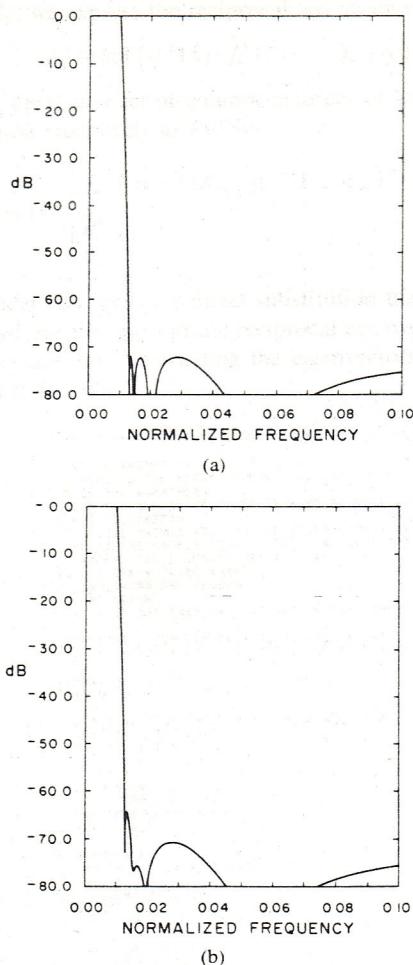


Fig. 7. Magnitude of transfer function of eighth-order low-pass filter realized as a normal lattice. (a) Multiplier coefficients unrounded. (b) Multiplier coefficients rounded to sixteen bits.

TABLE I  
COMPARISON OF ROUND OFF NOISE VARIANCE FOR  
TENTH-ORDER ELLIPTIC FILTER

REALIZATION	PARALLEL NORMAL	CASCADE NORMAL	NORMAL LATTICE
$\sigma^2/\sigma_0^2$ (dB)			
ROUNDING AFTER SUMMATION	7.69		7.16
$\sigma^2/\sigma_0^2$ (dB)			
ROUNDING BEFORE SUMMATION	14	16.9	13.7

$$\mathbf{b} = \text{col}(0.14743664, 0.27393837, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0) \quad (32)$$

$$\begin{aligned} \mathbf{c} = \text{col}(0.18535528, -0.099731184, 0.034813461, -0.070817758, \\ -0.06813572, -0.030977459, -0.033505959, 0.029857621) \end{aligned} \quad (33)$$

This design was scaled for an  $L_2$  norm of unity from the input to each of the state variable nodes. The computed noise gain from (30) was given by

$$\mu = 2.71.$$

Fig. 7 is a plot of the magnitude of the transfer function for unrounded multiplier coefficients, and with all multiplier coefficients rounded to 16 bits. For this example, we see that the effects of coefficient quantization are small.

In Table I we compare the noise performance of the normal lattice with the noise performance of a normal cascade (i.e., a single-input single-output cascade) for the tenth-order elliptic filter used by Jackson, Lindgren, and Kim [3], using their performance figure for the normal cascade. For this case, we see that the difference between the parallel normal and normal lattice realization is insignificant, whereas the normal lattice is about 3 dB better than the normal cascade.

## V. CONCLUDING REMARKS

Second-order normal modules can be combined in parallel or cascade to realize arbitrary rational transfer functions. The distinctive feature of the second-order normal module is that it is a two-input two-output structure. This leads to a unique form of cascade structure that we call the normal lattice. We have presented here an analysis of the normal lattice, an iterative algorithm for constructing the eigenvectors and reciprocal eigenvectors of the system matrix, an algorithm for scaling the structure, and an explicit expression for the roundoff noise gain. These relations provide an explicit procedure for designing digital filters in normal lattice form.

Filters realized in normal lattice form possess all of the desirable properties associated with the normal form; viz., no overflow limit cycles, low roundoff noise, and low coefficient sensitivity. They require the same number of multipliers, adders, and delays as the parallel normal form.

The noise performance of the normal lattice has been compared with the noise performances of other filter structures for a specific group of a low-pass filter transfer functions. Over the range of parameters considered, the normal lattice exhibited the lowest noise gain. In addition, a specific narrow-band low-pass filter design, realized as a normal lattice, was used to illustrate the relative insensitivity of this structure to coefficient quantization and to illustrate that, as in other low-noise structures, the multiplier coefficients do not take on extreme values for narrow-band filters.

