

JOINT OPTIMIZATION OF ERROR FEEDBACK AND COORDINATE TRANSFORMATION FOR ROUND-OFF NOISE MINIMIZATION IN 2-D STATE-SPACE DIGITAL FILTERS

Takao Hinamoto, Hiroaki Ohnishi and Wu-Sheng Lu[†]

Hiroshima University, Higashi-Hiroshima 739-8527, Japan. hinamoto@hiroshima-u.ac.jp

[†]University of Victoria, Victoria, BC, Canada V8W 3P6. wslu@ece.uvic.ca

ABSTRACT

The joint optimization of an error-feedback matrix and a coordinate-transformation matrix in 2-D state-space digital filters for roundoff noise minimization subject to L_2 -norm dynamic-range scaling constraints is investigated. Using linear-algebraic techniques, the problem at hand is converted into an unconstrained optimization problem, and the unconstrained problem obtained is then solved by applying an efficient quasi-Newton algorithm.

I. INTRODUCTION

When implementing IIR digital filters in fixed-point arithmetic, the problem of reducing the effects of roundoff noise (RN) at the filter output is of critical importance. Error feedback (EF) is a useful tool for reducing finite-word-length (FWL) effects in IIR digital filters. Many EF techniques have been reported for 2-D IIR digital filters [1]–[5]. Another useful approach is to construct the 2-D state-space filter structure for the RN gain to be minimized by applying a linear transformation to the state-space coordinates subject to L_2 -norm dynamic-range scaling constraints [6],[7]. As a natural extension of the fore-mentioned methods, efforts have been made to develop new methods that combine the EF and coordinate transformation for achieving better performance in the RN reduction. In [8], jointly-optimized iterative algorithms have also been developed for the filter with a scalar or general EF matrix.

This paper investigates the problem of jointly optimizing the EF and the coordinate transformation in 2-D state-space digital filters so as to minimize the RN subject to L_2 -norm dynamic-range scaling constraints. A jointly-optimized iterative technique, relying on an efficient quasi-Newton algorithm [9], is developed for RN minimization subject to the scaling constraints. The proposed technique can be applied to the cases where the error-feedback matrix is a scalar, diagonal, block-diagonal, or general matrix.

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

Consider the following local state-space (LSS) model $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{m,n}$ for 2-D IIR digital filters:

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

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 matrices of appropriate dimensions. The LSS model (1) is assumed stable, separately locally controllable and observable.

Due to finite register sizes, FWL constraints are imposed on the local state vector $\mathbf{x}(i, j)$, input, output, and coefficients in the realization $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{m,n}$. Assuming that the quantization is performed before matrix-vector multiplication, the actual FWL filter of (1) with EF can be implemented as

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

where \mathbf{D} is an $(m+n) \times (m+n)$ constant matrix referred to as *error feedback (EF) matrix*,

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

and each component of matrices $\mathbf{A}, \mathbf{b}, \mathbf{c}$, and d assumes an exact fractional B_c -bit representation. The FWL local state vector $\tilde{\mathbf{x}}(i, j)$ and 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 the multiplications and additions, where the sign bit is not counted. The quantization error $\mathbf{e}(i, j)$ is modeled as a zero-mean noise process of covariance $\sigma^2 \mathbf{I}_{m+n}$ with

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

Subtracting (2) from (1) yields

$$\begin{aligned} \Delta \mathbf{x}_{11}(i, j) &= \mathbf{A} \Delta \mathbf{x}(i, j) + (\mathbf{A} - \mathbf{D}) \mathbf{e}(i, j) \\ \Delta y(i, j) &= \mathbf{c} \Delta \mathbf{x}(i, j) + \mathbf{c} \mathbf{e}(i, j) \end{aligned} \quad (3)$$

where

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

The 2-D transfer function from the quantization error $\mathbf{e}(i, j)$ to the filter output $\Delta y(i, j)$ is given by

$$\mathbf{G}_D(z_1, z_2) = \mathbf{c}(\mathbf{Z} - \mathbf{A})^{-1}(\mathbf{A} - \mathbf{D}) + \mathbf{c}. \quad (4)$$

For the filter (3) with EF, the noise gain defined by $I(\mathbf{D}) = \sigma_{out}^2 / \sigma^2$ can be evaluated as

$$\begin{aligned} I(\mathbf{D}) &= \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{G}_D(z_1, z_2) \mathbf{G}_D^*(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \\ &= \text{tr}[\mathbf{W}_D] \end{aligned} \quad (5)$$

