



JOINTLY OPTIMIZED ERROR FEEDBACK AND REALIZATION FOR ROUNDOFF NOISE MINIMIZATION IN TWO-DIMENSIONAL STATE-SPACE DIGITAL FILTERS

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

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

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

ABSTRACT

The minimization of roundoff noise subject to l_2 -norm dynamic-range scaling constraints in two-dimensional (2-D) state-space digital filters is considered by using joint error feedback and coordinate transformation optimization. An iterative approach for minimizing the roundoff noise under l_2 -norm dynamic-range scaling constraints is developed by jointly optimizing a scalar error-feedback matrix and a coordinate transformation matrix. A numerical example is presented to illustrate the utility of the proposed technique.

I. INTRODUCTION

When a given transfer function is realized through hardware implementation using fixed-point arithmetic, the internal noise caused by finite-word-length (FWL) registers may be the most serious concern to be dealt with. One of the primary FWL register effects in fixed-point digital filters is the roundoff noise caused by the rounding of products/summations within the realization. The synthesis of state-space digital filter structures with minimum roundoff noise under l_2 -norm dynamic-range scaling constraints has been investigated in two-dimensional (2-D) state-space digital filters [1],[2]. Another technique for the reduction of roundoff noise at the filter output is to use error feedback. Some techniques for error feedback have been presented in the past for 2-D digital filters [3]-[7].

In this paper, an iterative noise reduction technique for 2-D state-space digital filters is developed by jointly optimizing a scalar error feedback matrix and a coordinate transformation matrix. A numerical example is presented to illustrate the algorithm proposed and to demonstrate its performance.

Throughout the paper, the i th diagonal element of a square matrix \mathbf{A} is denoted by $(\mathbf{A})_{ii}$.

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 controllable and separately locally observable:

$$\begin{aligned} \mathbf{x}_{11}(i, j) &= \mathbf{Ax}(i, j) + \mathbf{bu}(i, j) \\ y(i, j) &= \mathbf{cx}(i, j) + du(i, j) \end{aligned} \quad (1)$$

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

$$\begin{aligned} \tilde{\mathbf{x}}_{11}(i, j) &= \mathbf{AQ}[\tilde{\mathbf{x}}(i, j)] + \mathbf{bu}(i, j) \\ \tilde{y}(i, j) &= \mathbf{cQ}[\tilde{\mathbf{x}}(i, j)] + du(i, j) \end{aligned} \quad (2)$$

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

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

III. JOINT OPTIMIZATION OF ERROR FEEDBACK AND COORDINATE TRANSFORMATION

An equivalent description of (1), $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{m+n}$, obtained via a coordinate transformation $\bar{\mathbf{x}}(k) = \mathbf{T}^{-1} \mathbf{x}(k)$ with $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ is characterized by

$$\begin{aligned} \bar{\mathbf{A}} &= \mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c} \mathbf{T} \\ \bar{\mathbf{W}}_o &= \mathbf{T}^T \mathbf{W}_o \mathbf{T}, \quad \bar{\mathbf{K}}_c = \mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T}. \end{aligned} \quad (12)$$

The purpose of this section is to investigate the joint optimization of scalar error-feedback matrices $\mathbf{D}_1 = \alpha \mathbf{I}_m$ and $\mathbf{D}_4 = \beta \mathbf{I}_n$, and coordinate transformation matrices \mathbf{T}_1 and \mathbf{T}_4 for roundoff noise minimization under l_2 -norm dynamic-range scaling constraints:

$$(\bar{\mathbf{K}}_c)_{ii} = (\mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T})_{ii} = 1, \quad i = 1, 2, \dots, n. \quad (13)$$

Let matrix \mathbf{W}_D in (8) with $\mathbf{D} = \alpha \mathbf{I}_m \oplus \beta \mathbf{I}_n$ and matrix \mathbf{K}_c in (11) be denoted by

$$\mathbf{W}_D = \begin{bmatrix} \mathbf{W}_{1\alpha} & \mathbf{W}_{\alpha\beta}^T \\ \mathbf{W}_{\alpha\beta} & \mathbf{W}_{4\beta} \end{bmatrix}, \quad \mathbf{K}_c = \begin{bmatrix} \mathbf{K}_{c1} & \mathbf{K}_{c2} \\ \mathbf{K}_{c3} & \mathbf{K}_{c4} \end{bmatrix},$$

respectively.

