

ROUND OFF NOISE MINIMIZATION IN TWO-DIMENSIONAL STATE-SPACE DIGITAL FILTERS USING ERROR FEEDBACK

Takao Hinamoto, Keisuke Higashi and Wu-Sheng Lu[†]

Graduate School of Engineering, Hiroshima University, Japan
hinamoto@hiroshima-u.ac.jp

[†]Dept. Elect. Comput. Eng., University of Victoria, Canada
wslu@ece.uvic.ca

ABSTRACT

This paper considers the problem of minimizing round-off noise in two-dimensional (2-D) state-space digital filters subject to L_2 -norm dynamic-range scaling constraints. The minimization will be achieved by using error feedback. Several techniques for the determination of the optimal full-scale, block-diagonal, diagonal, and scalar error-feedback matrices for a given 2-D state-space digital filter are proposed. A numerical example is presented to illustrate the utility of the proposed techniques.

I. INTRODUCTION

One of the primary finite-word-length (FWL) register effects in fixed-point digital filters is the roundoff noise caused by the rounding of products/summations within the realization. One can reduce the roundoff noise at the filter output using error feedback, which is achieved by extracting the quantization error after multiplications and additions, and then feeding the error signal back to a certain point through a simple circuit. Several techniques for error feedback have been presented in the past for 1-D digital filters [1]-[5], and more recently for 2-D digital filters [6]-[9].

This paper proposes several new algorithms for the reduction of roundoff noise in 2-D state-space digital filters. Several closed-form formulas for evaluating the optimal full-scale, block-diagonal, diagonal, and scalar error-feedback matrices for a given 2-D state-space digital filter are derived. A numerical example is presented to illustrate the algorithms proposed and to demonstrate their performance.

II. 2-D STATE-SPACE DIGITAL FILTERS WITH ERROR FEEDBACK

Consider the Roesser local state-space (LSS) model $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{m,n}$ which is stable, separately locally con-

trollable and separately locally observable:

$$\begin{aligned} \mathbf{x}_{11}(i, j) &= \mathbf{A}\mathbf{x}(i, j) + \mathbf{b}u(i, j) \\ y(i, j) &= \mathbf{c}\mathbf{x}(i, j) + du(i, j) \end{aligned} \quad (1)$$

where

$$\begin{aligned} \mathbf{x}_{11}(i, j) &= \begin{bmatrix} \mathbf{x}^h(i+1, j) \\ \mathbf{x}^v(i, j+1) \end{bmatrix}, \quad \mathbf{x}(i, j) = \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} \\ \mathbf{A} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}, \quad \mathbf{c} = [\mathbf{c}_1 \quad \mathbf{c}_2]. \end{aligned}$$

Here, $\mathbf{x}^h(i, j)$ is an $m \times 1$ horizontal state vector, $\mathbf{x}^v(i, j)$ is an $n \times 1$ vertical state vector, $u(i, j)$ is a scalar input, $y(i, j)$ is a scalar output, and $\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3, \mathbf{A}_4, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}_1, \mathbf{c}_2$, and d are real constant matrices of appropriate dimensions.

Carrying out the quantization before matrix-vector multiplication, an FWL implementation of (1) can be expressed as

$$\begin{aligned} \tilde{\mathbf{x}}_{11}(i, j) &= \mathbf{A}\mathbf{Q}[\tilde{\mathbf{x}}(i, j)] + \mathbf{b}u(i, j) \\ \tilde{y}(i, j) &= \mathbf{c}\mathbf{Q}[\tilde{\mathbf{x}}(i, j)] + du(i, j) \end{aligned} \quad (2)$$

where each component of $\mathbf{A}, \mathbf{b}, \mathbf{c}$, and d assumes an exact fractional B_c bit representation. The FWL local state vector $\tilde{\mathbf{x}}(i, j)$ and the output $\tilde{y}(i, j)$ all have a B bit fractional representation, while the input $u(i, j)$ is a $(B - B_c)$ bit fraction.