## APPENDIX A

### RECURSIVE COMPUTATION OF THE EIGENVECTORS AND RECIPROCAL EIGENVECTORS OF A

We consider the system matrix given by (9)

Define

$$\begin{aligned} \mathbf{u}^1 &= 2^{-1/2} \text{col}(1, i) \\ \mathbf{u}^2 &= 2^{-1/2} \text{col}(1, -i). \end{aligned}$$

We observe that

$$\mathbf{A}_n \mathbf{u}^k = \lambda_n^k \mathbf{u}^k$$

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

We now express the eigenvectors of  $\mathbf{A}$  in the form

$$\phi_n^k = \text{col}\{\mathbf{e}_n^k(1), \mathbf{e}_n^k(2), \dots, \mathbf{e}_n^k(N)\}$$

where  $\{\mathbf{e}_n^k(m)\}$  is a set of column matrices of length 2 that can be computed recursively as follows:

$$\mathbf{e}_n^k(m) = \begin{cases} \mathbf{0}, & m < n \\ \mathbf{u}^k, & m = n \\ (\lambda_n^k \mathbf{I} - \mathbf{A}_m)^{-1} \mathbf{K}_m \mathbf{e}_n^k(m-1), & m > n. \end{cases}$$

The reader can verify by direct substitution that the vectors so constructed are the appropriate eigenvectors of  $\mathbf{A}$ .

Similarly, we express the reciprocal eigenvectors in the form

$$\psi_n^k = \text{col} \{ f_n^k(1), f_n^k(2), \dots, f_n^k(N) \}$$

where  $\{f_n^k(m)\}$  is a set of column matrices of length 2 that can be computed recursively as follows:

$$f_n^{k+}(m) = \begin{cases} f_n^{k+}(m+1)K_{m+1}(\lambda_n^k I - A_m)^{-1}, & m < n \\ u^{k+}, & m = n \\ 0, & m > n. \end{cases}$$

The reader can verify by direct substitution that the vectors so constructed are the appropriate reciprocal eigenvectors of  $A$ .

Fortran code for constructing the eigenvectors and reciprocal eigenvectors follows.

```

00100 C SUBROUTINE TO COMPUTE EIGENVECTORS AND RECIPROCAL EIGENVECTORS
00200 C
00300 C INPUT PARAMETERS :
00400 C   LAMBDA IS THE ARRAY STORING THE COORDINATES OF THE POLES
00500 C   A IS THE A-MATRIX
00600 C   N IS THE NUMBER OF POLES (N MUST BE EVEN AND LESS THAN 20)
00700 C
00800 C OUTPUT PARAMETERS :
00900 C   T IS THE MATRIX STORING THE EIGENVECTORS AS COLUMNS
01000 C   TPI IS THE MATRIX STORING THE CONJUGATE TRANSPOSE OF THE
01100 C   RECIPROCAL EIGENVECTORS AS ROWS
01200 C   EIG1 IS THE WORKING ARRAY
01300 C   IER INDICATES THE ERRORS :
01400 C     0= NO ERROR
01500 C     1= N IS NOT EVEN OR N IS NOT IN THE RANGE OF 2 TO 20
01600 C
01700 C SUBROUTINE EIGEN(LAMBDA,R,T,TPI,EIG1,N,IER)
01800 C COMPLEX LAMBDA(N),T(20,N),TPI(20,N),EIG1(N),POL,DET,DIF,
01900 C F1,F2
02000 C DIMENSION A(20,N)
02100 C
02200 C   ERROR CHECK
02300 C   IER=2
02400 C   IF(N-N/2.GT.0.OR.(N.LT.2.AND.N.GT.20)) IER=1
02500 C
02600 C   INITIALIZATION
02700 C   DO 5 J=1,N
02800 C   DO 5 I=1,N
02900 C   T(I,J)=CMPLX(0.,0.)
03000 C   COMPUTE EIGENVECTORS
03100 C   DO 50 I=1,N,2
03200 C   K=I
03300 C   POL=LAMBDA(K)
03400 C   COMPUTE THE I-TH EIGENVECTOR
03500 C   DO 10 J=1,N
03600 C   EIG1(J)=CMPLX(0.,0.)
03700 C   ASSIGN VALUE TO I-TH SECTION
03800 C   EIG1(J)=CMPLX(1.,0.)/SQRT(2.)
03900 C   K1=k+1
04000 C   EIG1(K1)=CMPLX(0.,1.)/SQRT(2.)
04100 C   IF(K1.EQ.N) GO TO 25
04200 C   K2=k-2
04300 C   COMPUTE RECURSIVELY FOR THOSE SECTIONS AFTER I-TH SECTION
04400 C   DO 20 L=K2,N,2
04500 C   SIGA(L,L)
04600 C   OMEGA(L,L+1)
04700 C   S1=A(L,L+2)
04800 C   S2=A(L+1,L+1)
04900 C   DIF=POL-SIG
05000 C   DET=DIF*DIF+OME*OME
05100 C   F1=EIG1(L-2)
05200 C   F2=EIG1(L-1)
05300 C   EIG1(L)=(DIF*S1+F1+F2*S2*OME)/DET
05400 C   L1=l+1
05500 C   EIG1(L1)=(-F1*OME*S1+F2*DIF*S2)/DET
05600 C   ASSIGN I-TH EIGENVECTOR TO THE I-TH COLUMN OF T-MATRIX
05700 C   DO 30 M=1,N
05800 C   T(M,K)=EIG1(M)
05900 C   T(M,K1)=CONJG(EIG1(M))
06000 C   CONTINUE
06100 C
06200 C   COMPUTE THE CONJUGATE TRANSPOSE RECIPROCAL EIGENVECTOR
06300 C   DO 40 L=N-1,1,-2
06400 C   K=I
06500 C   POL=LAMBDA(K)
06600 C   COMPUTE THE I-TH CONJUGATE TRANSPOSE OF RECIPROCAL EIGENVECTOR
06700 C   DO 70 J=1,N
06800 C   EIG1(J)=CMPLX(0.,0.)
06900 C   ASSIGN VALUE TO I-TH SECTION
07000 C   EIG1(K)=CMPLX(1.,0.)/SQRT(2.)
07100 C   K1=k+1
07200 C   EIG1(K1)=CMPLX(0.,-1.)/SQRT(2.)
07300 C   IF(K1.EQ.2) GO TO 90
07400 C   K2=k-2
07500 C   COMPUTE RECURSIVELY FOR THOSE SECTIONS AFTER I-TH SECTION
07600 C   DO 80 L=K2,1,-2
07700 C   SIGA(L,L)
07800 C   OMEGA(L,L+1)
07900 C   S1=A(L+2,L)
08000 C   S2=A(L+3,L+1)
08100 C   DIF=POL-SIG
08200 C   DET=DIF*DIF+OME*OME
08300 C   F1=EIG1(L+2)
08400 C   F2=EIG1(L+3)
08500 C   EIG1(L)=(F1*S1*DIF-F2*S2*OME)/DET
08600 C   L1=l+1
08700 C   EIG1(L1)=(F1*S1*OME+F2*S2*DIF)/DET
08800 C   ASSIGN THE I-TH CONJUGATE TRANSPOSE OF RECIPROCAL EIGENVECTOR
08900 C   TO THE I-TH ROW OF TPI-MATRIX
09000 C   DO 100 M=1,N
09100 C   TPI(K,M)=EIG1(M)
09200 C   TPI(K1,M)=CONJG(EIG1(M))
09300 C   CONTINUE
09400 C   RETURN
END