Under the joint application of scalar error-feedback and coordinate transformation, we minimize $I(\mathbf{D})$ (with α and β temporarily fixed) over an $m \times m$ nonsingular matrix \mathbf{T}_1 and an $n \times n$ nonsingular matrix \mathbf{T}_4 subject to the constraints stated in (13). To this end, we define the Lagrange function

$$\begin{aligned} J(\alpha, \beta, \mathbf{P}) &= \text{tr}[\mathbf{W}_{1\alpha} \mathbf{P}_1] + \lambda_1 (\text{tr}[\mathbf{K}_{c1} \mathbf{P}_1^{-1}] - m) \\ &\quad + \text{tr}[\mathbf{W}_{4\beta} \mathbf{P}_4] + \lambda_4 (\text{tr}[\mathbf{K}_{c4} \mathbf{P}_4^{-1}] - n) \end{aligned} \quad (14)$$

where $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$, $\mathbf{P}_i = \mathbf{T}_i \mathbf{T}_i^T$ for $i = 1, 4$, and λ_i for $i = 1, 4$ are Lagrange multipliers. We compute

$$\begin{aligned} \frac{\partial J(\alpha, \beta, \mathbf{P})}{\partial \mathbf{P}_1} &= \mathbf{W}_{1\alpha} - \lambda_1 \mathbf{P}_1^{-1} \mathbf{K}_{c1} \mathbf{P}_1^{-1} \\ \frac{\partial J(\alpha, \beta, \mathbf{P})}{\partial \mathbf{P}_4} &= \mathbf{W}_{4\beta} - \lambda_4 \mathbf{P}_4^{-1} \mathbf{K}_{c4} \mathbf{P}_4^{-1} \\ \frac{\partial J(\alpha, \beta, \mathbf{P})}{\partial \lambda_1} &= \text{tr}[\mathbf{K}_{c1} \mathbf{P}_1^{-1}] - m \\ \frac{\partial J(\alpha, \beta, \mathbf{P})}{\partial \lambda_4} &= \text{tr}[\mathbf{K}_{c4} \mathbf{P}_4^{-1}] - n. \end{aligned} \quad (15)$$

Let $\partial J(\alpha, \beta, \mathbf{P})/\partial \mathbf{P}_i = \mathbf{0}$ and $\partial J(\alpha, \beta, \mathbf{P})/\partial \lambda_i = 0$ for $i = 1, 4$. Then, it is derived that

$$\begin{aligned} \mathbf{P}_1 &= \sqrt{\lambda_1} \mathbf{W}_{1\alpha}^{-\frac{1}{2}} [\mathbf{W}_{1\alpha}^{\frac{1}{2}} \mathbf{K}_{c1} \mathbf{W}_{1\alpha}^{\frac{1}{2}}]^{\frac{1}{2}} \mathbf{W}_{1\alpha}^{-\frac{1}{2}} \\ \mathbf{P}_4 &= \sqrt{\lambda_4} \mathbf{W}_{4\beta}^{-\frac{1}{2}} [\mathbf{W}_{4\beta}^{\frac{1}{2}} \mathbf{K}_{c4} \mathbf{W}_{4\beta}^{\frac{1}{2}}]^{\frac{1}{2}} \mathbf{W}_{4\beta}^{-\frac{1}{2}} \\ \frac{1}{\sqrt{\lambda_1}} \text{tr}[\mathbf{K}_{c1} \mathbf{W}_{1\alpha}]^{\frac{1}{2}} &= \frac{1}{\sqrt{\lambda_1}} \left(\sum_{i=1}^m \mu_i \right) = m \quad (16) \\ \frac{1}{\sqrt{\lambda_4}} \text{tr}[\mathbf{K}_{c4} \mathbf{W}_{4\beta}]^{\frac{1}{2}} &= \frac{1}{\sqrt{\lambda_4}} \left(\sum_{i=1}^n \nu_i \right) = n \end{aligned}$$