The quantizer $\mathbf{Q}[\cdot]$ in (2) rounds the B bit fraction $\tilde{\mathbf{x}}(i, j)$ to $(B - B_c)$ bits after multiplications and additions, where the sign bit is not counted. The quantization error

$$\mathbf{e}(i, j) = \tilde{\mathbf{x}}(i, j) - \mathbf{Q}[\tilde{\mathbf{x}}(i, j)] \quad (3)$$

coincides with the residue left in the lower part of $\tilde{\mathbf{x}}(i, j)$. The roundoff error $\mathbf{e}(i, j)$ is modeled as a zero-mean noise process of covariance $\sigma^2 \mathbf{I}_{m+n}$ with

$$\sigma^2 = \frac{1}{12} 2^{-2(B-B_c)}.$$

To reduce the filter's roundoff noise, the quantization error $\mathbf{e}(i, j)$ is fed back to each input of delay operators through an $(m+n) \times (m+n)$ constant matrix \mathbf{D} in the FWL filter (2). The 2-D filter with error feedback can be characterized by

$$\begin{aligned}\tilde{\mathbf{x}}_{11}(i, j) &= \mathbf{A}\mathbf{Q}[\tilde{\mathbf{x}}(i, j)] + \mathbf{b}u(i, j) + \mathbf{D}\mathbf{e}(i, j) \\ \tilde{y}(i, j) &= \mathbf{c}\mathbf{Q}[\tilde{\mathbf{x}}(i, j)] + du(i, j)\end{aligned}\quad (4)$$

where \mathbf{D} is referred to as an *error-feedback matrix*.

Subtracting (4) from (1) yields

$$\begin{aligned}\Delta\mathbf{x}_{11}(i, j) &= \mathbf{A}\Delta\mathbf{x}(i, j) + (\mathbf{A} - \mathbf{D})\mathbf{e}(i, j) \\ \Delta y(i, j) &= \mathbf{c}\Delta\mathbf{x}(i, j) + \mathbf{c}\mathbf{e}(i, j)\end{aligned}\quad (5)$$

where

$$\begin{aligned}\Delta\mathbf{x}(i, j) &= \mathbf{x}(i, j) - \tilde{\mathbf{x}}(i, j) \\ \Delta\mathbf{x}_{11}(i, j) &= \mathbf{x}_{11}(i, j) - \tilde{\mathbf{x}}_{11}(i, j) \\ \Delta y(i, j) &= y(i, j) - \tilde{y}(i, j).\end{aligned}$$

Let $\mathbf{G}_D(z_1, z_2)$ be the 2-D transfer function from the quantization error, $\mathbf{e}(i, j)$, to the filter output, $\Delta y(i, j)$. Then, we obtain

$$\mathbf{G}_D(z_1, z_2) = \mathbf{c}(\mathbf{Z} - \mathbf{A})^{-1}(\mathbf{A} - \mathbf{D}) + \mathbf{c} \quad (6)$$

where $\mathbf{Z} = z_1\mathbf{I}_m \oplus z_2\mathbf{I}_n$. The noise variance gain $I(\mathbf{D}) = \sigma_{out}^2/\sigma^2$ is then defined by

$$I(\mathbf{D}) = \text{tr}[\mathbf{W}_D] \quad (7)$$

where σ_{out}^2 denotes noise variance at the output, and

$$\mathbf{W}_D = \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{G}_D^*(z_1, z_2) \mathbf{G}_D(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2}$$

with $\Gamma_i = \{z_i : |z_i| = 1\}$ for $i = 1, 2$. By applying the 2-D Cauchy integral theorem, we obtain

$$\mathbf{W}_D = (\mathbf{A} - \mathbf{D})^T \mathbf{W}_o (\mathbf{A} - \mathbf{D}) + \mathbf{c}^T \mathbf{c} \quad (8)$$

where \mathbf{W}_o is called the *local observability Gramian* of the 2-D filter, and is defined by

$$\begin{aligned}\mathbf{W}_o &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} (\mathbf{Z}^* - \mathbf{A}^T)^{-1} \mathbf{c}^T \mathbf{c} (\mathbf{Z} - \mathbf{A})^{-1} \\ &\cdot \frac{dz_1 dz_2}{z_1 z_2} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{g}(i, j)^T \mathbf{g}(i, j).\end{aligned}\quad (9)$$