```

## APPENDIX B

### SCALING ALGORITHM

Each state variable summing node is  $L_2$ -scaled by choosing  $b_1, K_2, K_3, \dots, K_N$  such that for unit variance white noise input,

the variance of each state variable is equal to be unity. The transfer function from the input to the  $m$ th-state variable pair is given by

$$G_m(z) = (zI - A_m)^{-1} K_m (zI - A_{m-1})^{-1} \cdots K_2 (zI - A_2)^{-1} K_2 (zI - A_1)^{-1} b_1$$

where

$$b_1 = \begin{pmatrix} b_1^1 \\ b_1^2 \end{pmatrix}$$

and

$$K_m = \begin{pmatrix} b_m^1 & 0 \\ 0 & b_m^2 \end{pmatrix}, \quad \text{for } m = 2, \dots, N.$$

For unit variance white noise input, the covariance matrix of the  $m$ th-state variable pair is

$$\begin{aligned} \Lambda_m &= \frac{1}{2\pi i} \oint_{\text{unit circle}} \frac{dz}{z} G_m(z) G_m^\dagger(z) \\ &= \sum_j \sum_k \Lambda_m^{jk} \end{aligned}$$

where  $\Lambda_m^{jk}$  is the matrix residue of  $\Lambda_m$  at  $z = \lambda_j^k$ . This matrix residue can be expressed in the form

$$\Lambda_m^{jk} = H_m^{jk} K_m \Delta_{m-1}^{jk} K_m E_m^{jk}$$

where

$$H_m^{jk} = \begin{cases} (\lambda_j^k)^{-1} (\lambda_j^k I - A_m)^{-1}, & j \neq m \\ (z^{-1} (zI - A_m)^{-1} (z - \lambda_j^k))|_{z=\lambda_j^k}, & j = m \end{cases}$$

$$E_m^{jk} = ((\lambda_j^k)^{-1} I - A_m^\dagger)^{-1}$$

and

$$\Delta_{m-1}^{jk} = \begin{cases} (G_{m-1}(z) G_{m-1}^\dagger(z) (z - \lambda_j^k))|_{z=\lambda_j^k}, & j \neq m \\ G_{m-1}(\lambda_j^k) G_{m-1}^\dagger(\lambda_j^k), & j = m. \end{cases}$$