where μ_i^2 for $i = 1, 2, \dots, m$ and ν_i^2 for $i = 1, 2, \dots, n$ are the eigenvalues of $\mathbf{K}_{c1} \mathbf{W}_{1\alpha}$ and $\mathbf{K}_{c4} \mathbf{W}_{4\beta}$, respectively. Therefore, we obtain

$$\begin{aligned} \mathbf{P}_1 &= \frac{1}{m} \left(\sum_{i=1}^m \mu_i \right) \mathbf{W}_{1\alpha}^{-\frac{1}{2}} [\mathbf{W}_{1\alpha}^{\frac{1}{2}} \mathbf{K}_{c1} \mathbf{W}_{1\alpha}^{\frac{1}{2}}]^{\frac{1}{2}} \mathbf{W}_{1\alpha}^{-\frac{1}{2}} \\ \mathbf{P}_4 &= \frac{1}{n} \left(\sum_{i=1}^n \nu_i \right) \mathbf{W}_{4\beta}^{-\frac{1}{2}} [\mathbf{W}_{4\beta}^{\frac{1}{2}} \mathbf{K}_{c4} \mathbf{W}_{4\beta}^{\frac{1}{2}}]^{\frac{1}{2}} \mathbf{W}_{4\beta}^{-\frac{1}{2}}. \end{aligned} \quad (17)$$

Substituting (17) into (14) yields the minimum value of $J(\alpha, \beta, \mathbf{P})$ for fixed α and β as

$$\min_{\mathbf{P}} J(\alpha, \beta, \mathbf{P}) = \frac{1}{m} \left(\sum_{i=1}^m \mu_i \right)^2 + \frac{1}{n} \left(\sum_{i=1}^n \nu_i \right)^2. \quad (18)$$

An iterative procedure for minimizing (14) with respect to scalar parameters λ_i for $i = 1, 4$, α and β as well as an $(m+n) \times (m+n)$ symmetric positive-definite matrix $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$ can be summarized as follows:

1) Set $i = 1$ and

$$\mathbf{P}(0) = \text{diag}\{(\mathbf{K}_c)_{11}, (\mathbf{K}_c)_{22}, \dots, (\mathbf{K}_c)_{m+n, m+n}\}.$$

2) Compute scalars $\alpha(i)$ and $\beta(i)$ using

$$\begin{aligned} \alpha(i) &= \frac{\text{tr}[(\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3) \mathbf{P}_1(i-1)]}{\text{tr}[\mathbf{W}_{o1} \mathbf{P}_1(i-1)]} \\ \beta(i) &= \frac{\text{tr}[(\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4) \mathbf{P}_4(i-1)]}{\text{tr}[\mathbf{W}_{o4} \mathbf{P}_4(i-1)]}. \end{aligned}$$

3) Compute $I_{min}(\alpha(i) \mathbf{I}_m \oplus \beta(i) \mathbf{I}_n) = (1 - \alpha(i)^2) \cdot \text{tr}[\mathbf{W}_{o1} \mathbf{P}_1(i-1)] + (1 - \beta(i)^2) \text{tr}[\mathbf{W}_{o4} \mathbf{P}_4(i-1)]$.