If there is no error feedback in the 2-D filter, then the noise variance gain $I(\mathbf{D})$ with $\mathbf{D} = \mathbf{0}$ becomes

$$\begin{aligned}I(\mathbf{0}) &= \text{tr}[\mathbf{A}^T \mathbf{W}_o \mathbf{A} + \mathbf{c}^T \mathbf{c}] \\ &= \text{tr}[\mathbf{W}_o].\end{aligned}\quad (10)$$

The l_2 -norm dynamic-range scaling constraints on the local state vector involves the *local controllability Gramian* of the 2-D filter, which is defined by

$$\begin{aligned}\mathbf{K}_c &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} (\mathbf{Z} - \mathbf{A}^T)^{-1} \mathbf{b} \mathbf{b}^T (\mathbf{Z}^* - \mathbf{A}^T)^{-1} \\ &\cdot \frac{dz_1 dz_2}{z_1 z_2} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{f}(i, j) \mathbf{f}(i, j)^T.\end{aligned}\quad (11)$$

The problem considered is to design the error-feedback matrix \mathbf{D} that minimizes (7), where matrix \mathbf{W}_D is specified by (8), subject to that all the diagonal elements of \mathbf{K}_c equal unity.

III. DETERMINATION OF OPTIMAL ERROR FEEDBACK MATRICES

In this section, we derive closed-form formulas for the determination of the optimal full-scale, block-diagonal, diagonal, and scalar error-feedback matrix \mathbf{D} to minimize $I(\mathbf{D}) = \text{tr}[\mathbf{W}_D]$ for a given 2-D state-space digital filter.

Case 1: \mathbf{D} is a general matrix

Substituting (8) into (7), we obtain

$$\begin{aligned}I(\mathbf{D}) &= \text{tr}[\mathbf{c}^T \mathbf{c} + (\mathbf{A} - \mathbf{D})^T \mathbf{W}_o (\mathbf{A} - \mathbf{D})] \\ &= \text{tr}[\mathbf{W}_o] + \text{tr}[\mathbf{D}^T \mathbf{W}_o \mathbf{D}] - 2 \text{tr}[\mathbf{D}^T \mathbf{W}_o \mathbf{A}].\end{aligned}\quad (12)$$

Differentiating (12) with respect to the error-feedback matrix \mathbf{D} yields

$$\frac{\partial I(\mathbf{D})}{\partial \mathbf{D}} = 2\mathbf{W}_o (\mathbf{D} - \mathbf{A}). \quad (13)$$

By choosing the error-feedback matrix as $\mathbf{D} = \mathbf{A}$, the noise gain $I(\mathbf{D})$ in (12) achieves its minimum value

$$\begin{aligned}I_{min}(\mathbf{D}) &= \text{tr}[\mathbf{W}_o] - \text{tr}[\mathbf{A}^T \mathbf{W}_o \mathbf{A}] \\ &= \text{tr}[\mathbf{c}^T \mathbf{c}].\end{aligned}\quad (14)$$

Case 2: \mathbf{D} is a block-diagonal matrix

In this case, matrix \mathbf{D} assumes the form

$$\mathbf{D} = \mathbf{D}_1 \oplus \mathbf{D}_4 \quad (15)$$

where \mathbf{D}_1 and \mathbf{D}_4 are $m \times m$ and $n \times n$ matrices, respectively, which leads (12) to

$$\begin{aligned}I(\mathbf{D}) &= \text{tr}[\mathbf{W}_o] + \text{tr}[\mathbf{D}_1^T \mathbf{W}_{o1} \mathbf{D}_1] + \text{tr}[\mathbf{D}_4^T \mathbf{W}_{o4} \mathbf{D}_4] \\ &\quad - 2 \text{tr}[\mathbf{D}_1^T (\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3)] \\ &\quad - 2 \text{tr}[\mathbf{D}_4^T (\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4)]\end{aligned}\quad (16)$$

where

$$\mathbf{W}_o = \begin{bmatrix} \mathbf{W}_{o1} & \mathbf{W}_{o2} \\ \mathbf{W}_{o3} & \mathbf{W}_{o4} \end{bmatrix}.$$