The diagonal elements of the matrix residues are given by

$$\begin{aligned} (\Lambda_m^{jk})_{11} &= (H_m^{jk})_{11} (E_m^{jk})_{11} (\Delta_{m-1}^{jk})_{11} (b_m^1)^2 \\ &\quad + (H_m^{jk})_{12} (E_m^{jk})_{11} (\Delta_{m-1}^{jk})_{21} b_m^1 b_m^2 \\ &\quad + (H_m^{jk})_{11} (E_m^{jk})_{21} (\Delta_{m-1}^{jk})_{12} b_m^1 b_m^2 \\ &\quad + (H_m^{jk})_{12} (E_m^{jk})_{21} (\Delta_{m-1}^{jk})_{22} (b_m^2)^2 \end{aligned}$$

and

$$\begin{aligned} (\Lambda_m^{jk})_{22} &= (H_m^{jk})_{21} (E_m^{jk})_{12} (\Delta_{m-1}^{jk})_{11} (b_m^1)^2 \\ &\quad + (H_m^{jk})_{22} (E_m^{jk})_{12} (\Delta_{m-1}^{jk})_{21} b_m^1 b_m^2 \\ &\quad + (H_m^{jk})_{21} (E_m^{jk})_{22} (\Delta_{m-1}^{jk})_{12} b_m^1 b_m^2 \\ &\quad + (H_m^{jk})_{22} (E_m^{jk})_{22} (\Delta_{m-1}^{jk})_{22} (b_m^2)^2. \end{aligned}$$

Summing up all of the residues and we obtain

$$\begin{aligned} (\Lambda_m)_{11} &= \sum_{j=1}^m \sum_{k=1}^2 (\Lambda_m^{jk})_{11} \\ &= \epsilon_m^1 (b_m^1)^2 + \xi_m^1 b_m^1 b_m^2 + \eta_m^1 (b_m^2)^2 \end{aligned}$$

$$(\Lambda_m^2)_{22} = \sum_{j=1}^m \sum_{k=1}^2 (\Lambda_m^{jk})_{22}$$

$$= \epsilon_m^2 (b_m^1)^2 + \xi_m^2 b_m^1 b_m^2 + \eta_m^2 (b_m^2)^2$$

where

$$\epsilon_m^1 = \sum_{j=1}^m \sum_{k=1}^2 (H_m^{jk})_{11} (E_m^{jk})_{11} (\Delta_{m-1}^{jk})_{11}$$

$$\xi_m^1 = \sum_{j=1}^m \sum_{k=1}^2 (H_m^{jk})_{12} (E_m^{jk})_{11} (\Delta_{m-1}^{jk})_{21}$$

$$+ (H_m^{jk})_{11} (E_m^{jk})_{21} (\Delta_{m-1}^{jk})_{12}$$

$$\eta_m^1 = \sum_{j=1}^m \sum_{k=1}^2 (H_m^{jk})_{12} (E_m^{jk})_{21} (\Delta_{m-1}^{jk})_{22}$$

$$\epsilon_m^2 = \sum_{j=1}^m \sum_{k=1}^2 (H_m^{jk})_{21} (E_m^{jk})_{12} (\Delta_{m-1}^{jk})_{11}$$

$$\xi_m^2 = \sum_{j=1}^m \sum_{k=1}^2 (H_m^{jk})_{22} (E_m^{jk})_{12} (\Delta_{m-1}^{jk})_{21}$$

$$+ (H_m^{jk})_{21} (E_m^{jk})_{22} (\Delta_{m-1}^{jk})_{12}$$

$$\eta^2 m = \sum_{j=1}^m \sum_{k=1}^2 (H_m^{jk})_{22} (E_m^{jk})_{22} (\delta^{jk} m - 1)_{22}.$$

For  $L_2$  scaling,

$$(\epsilon_m^1 + \epsilon_m^2)(b_m^1)^2 + (\xi_m^1 + \xi_m^2)b_m^1 b_m^2 + (\eta_m^1 + \eta_m^2)(b_m^2)^2 = 2$$

$$(\epsilon_m^1 - \epsilon_m^2)(b_m^1)^2 + (\xi_m^1 - \xi_m^2)b_m^1 b_m^2 + (\eta_m^1 - \eta_m^2)(b_m^2)^2 = 0.$$

We obtain

$$b_m^1 = f b_m^2$$

where

$$f = \frac{-(\xi_m^1 - \xi_m^2) \pm \sqrt{(\xi_m^1 - \xi_m^2)^2 - 4(\epsilon_m^1 - \epsilon_m^2)(\eta_m^1 - \eta_m^2)}}{2(\epsilon_m^1 - \epsilon_m^2)}$$

and thus

$$b_m^2 = \sqrt{\frac{2}{(\epsilon_m^1 + \epsilon_m^2)f^2 + (\xi_m^1 + \xi_m^2)f + (\eta_m^1 + \eta_m^2)}}.$$

repeat the algorithm for every section of the normal lattice to generate all of the scaling multipliers.