4) Replace $\mathbf{W}_{1\alpha}$ and $\mathbf{W}_{4\beta}$ by $\mathbf{W}_{1\alpha(i)}$ and $\mathbf{W}_{4\beta(i)}$ computed using

$$\begin{aligned} \mathbf{W}_{1\alpha(i)} &= (1 + \alpha(i)^2) \mathbf{W}_{o1} - \alpha(i) [(\mathbf{W}_{o1} \mathbf{A}_1 \\ &\quad + \mathbf{W}_{o2} \mathbf{A}_3)^T + \mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3] \\ \mathbf{W}_{4\beta(i)} &= (1 + \beta(i)^2) \mathbf{W}_{o4} - \beta(i) [(\mathbf{W}_{o4} \mathbf{A}_4 \\ &\quad + \mathbf{W}_{o3} \mathbf{A}_2)^T + \mathbf{W}_{o4} \mathbf{A}_4 + \mathbf{W}_{o3} \mathbf{A}_2], \end{aligned}$$

respectively.

5) Derive $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$ from (17), and take the resulting matrix $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$ as $\mathbf{P}(i) = \mathbf{P}_1(i) \oplus \mathbf{P}_4(i)$.

6) Compute $\text{tr}[\mathbf{W}_{1\alpha(i)} \mathbf{P}_1(i)] + \text{tr}[\mathbf{W}_{4\beta(i)} \mathbf{P}_4(i)]$.

7) Update $i := i + 1$ and repeat from Step 2) until the change in either $I(\alpha(i) \mathbf{I}_m \oplus \beta(i) \mathbf{I}_n)$ or $\text{tr}[\mathbf{W}_{1\alpha(i)} \mathbf{P}_1(i)] + \text{tr}[\mathbf{W}_{4\beta(i)} \mathbf{P}_4(i)]$ becomes insignificant compared to a prescribed tolerance.

From (17), the optimal coordinate transformation matrices \mathbf{T}_1 and \mathbf{T}_4 that minimize (14) can be obtained in closed form as

$$\begin{aligned} \mathbf{T}_1 &= \frac{1}{\sqrt{m}} \left(\sum_{i=1}^m \mu_i \right)^{\frac{1}{2}} \mathbf{W}_{1\alpha}^{-\frac{1}{2}} [\mathbf{W}_{1\alpha}^{\frac{1}{2}} \mathbf{K}_{c1} \mathbf{W}_{1\alpha}^{\frac{1}{2}}]^{\frac{1}{4}} \mathbf{U}_1 \\ \mathbf{T}_4 &= \frac{1}{\sqrt{n}} \left(\sum_{i=1}^n \nu_i \right)^{\frac{1}{2}} \mathbf{W}_{4\beta}^{-\frac{1}{2}} [\mathbf{W}_{4\beta}^{\frac{1}{2}} \mathbf{K}_{c4} \mathbf{W}_{4\beta}^{\frac{1}{2}}]^{\frac{1}{4}} \mathbf{U}_4 \end{aligned} \quad (19)$$

where \mathbf{U}_1 and \mathbf{U}_4 are $m \times m$ and $n \times n$ orthogonal matrices, respectively, which are obtained by applying a method reported in [7].

Once the optimal coordinate transformation matrix $\mathbf{T}(N) = \mathbf{T}_1(N) \oplus \mathbf{T}_4(N)$ is computed after N iterations, the diagonal error-feedback matrix $\mathbf{D} = \mathbf{D}_1 \oplus \mathbf{D}_4$ with $\mathbf{D}_1 = \text{diag}\{\alpha_1, \alpha_2, \dots, \alpha_m\}$ and $\mathbf{D}_4 = \text{diag}\{\beta_1, \beta_2, \dots, \beta_n\}$ that minimizes $I(\mathbf{D})$, (7), in a new realization characterized by (12) is given by

$$\begin{aligned} \alpha_i &= \frac{(\mathbf{T}_1^T(N)(\mathbf{W}_{o1} \mathbf{A}_1 + \mathbf{W}_{o2} \mathbf{A}_3) \mathbf{T}_1(N))_{ii}}{(\mathbf{T}_1^T(N) \mathbf{W}_{o1} \mathbf{T}_1(N))_{ii}} \\ \beta_i &= \frac{(\mathbf{T}_4^T(N)(\mathbf{W}_{o3} \mathbf{A}_2 + \mathbf{W}_{o4} \mathbf{A}_4) \mathbf{T}_4(N))_{ii}}{(\mathbf{T}_4^T(N) \mathbf{W}_{o4} \mathbf{T}_4(N))_{ii}}. \end{aligned} \quad (20)$$