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 $i = 1, 2$. Utilizing the *2-D Cauchy integral theorem*, the matrix \mathbf{W}_D in (5) can be expressed in closed form

$$\mathbf{W}_D = (\mathbf{A} - \mathbf{D})^T \mathbf{W}_o (\mathbf{A} - \mathbf{D}) + \mathbf{c}^T \mathbf{c} \quad (6)$$

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

$$\mathbf{W}_o = \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} (\mathbf{Z}^* - \mathbf{A}^T)^{-1} \mathbf{c}^T \mathbf{c} (\mathbf{Z} - \mathbf{A})^{-1} \frac{dz_1 dz_2}{z_1 z_2}. \quad (7)$$

Matrix \mathbf{W}_o in (7) is referred to as the *unit noise matrix* for the 2-D filter (2) with $\mathbf{D} = \mathbf{0}$, and \mathbf{W}_D is viewed as the *unit noise matrix* for the 2-D filter (2) with EF specified by matrix \mathbf{D} . In the case where there is no

EF in the 2-D filter, the noise gain $I(\mathbf{D})$ with $\mathbf{D} = \mathbf{0}$ is expressed as

$$I(\mathbf{0}) = \text{tr}[\mathbf{A}^T \mathbf{W}_o \mathbf{A} + \mathbf{c}^T \mathbf{c}] = \text{tr}[\mathbf{W}_o]. \quad (8)$$

It is noted that the L_2 -norm scaling constraints on the local state vector $\mathbf{x}(i, j)$ involves the local controllability Gramian \mathbf{K}_c of the 2-D filter, defined by

$$\mathbf{K}_c = \frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} (\mathbf{Z} - \mathbf{A})^{-1} \mathbf{b} \mathbf{b}^T (\mathbf{Z}^* - \mathbf{A}^T)^{-1} \frac{dz_1 dz_2}{z_1 z_2}. \quad (9)$$

III. JOINT ERROR FEEDBACK AND REALIZATION OPTIMIZATION

A. Problem Statement

Applying a coordinate transformation defined by $\bar{\mathbf{x}}(i, j) = \mathbf{T}^{-1} \mathbf{x}(i, j)$ with $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ transforms the filter $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{m,n}$ to $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{m,n}$ where

$$\bar{\mathbf{A}} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c} \mathbf{T}. \quad (10)$$

The local controllability Gramian $\bar{\mathbf{K}}_c$ and local observability Gramian $\bar{\mathbf{W}}_o$ in the new realization then satisfy

$$\bar{\mathbf{K}}_c = \mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T}, \quad \bar{\mathbf{W}}_o = \mathbf{T}^T \mathbf{W}_o \mathbf{T}. \quad (11)$$

If the L_2 -norm dynamic-range scaling constraints

$$(\bar{\mathbf{K}}_c)_{ii} = (\mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T})_{ii} = 1, \quad i = 1, 2, \dots, m+n \quad (12)$$

are imposed on the new realization, then it is known that [16],[17]

$$\min_{\mathbf{T}} \text{tr}[\bar{\mathbf{W}}_o] = \frac{1}{m} \left(\sum_{i=1}^m \sigma_{1i} \right)^2 + \frac{1}{n} \left(\sum_{i=1}^n \sigma_{4i} \right)^2 \quad (13)$$

where σ_{1i}^2 for $i = 1, 2, \dots, m$ and σ_{4i}^2 for $i = 1, 2, \dots, n$ are the eigenvalues of matrices $\mathbf{K}_{1c} \mathbf{W}_{1o}$ and $\mathbf{K}_{4c} \mathbf{W}_{4o}$, respectively, and

$$\mathbf{K}_c = \begin{bmatrix} \mathbf{K}_{1c} & \mathbf{K}_{2c} \\ \mathbf{K}_{3c} & \mathbf{K}_{4c} \end{bmatrix}.$$

The LSS model $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{m,n}$ satisfying (12) and (13) simultaneously is known as the *optimal realization*.

If a coordinate transformation $\bar{\mathbf{x}}(i, j) = \mathbf{T}^{-1} \mathbf{x}(i, j)$ with $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ is applied to the LSS model (2), then the 2-D filter with EF can be characterized by

$$\begin{aligned} \tilde{\mathbf{x}}_{11}(i, j) &= \bar{\mathbf{A}} \mathbf{Q}[\tilde{\mathbf{x}}(i, j)] + \bar{\mathbf{b}} u(i, j) + \mathbf{D} \mathbf{e}(i, j) \\ \tilde{y}(i, j) &= \bar{\mathbf{c}} \mathbf{Q}[\tilde{\mathbf{x}}(i, j)] + d u(i, j). \end{aligned} \quad (14)$$

This corresponds to (2) in the original realization. In this case, the noise gain $I(\mathbf{D}, \mathbf{T})$ can be expressed as a function of matrices \mathbf{D} and $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ in the form

$$I(\mathbf{D}, \mathbf{T}) = \text{tr}[\bar{\mathbf{W}}_D] \quad (15)$$

where

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

The problem of RN minimization is to obtain matrices \mathbf{D} and $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ which minimize (15) subject to the scaling constraints in (12).