Fortran code for the scaling algorithm follows.

```

00100 C SCALING ROUTINE IS TO DETERMINE THE SCALING MULTIPLIERS
00200 C
00300 C INPUT PARAMETERS:
00400 C   POLE IS THE ARRAY STORING THE COORDINATES OF THE POLES
00500 C   N IS THE NUMBER OF POLES (N MUST BE EVEN AND LESS THAN 20)
00600 C   W1, W2, W3 ARE 2 BY 2 WORKING ARRAYS
00700 C
00800 C OUTPUT PARAMETERS:
00900 C   B IS THE ARRAY STORING THE SCALING MULTIPLIERS (IN POLE ORDER)
01000 C   IER INDICATES THE ERROR
01100 C   0=NO ERROR
01200 C   1=N IS NOT EVEN OR N IS NOT IN THE RANGE OF 2 TO 20
01300 C   SUBROUTINE CALLED : PHI ( TO EVALUATE THE COVARIANCE MATRIX )
01400 C
01500 C   SUBROUTINE SCALE(POLE,B,N,W1,W2,W3,IER)
01600 C   COMPLEX POLE(N),W1(2,2),W2(2,2),W3(2,2),D,DPI,
01700 C   1 ALPHA1,ALPHA2,BETA1,BETA2,GEMA1,GEMA2,A,C,E,G
01800 C   DIMENSION B(2,2)
01900 C   ERROR CHECK
02000 C   IER=0
02100 C   IF(N-N/2.GT.0.OR.(N.LT.2.AND.N.GT.20)) IER=1
02200 C   IF(IER.GT.0) RETURN
02300 C   DO 200 L=1,N/2
02400 C   1 INITIALIZATION
02500 C   L2=L*2
02600 C   L1=L2-1
02700 C   ALPHA1=CMPLX(0.,0.)
02800 C   ALPHA2=CMPLX(0.,0.)
02900 C   BETA1=CMPLX(0.,0.)
03000 C   BETA2=CMPLX(0.,0.)
03100 C   GEMA1=CMPLX(0.,0.)
03200 C   GEMA2=CMPLX(0.,0.)
03300 C   2 START TO COMPUTE THE SCALING PARAMETERS FOR EACH SECTION
03400 C   DO 100 I=1,L2
03500 C   NI=I
03600 C   W1(1,1)=CMPLX(1.,0.)
03700 C   W1(1,2)=CMPLX(1.,0.)
03800 C   W1(2,1)=CMPLX(1.,0.)
03900 C   W1(2,2)=CMPLX(1.,0.)
04000 C   I1=L-1