Letting $\partial I(\mathbf{D})/\partial \mathbf{D}_i = \mathbf{0}$ for $i = 1, 4$ yields

$$\begin{aligned} \mathbf{D}_1 &= \mathbf{A}_1 + \mathbf{W}_{o1}^{-1} \mathbf{W}_{o2} \mathbf{A}_3 \\ \mathbf{D}_4 &= \mathbf{A}_4 + \mathbf{W}_{o4}^{-1} \mathbf{W}_{o3} \mathbf{A}_2. \end{aligned} \quad (17)$$

By substituting (17) into (16), we obtain the minimum value of the noise variance gain $I(\mathbf{D})$ as

$$\begin{aligned} I_{min}(\mathbf{D}) &= \text{tr}[\mathbf{W}_o] - \text{tr}[\mathbf{D}_1^T (\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3)] \\ &\quad - \text{tr}[\mathbf{D}_4^T (\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4)]. \end{aligned} \quad (18)$$

Case 3: \mathbf{D} is a diagonal matrix

In this case, matrix \mathbf{D} assumes the form

$$\begin{aligned} \mathbf{D}_1 &= \text{diag}\{\alpha_1, \alpha_2, \dots, \alpha_m\} \\ \mathbf{D}_4 &= \text{diag}\{\beta_1, \beta_2, \dots, \beta_n\} \end{aligned} \quad (19)$$

which leads (16) to

$$\begin{aligned} I(\mathbf{D}) &= \text{tr}[\mathbf{W}_o] + \text{tr}[\mathbf{D}_1^2 \mathbf{W}_{o1}] + \text{tr}[\mathbf{D}_4^2 \mathbf{W}_{o4}] \\ &\quad - 2 \text{tr}[\mathbf{D}_1 (\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3)] \\ &\quad - 2 \text{tr}[\mathbf{D}_4 (\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4)]. \end{aligned} \quad (20)$$

This implies that if α_i 's and β_i 's satisfy

$$\begin{aligned} \alpha_i \left(\alpha_i - 2 \frac{(\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3)_{ii}}{(\mathbf{W}_{o1})_{ii}} \right) &< 0, \quad i = 1, 2, \dots, m \\ \beta_i \left(\beta_i - 2 \frac{(\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4)_{ii}}{(\mathbf{W}_{o4})_{ii}} \right) &< 0, \quad i = 1, 2, \dots, n \end{aligned} \quad (21)$$

then $I(\mathbf{D}) = \text{tr}[\mathbf{W}_D] < \text{tr}[\mathbf{W}_o]$. Letting $\partial I(\mathbf{D})/\partial \alpha_i = 0$ for $i = 1, 2, \dots, m$ and letting $\partial I(\mathbf{D})/\partial \beta_i = 0$ for $i = 1, 2, \dots, n$ gives

$$\begin{aligned} \alpha_i &= \frac{(\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3)_{ii}}{(\mathbf{W}_{o1})_{ii}}, \quad i = 1, 2, \dots, m \\ \beta_i &= \frac{(\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4)_{ii}}{(\mathbf{W}_{o4})_{ii}}, \quad i = 1, 2, \dots, n \end{aligned} \quad (22)$$

at which $I(\mathbf{D})$ achieves its minimum as

$$\begin{aligned} I_{min}(\mathbf{D}) &= \text{tr}[\mathbf{W}_o] - \sum_{i=1}^m \frac{(\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3)_{ii}^2}{(\mathbf{W}_{o1})_{ii}} \\ &\quad - \sum_{i=1}^n \frac{(\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4)_{ii}^2}{(\mathbf{W}_{o4})_{ii}} \end{aligned} \quad (23)$$

where $(\mathbf{A})_{ii}$ denotes the i th diagonal element of a square matrix \mathbf{A} .

Case 4: \mathbf{D}_1 and \mathbf{D}_4 are scalar matrices $\alpha \mathbf{I}_m$ and $\beta \mathbf{I}_n$

If $\mathbf{D}_1 = \alpha \mathbf{I}_m$ and $\mathbf{D}_4 = \beta \mathbf{I}_n$ with scalars α and β , then (20) becomes

$$\begin{aligned} I(\mathbf{D}) &= \text{tr}[\mathbf{W}_o] + \text{tr}[\mathbf{W}_{o1}] \alpha^2 + \text{tr}[\mathbf{W}_{o4}] \beta^2 \\ &\quad - 2 \text{tr}[\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3] \alpha \\ &\quad - 2 \text{tr}[\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4] \beta. \end{aligned} \quad (24)$$