This diagonal error-feedback matrix $\mathbf{D} = \mathbf{D}_1 \oplus \mathbf{D}_4$ leads to further reduction of the noise variance gain, i.e.,

$$I_{min}(\mathbf{D}) < I_{min}(\alpha(N) \mathbf{I}_m \oplus \beta(N) \mathbf{I}_n). \quad (21)$$

◀

▶

IV. A NUMERICAL EXAMPLE

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

Now we apply the iterative optimization procedure to the above *optimal realization*. The proposed algorithm converges after nine iterations to $\alpha = 0.845965$, $\beta = 0.845177$, and a coordinate transformation matrix $\mathbf{T}(9) = \mathbf{T}_1(9) \oplus \mathbf{T}_4(9)$ with

$$\begin{aligned}\mathbf{T}_1(9) &= \begin{bmatrix} -0.915091 & 0.162455 & 0.058237 \\ -0.470675 & 0.116512 & 0.581377 \\ -0.094509 & 0.398443 & 0.657685 \end{bmatrix} \\ \mathbf{T}_4(9) &= \begin{bmatrix} 0.157780 & -0.908860 & -0.069259 \\ 0.117342 & -0.463086 & -0.587994 \\ 0.403126 & -0.087756 & -0.657977 \end{bmatrix}\end{aligned}$$

which yield the noise variance gain $I(\alpha \mathbf{I}_3 \oplus \beta \mathbf{I}_3) = 1.665203$. After 3-bit quantization (integer quantization), this scalar error-feedback matrix gives $I_{min}(\alpha \mathbf{I}_3 \oplus \beta \mathbf{I}_3) = 1.670002$ (1.803185).

Next, a refined solution which offers further reduced noise variance gain is deduced by calculating an optimal diagonal error-feedback matrix for the optimal realization $(\mathbf{T}^{o-1} \mathbf{A} \mathbf{T}^o, \mathbf{T}^{o-1} \mathbf{b}, \mathbf{c} \mathbf{T}^o, d)_{3,3}$. In this case, the optimal diagonal error-feedback matrix $\mathbf{D} =$

$\mathbf{D}_1 \oplus \mathbf{D}_4$ is obtained using (20) as

$$\mathbf{D}_1 = \text{diag}\{0.310843, 0.931481, 0.773711\}$$

$$\mathbf{D}_4 = \text{diag}\{0.931135, 0.307381, 0.772593\}$$

which yields $I_{min}(\mathbf{D}) = 1.497455$. After 3-bit quantization (integer quantization), the above diagonal error-feedback matrix gives $I(\mathbf{D}) = 1.512761$ (1.629652).

VI. CONCLUSION

The minimization of roundoff noise in 2-D state-space digital filters has been investigated by means of joint optimization of error feedback/coordinate transformation. An iterative procedure for minimizing the roundoff noise in a 2-D digital filter has also been developed by jointly optimizing a scalar error-feedback matrix and a coordinate transformation matrix subject to the usual l_2 -norm dynamic-range scaling constraints. Simulation results have been presented to illustrate the validity of our proposed algorithm.

1. REFERENCES

- [1] M. Kawamata and T. Higuchi, "A unified study on the roundoff noise in 2-D state-space digital filters," *IEEE Trans. Circuits Syst.*, vol. 33, pp. 724-730, July 1986.
- [2] W.-S. Lu and A. Antoniou, "Synthesis of 2-D state-space fixed-point digital filter structures with minimum roundoff noise," *IEEE Trans. Circuits Syst.*, vol. 33, pp. 965-973, Oct. 1986.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] T. Hinamoto and H. Ohnishi, "Minimization of roundoff noise in state-space digital filters using error feedback and coordinate transformation," *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS'02)*, vol. 5, pp. 709-712, May 2002.