```

```

04100 C   IF(I1) 100,100,30
04200 C   30 COMPUTE THE (L-1)-TH COVARIANCE MATRIX AT NI-TH POLE
04300 C   CALL PHI(POLE,B,NI,I1,N,W1,W2,W3)
04400 C   100 IF(L-(NI+1)/2) 120,110,120
04500 C   110 THE DENOMINATOR OF H-MATRIX WHEN I=L
04600 C   D=(POLE(NI)-CONJG(POLE(NI))*POLE(NI))
04700 C   GO TO 140
04800 C   120 THE DENOMINATOR OF H-MATRIX WHEN I<L
04900 C   D=(1/POLE(NI)-POLE(L1))*(1/POLE(NI)-POLE(L2))
05000 C   140 SIGN=1/(POLE(L1)-POLE(L2))
05100 C   OMEGA=A*IMAG(POLE(L1))
05200 C   150 COMPUTE THE COMPONENTS OF H-MATRIX AND E-MATRIX
05300 C   A=(CMPLX(1.,0.)/POLE(NI))-CMPLX(SIGMA,0.))/D
05400 C   C=CMPLX(OMEGA,0.)/D
05500 C   E=(CMPLX(1.,0.)/POLE(NI))-CMPLX(SIGMA,0.))/DPI
05600 C   G=CMPLX(OMEGA,0.)/DPI
05700 C   160 COMPUTE THE RESIDUES
05800 C   ALPHA1=ALPHA1+A*EW1(1,1)
05900 C   ALPHA2=ALPHA2+C*EW1(1,1)
06000 C   170 BETAI=BETAI-C*EW1(2,1)+A*GW1(1,2)
06100 C   BETAZ=BETAZ-A*GW1(2,1)+C*EW1(1,2)
06200 C   GEMA1=GEMA1-C*EW1(2,2)
06300 C   GEMA2=GEMA2+A*EW1(2,2)
06400 C   CONTINUE
06500 C   COMPUTE SCALING MULTIPLIERS OF THE L-TH BLOCK
06600 C   P1=REAL(ALPHA1+ALPHA2)/2.
06700 C   P2=REAL(ALPHA1-ALPHA2)/2.
06800 C   Q1=REAL(BETAI-BETAZ)/2.
06900 C   Q2=REAL(BETAI+BETAZ)/2.
07000 C   R1=REAL(GEMA1+GEMA2)/2.
07100 C   R2=REAL(GEMA1-GEMA2)/2.
07200 C   FRACT=(SQRT(Q2**2-R1**2+R2**2)-Q2)/2./P2
07300 C   CONST=P1*FRACT**2+Q1*FRACT*R1
07400 C   B(L2)=1/SQRT(CONST)
07500 C   B(L1)=FRACT*B(L2)
07600 C   CONTINUE
07700 C   200 RETURN
07800 C   END
07900 C
08000 C
08100 C   SUBROUTINE TO EVALUATE L-TH COVARIANCE MATRIX AT NI-TH POLE
08200 C   TO FACILITATE THE SCALE ROUTINE
08300 C
08400 C   INPUT PARAMETERS :
08500 C   POLE IS THE ARRAY STORING THE COORDINATES OF THE POLES
08600 C   B IS THE ARRAY STORING THE SCALING MULTIPLIERS UP TO L-TH
08700 C   SECTION
08800 C   W1 IS THE 2 BY 2 MATRIX CONSISTING UNIT ELEMENTS
08900 C   NI INDICATES THE POLE AT WHICH COVARIANCE MATRIX EVALUATED
09000 C   L INDICATES L-TH COVARIANCE MATRIX
09100 C   OUTPUT PARAMETERS :
09200 C   W1 IS ALSO STORING THE EVALUATED VALUES OF THE L-TH COVARIANCE
09300 C   MATRIX AT NI-TH POLE
09400 C   W2, W3 ARE WORKING MATRIX OF 2 BY 2
09500 C   SUBROUTINE CALLED : CGMPRD (COMPLEX MATRIX MULTIPLICATION)
09600 C
09700 C   SUBROUTINE PHI(POLE,B,NI,L,N,W1,W2,W3)
09800 C   COMPLEX POLE(N),W1(2,2),W2(2,2),W3(2,2),D,DPI,S,O
09900 C   DIMENSION B(4)
01000 C
01100 C   DO 100 K=1,L
01200 C   K2=K*2
01300 C   K1=K2-1
01400 C   100 IF(K-(NI+1)/2) 20,10,20
01500 C   20 DENOMINATOR OF K-TH SECTION WHEN K>NI
01600 C   D=(POLE(NI)-CONJG(POLE(NI))*POLE(NI))
01700 C   GO TO 30
01800 C   30 DENOMINATOR OF K-TH SECTION WHEN K<NI
01900 C   D=(POLE(NI)-POLE(K1))*POLE(NI)-POLE(K2))
02000 C   DPI=(1/POLE(NI)-POLE(K1))*(1/POLE(NI)-POLE(K2))
02100 C   S=(POLE(K1)-POLE(K2))/2.
02200 C   O=CMPLX(A*IMAG(POLE(K1)),0.)
02300 C   EVALUATE THE TRASFER FUNCTION TO K-TH SECTION AT NI-TH POLE
02400 C   W2(1,1)=(POLE(NI)-S)/D*PK1)
02500 C   W2(1,2)=O/D*PK2)
02600 C   W2(2,1)=W2(1,2)*B(K1)/B(K2)
02700 C   W2(2,2)=W2(1,1)*B(K2)/B(K1)
02800 C   2 BY 2 MATRIX MULTIPLICATION TO CONCATENATE THE K-TH SECTION
02900 C   TO THE (K-1)-TH COVARIANCE MATRIX
03000 C   CALL CGMPRD(W2,W1,W3,2,2,2)
03100 C   EVALUATE THE TRANSPOSE OF THE TRANSFER FUNCTION TO K-TH
03200 C   SECTION AT NI-TH POLE
03300 C   W2(1,1)=(1/(POLE(NI)-S))/DPI*B(K1)
03400 C   W2(2,1)=O/DPI*B(K2)
03500 C   W2(1,2)=-W2(2,1)*B(K1)/B(K2)
03600 C   W2(2,2)=W2(1,1)*B(K2)/B(K1)
03700 C   2 BY 2 MULTIPLICATION TO EVALUATE THE K-TH COVARIANCE MATRIX
03800 C   CALL CGMPRD(W3,W2,W1,2,2)
03900 C   CONTINUE
04000 C   100 RETURN
04100 C   END
04200 C
04300 C
04400 C   SUBROUTINE CGMPRD IS TO MULTIPLY TWO GENERAL COMPLEX MATRIX
04500 C   TO FORM A RESULTANT GENERAL COMPLEX MATRIX
04600 C
04700 C   INPUT PARAMETERS:
04800 C   A IS THE NAME OF THE FIRST INPUT MATRIX
04900 C   B IS THE NAME OF THE SECOND INPUT MATRIX
05000 C   OUTPUT PARAMETERS:
05100 C   C IS THE NAME OF THE OUTPUT MATRIX
05200 C   N IS THE NUMBERS OF ROWS IN A
05300 C   M IS THE NUMBER OF COLUMNS IN A AND ROWS IN B
05400 C   L IS THE NUMBER OF COLUMNS IN B
05500 C
05600 C   REMARKS: NUMBER OF COLUMNS OF MATRIX A MUST BE EQUAL TO NUMBER
05700 C   OF ROWS OF MATRIX B AND N CANNOT BE GREATER THAN 20
05800 C
05900 C   SUBROUTINE CGMPRD(A,B,R,N,M,L)
06000 C   COMPLEX A(20,M),B(20,L),R(20,L)
06100 C
06200 C   DO 20 I=1,N
06300 C   DO 20 J=1,L
06400 C   INITIALIZATION
06500 C   R(I,J)=CMPLX(0.,0.)
06600 C   DO 20 K=1,M
06700 C   R(I,J)=R(I,J)+A(I,K)*B(K,J)
06800 C   CONTINUE
06900 C   RETURN
07000 C   END