Hence, if α and β satisfy

$$\begin{aligned} \alpha \left(\alpha - 2 \frac{\text{tr}[\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3]}{\text{tr}[\mathbf{W}_{o1}]} \right) &< 0 \\ \beta \left(\beta - 2 \frac{\text{tr}[\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4]}{\text{tr}[\mathbf{W}_{o4}]} \right) &< 0 \end{aligned} \quad (25)$$

then $I(\mathbf{D}) = \text{tr}[\mathbf{W}_D] < \text{tr}[\mathbf{W}_o]$. Moreover, from $\partial I(\mathbf{D})/\partial \alpha = 0$ and $\partial I(\mathbf{D})/\partial \beta = 0$, it follows that

$$\begin{aligned} \alpha &= \frac{\text{tr}[\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3]}{\text{tr}[\mathbf{W}_{o1}]} \\ \beta &= \frac{\text{tr}[\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4]}{\text{tr}[\mathbf{W}_{o4}]} \end{aligned} \quad (26)$$

which lead (24) to

$$\begin{aligned} I_{min}(\mathbf{D}) &= \text{tr}[\mathbf{W}_o] - \frac{(\text{tr}[\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3])^2}{\text{tr}[\mathbf{W}_{o1}]} \\ &\quad - \frac{(\text{tr}[\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4])^2}{\text{tr}[\mathbf{W}_{o4}]}. \end{aligned} \quad (27)$$

IV. A NUMERICAL EXAMPLE

Let a 2-D state-space digital filter $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{3,3}$ with $d = 0.0$ be described by

$$\begin{aligned} \mathbf{A}_1 &= \begin{bmatrix} 0.621553 & 0.014666 & -0.476979 \\ -0.081625 & 0.621548 & -0.181986 \\ 0.181983 & 0.476990 & 0.663600 \end{bmatrix} \\ \mathbf{A}_2 &= \begin{bmatrix} 0.059369 & -0.004829 & -0.024002 \\ -0.646852 & 0.061969 & 0.227715 \\ -0.229635 & 0.021958 & 0.076674 \end{bmatrix} \\ \mathbf{A}_3 &= \begin{bmatrix} 0.000378 & 0.000763 & 0.001503 \\ -0.000463 & -0.001501 & 0.000812 \\ -0.000021 & -0.000219 & 0.000908 \end{bmatrix} \\ \mathbf{A}_4 &= \begin{bmatrix} 0.620418 & 0.016504 & -0.479313 \\ -0.083124 & 0.620420 & -0.181961 \\ 0.181967 & 0.479315 & 0.661692 \end{bmatrix} \\ \mathbf{b}_1 &= [-0.007708 \quad 0.081835 \quad 0.028969]^T \\ \mathbf{b}_2 &= [-0.079883 \quad 0.846271 \quad 0.294745]^T \\ \mathbf{c}_1 &= [-0.766526 \quad 0.072050 \quad 0.267706] \\ \mathbf{c}_2 &= [-0.074064 \quad 0.007031 \quad 0.026238]. \end{aligned}$$

which is stable, separately locally controllable and separately locally observable. This corresponds to the *optimal realization* with minimum roundoff noise $I(\mathbf{0}) = 4.927082$, subject to the l_2 -norm dynamic-range scaling constraints.

In case \mathbf{D} is allowed to be a general matrix, then (13) suggests that we should choose $\mathbf{D} = \mathbf{A}$ which yields $I_{\min}(\mathbf{D}) = 0.670643$. Suppose the elements of matrix \mathbf{D} are rounded to power-of-two quantization with 3 bits after binary point (integer quantization), then the noise gain is given by

$$I(\mathbf{D}) = 0.700468 \quad (1.726719).$$

If \mathbf{D} is constrained to be a block-diagonal matrix, then the optimal $\mathbf{D} = \mathbf{D}_1 \oplus \mathbf{D}_4$ is calculated using (17), which gives

$$\mathbf{D}_1 = \begin{bmatrix} 0.621550 & 0.014658 & -0.476971 \\ -0.081619 & 0.621565 & -0.181992 \\ 0.181976 & 0.476968 & 0.663613 \end{bmatrix}$$