```

## REFERENCES

- [1] C. W. Barnes and A. T. Fam, "Minimum norm recursive digital filters that are free of overflow limit cycles," *IEEE Trans. Circuits Syst.*, vol. CAS-24, pp. 569-574, Oct. 1977.
- [2] C. M. Rader and B. Gold, "Effects of parameter quantization on the poles of a digital filter," *Proc. IEEE*, vol. 55, pp. 688-689, May 1967.
- [3] L. B. Jackson, A. G. Lindgren, and Y. Kim, "Optimal synthesis of second-order state-space structures for digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-26, pp. 149-153, Mar. 1979.

- [4] C. W. Barnes, "Roundoff noise and overflow in normal digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-26, pp. 154-159, Mar. 1979.
- [5] W. L. Mills, C. T. Mullis, and R. A. Roberts, "Normal realizations of IIR digital filters," in *1979 IEEE Int. Conf. Acoustics, Speech, Signal Processing Rec.*, pp. 340-343, Apr. 1979.
- [6] M. Arjmand and R. A. Roberts, "Reduced multiplier, low roundoff noise digital filters," in *1979 IEEE Int. Conf. Acoustics, Speech, Signal Processing Rec.*, pp. 344-346, Apr. 1979.
- [7] T. H. Crystal and L. Erhman, "The design and application of digital filters with complex coefficients," *IEEE Trans. Audio Electroacoust.*, vol. AU-16, pp. 315-320, Sept. 1968.
- [8] L. B. Jackson, "Roundoff noise analysis for fixed-point digital filters realized in cascade or parallel form," *IEEE Trans. Audio Electroacoust.*, vol. AU-18, pp. 107-122, June 1970.
- [9] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, Sept. 1976.
- [10] A. H. Gray, Jr. and J. D. Markel, "A normalized digital filter structure," *IEEE Trans. Acoustics, Speech, Signal Processing*, vol. ASSP-23, pp. 268-277, June 1975.
- [11] J. D. Markel and A. H. Gray, Jr., "Roundoff noise characteristics of a class of orthogonal polynomial structures," *IEEE Trans. Acoustics, Speech, Signal Processing*, vol. ASSP-23, pp. 473-485, Oct. 1975.
- [12] J. D. Markel and A. H. Gray, Jr., "Fixed-point implementation algorithms for a class of orthogonal polynomial filter structures," *IEEE Trans. Acoustics, Speech, Signal Processing*, vol. ASSP-23, pp. 486-494, Oct. 1975.

## Conditions for Finiteness of a Constructive Algorithm for Determining Stability

JOHN N. TSITSIKLIS

**Abstract** — In earlier papers [1], [2], a methodology for deciding on the stability of a system of nonlinear differential equations was proposed. This methodology reduced the problem to a test of boundedness of all finite products of a given finite set of matrices. This paper investigates further the issue of whether the algorithm will terminate in finitely many steps and obtains some new conditions, concentrating on the case where the system being tested is unstable.

### I. INTRODUCTION

Brayton and Tong have proposed in [1] and [2] a constructive algorithm for determining the stability properties of a nonlinear time-invariant differential equation. They show that under certain conditions the problem may be reduced to the study of the "stability of a finite set of matrices", a concept to be defined below. They then proposed an algorithm that solves the latter problem and which is equivalent to constructing a Lyapunov function for the original differential equation. The range of potential applications of this algorithm is quite broad. It can be used, for example, in the study of switching systems, meaning linear systems such that the  $A$  matrix may undergo sudden changes (either random or deterministic).

Given the importance of this algorithm, it is natural to ask under what conditions it will be constructive (i.e., it terminates

Manuscript received May 11, 1981; revised November 13, 1981, and December 21, 1981.

The author is with the Laboratory for Information and Decision Systems and the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139.

after a finite number of steps) and is, therefore, implementable in a computer. Some sufficient conditions are given in [1] and [2]. In this paper we obtain more conditions for the finiteness of the algorithm.

The termination of the algorithm is very closely linked to the properties of subsets of  $\mathcal{C}^n$  that are left invariant under multiplication by a set of matrices. Most of the results in this paper involve characterizations of such invariant subsets.

**Definition 1:** A set of  $n \times n$  complex matrices  $\mathcal{Q}$  is stable if for every neighborhood of the origin  $U \subseteq \mathcal{C}^n$ , there exists another neighborhood of the origin  $V$  such that  $MV \subseteq U, \forall M \in \mathcal{Q}^*$ , where  $\mathcal{Q}^*$  is the semigroup of matrices generated by  $\mathcal{Q}$  (i.e., the set of all finite products of matrices in  $\mathcal{Q}$ ).

If  $W \subseteq \mathcal{C}^n$ , let  $\mathcal{H}[W]$  be the convex hull of  $W$ .

**Theorem 1:** Given a finite set  $\mathcal{Q} = \{M_0, M_1, \dots, M_{m-1}\}$  of  $m$  distinct  $n \times n$  complex matrices. Let  $W_0 \subseteq \mathcal{C}^n$  be a bounded, convex and symmetric polyhedral neighborhood of the origin. Define, for  $k > 0$

$$W_k = \mathcal{H} \left[ \bigcup_{i=0}^{\infty} M_j^i W_{k-i} \right] \quad (1.1)$$

where  $j \equiv (k-1)(\text{mod } m)$ . Then  $\mathcal{Q}$  is stable if and only if  $W^* \equiv \bigcup_{i=0}^{\infty} W_i$  is bounded.

*Proof:* See [1], p. 227.

Brayton and Tong's algorithm consists of constructing the set  $W^*$ . We should point out that if  $W$  is a convex polyhedron it is uniquely determined by the finite set  $E[W]$  of its extreme points. The algorithm is finite if and only if  $W^*$  can be constructed in a finite number of steps. In particular, a necessary condition for finiteness is that  $W^*$ , as well as  $W_i, i=1, 2, \dots$  are convex polyhedra. Concerning the first step of the algorithm (equation (1.1)), the following result is proved in [1]:

**Theorem 2:** Let  $M$  be a matrix whose eigenvalues have magnitude  $|\lambda(M)| < 1$ ; then, for any bounded neighborhood of the origin  $R$ , there exists some  $J$  such that  $\bigcup_{i=0}^{\infty} M^i R = \bigcup_{i=0}^J M^i R$ .

In order to complete the discussion, we must consider what happens if  $M$  is stable but has eigenvalues with  $|\lambda_i| = 1$ . In that case, it is necessary for finiteness that  $M$  has no eigenvalue equal to  $e^{i\theta} (\theta \in [0, 2\pi])$  with  $\theta/2\pi$  irrational. However, this condition is not sufficient. (For an example, see [2, p. 1124].)

Concerning the second step of the algorithm, we have, when  $\mathcal{Q}$  is stable [2]:

**Theorem 3:** If a set  $\mathcal{Q}$  of matrices is stable and there exists a positive  $k$  such that  $|\lambda_i(M)| \leq k < 1$  for all  $M \in \mathcal{Q}^*$ , then the algorithm will terminate "stable" in a finite number of steps.

Similarly with the case of the first step, it is easy to prove the following necessary condition for finiteness: If there exists some  $M \in \mathcal{Q}^*$ , and an eigenvalue  $\lambda$  of  $M$  such that  $\lambda = \exp[i\theta]$  and  $\theta/2\pi$  is irrational, then the algorithm is not finite.

Now define a set  $\Theta$  as follows:  $\theta \in \Theta$  if and only if  $0 \leq \theta \leq 2\pi$  and  $\exp[i\theta]$  is an eigenvalue of some  $M \in \mathcal{Q}^*$ . Then, if the algorithm terminates,  $\theta/2\pi$  has to be rational,  $\forall \theta \in \Theta$ . We also have the following result:

**Proposition 1:** If the algorithm is finite, then  $\Theta$  is a finite set.

*Proof:* Suppose that  $\theta \in \Theta$  and  $\theta/2\pi$  is rational and has been written as an irreducible ratio  $n/d$  of two integers. It is then easy to show that  $W^*$  has at least  $d$  extreme points. On the other hand, finiteness of the algorithm means that  $W^*$  has finitely many extreme points. Because there are only finitely many rational numbers in  $[0, 1]$  whose denominator is less or equal than a finite number, it follows that  $\Theta$  is a finite set. ■