$$\mathbf{D}_4 = \begin{bmatrix} 0.618351 & 0.016703 & -0.478595 \\ -0.082737 & 0.620382 & -0.182106 \\ 0.181250 & 0.479384 & 0.661937 \end{bmatrix}$$

$$I_{\min}(\mathbf{D}_1 \oplus \mathbf{D}_4) = 1.331653.$$

After 3-bit quantization (integer quantization), this block-diagonal matrix $\mathbf{D} = \mathbf{D}_1 \oplus \mathbf{D}_4$ gives

$$I(\mathbf{D}_1 \oplus \mathbf{D}_4) = 1.351427 \quad (2.247670).$$

If \mathbf{D} is constrained to be a diagonal error-feedback matrix, then it can be calculated using (22) as

$$\mathbf{D}_1 = \text{diag}\{0.523405, 0.934410, 0.859620\}$$

$$\mathbf{D}_4 = \text{diag}\{0.521174, 0.934021, 0.858810\}$$

which yields $I_{\min}(\mathbf{D}) = 1.833208$. After 3-bit quantization (integer quantization), this diagonal matrix $\mathbf{D} = \mathbf{D}_1 \oplus \mathbf{D}_4$ yields

$$I(\mathbf{D}_1 \oplus \mathbf{D}_4) = 1.840195 \quad (2.247670).$$

If a scalar error-feedback matrix is calculated using (26), then we obtain

$$\alpha = 0.772479, \quad \beta = 0.771335$$

which yields $I_{\min}(\mathbf{D}) = 1.991329$. After 3-bit quantization (integer quantization), this scalar matrix \mathbf{D} results in

$$I(\mathbf{D}) = 1.993695 \quad (2.247670).$$

From these results, it is observed that the utilization of an optimal error feedback matrix leads to considerable reduction in roundoff noise, even when a scalar $\mathbf{D} = \alpha \mathbf{I}_m \oplus \beta \mathbf{I}_n$ with quantized α and β .

V. CONCLUSION

In this paper, the problem of minimizing roundoff noise in 2-D state-space digital filters has been investigated by means of error feedback. General, block-diagonal, diagonal, and scalar error-feedback matrices for minimizing the noise variance gain in a given 2-D state-space digital filter have been derived. Simulation results have been presented to illustrate the validity of our theoretical analysis and proposed algorithms.

1. REFERENCES

- [1] W. E. Higgins and D. C. Munson, "Noise reduction strategies for digital filters: Error spectrum shaping versus the optimal linear state-space formulation," *IEEE Trans. Acoust. Speech, Signal Processing*, vol. 30, pp. 963-973, Dec. 1982.
- [2] W. E. Higgins and D. C. Munson, "Optimal and suboptimal error-spectrum shaping for cascade-form digital filters," *IEEE Trans. Circuits Syst.*, vol. 31, pp. 429-437, May 1984.
- [3] T. I. Laakso and I. O. Hartimo, "Noise reduction in recursive digital filters using high-order error feedback," *IEEE Trans. Signal Processing*, vol. 40, pp. 1096-1107, May 1992.
- [4] P. P. Vaidyanathan, "On error-spectrum shaping in state-space digital filters," *IEEE Trans. Circuits Syst.*, vol. 32, pp. 88-92, Jan. 1985.
- [5] D. Williamson, "Roundoff noise minimization and pole-zero sensitivity in fixed-point digital filters using residue feedback," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 34, pp. 1210-1220, Oct. 1986.
- [6] T. Hinamoto, S. Karino and N. Kuroda, "Error spectrum shaping in 2-D digital filters," *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS'95)*, vol. 1, pp. 348-351, May 1995.
- [7] P. Agathoklis and C. Xiao, "Low roundoff noise structures for 2-D filters," *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS'96)*, vol. 2, pp. 352-355, May 1996.
- [8] T. Hinamoto, S. Karino and N. Kuroda, "2-D state-space digital filters with error spectrum shaping," *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS'96)*, vol. 2, pp. 766-769, May 1996.
- [9] T. Hinamoto, S. Karino, N. Kuroda and T. Kuma, "Error spectrum shaping in two-dimensional recursive digital filters," *IEEE Trans. Circuits Syst.*, vol. 46, pp. 1203-1215, Oct. 1999.