B. Problem Relaxation and Conversion

In order to reduce solution sensitivity, the objective function in (15) is modified to

$$J(\mathbf{D}, \mathbf{T}) = \text{tr}[(1 - \mu)\overline{\mathbf{W}}_D + \mu\overline{\mathbf{W}}_o] \quad (16)$$

where $0 \leq \mu \leq 1$ is a scalar that weights the importance of reducing $\text{tr}[\overline{\mathbf{W}}_o]$ relative to reducing $\text{tr}[\overline{\mathbf{W}}_D]$. Defining

$$\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4 = (\mathbf{T}_1 \oplus \mathbf{T}_4)^T (\mathbf{K}_{1c} \oplus \mathbf{K}_{4c})^{-\frac{1}{2}} \quad (17)$$

it follows that

$$\overline{\mathbf{K}}_c = \hat{\mathbf{T}}^{-T} \begin{bmatrix} \mathbf{I}_m & \mathbf{K}_{1c}^{-\frac{1}{2}} \mathbf{K}_{2c} \mathbf{K}_{4c}^{-\frac{1}{2}} \\ \mathbf{K}_{4c}^{-\frac{1}{2}} \mathbf{K}_{3c} \mathbf{K}_{1c}^{-\frac{1}{2}} & \mathbf{I}_n \end{bmatrix} \hat{\mathbf{T}}^{-1}. \quad (18)$$

This reduces the scaling constraints in (12) to

$$\begin{aligned} (\hat{\mathbf{T}}_1^{-T} \hat{\mathbf{T}}_1^{-1})_{ii} &= 1, & i &= 1, 2, \dots, m \\ (\hat{\mathbf{T}}_4^{-T} \hat{\mathbf{T}}_4^{-1})_{kk} &= 1, & k &= 1, 2, \dots, n. \end{aligned} \quad (19)$$

The constraints in (19) simply state that each column in matrices $\hat{\mathbf{T}}_1^{-1}$ and $\hat{\mathbf{T}}_4^{-1}$ must be a unity vector. These are satisfied if $\hat{\mathbf{T}}_1^{-1}$ and $\hat{\mathbf{T}}_4^{-1}$ assume the forms

$$\begin{aligned} \hat{\mathbf{T}}_1^{-1} &= \left[\frac{\mathbf{t}_{11}}{\|\mathbf{t}_{11}\|}, \frac{\mathbf{t}_{12}}{\|\mathbf{t}_{12}\|}, \dots, \frac{\mathbf{t}_{1m}}{\|\mathbf{t}_{1m}\|} \right] \\ \hat{\mathbf{T}}_4^{-1} &= \left[\frac{\mathbf{t}_{41}}{\|\mathbf{t}_{41}\|}, \frac{\mathbf{t}_{42}}{\|\mathbf{t}_{42}\|}, \dots, \frac{\mathbf{t}_{4n}}{\|\mathbf{t}_{4n}\|} \right] \end{aligned} \quad (20)$$

where \mathbf{t}_{1i} for $i = 1, 2, \dots, m$ and \mathbf{t}_{4j} for $j = 1, 2, \dots, n$ are $m \times 1$ and $n \times 1$ real vectors, respectively. In such a case, matrix $\overline{\mathbf{W}}_D$ in (15) can be written as

$$\overline{\mathbf{W}}_D = \hat{\mathbf{T}} [(\hat{\mathbf{A}} - \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T})^T \hat{\mathbf{W}}_o (\hat{\mathbf{A}} - \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T}) + \hat{\mathbf{C}}] \hat{\mathbf{T}}^T \quad (21)$$

where $\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4$ and

$$\begin{aligned} \hat{\mathbf{A}} &= (\mathbf{K}_{1c} \oplus \mathbf{K}_{4c})^{-\frac{1}{2}} \mathbf{A} (\mathbf{K}_{1c} \oplus \mathbf{K}_{4c})^{\frac{1}{2}} \\ \hat{\mathbf{C}} &= (\mathbf{K}_{1c} \oplus \mathbf{K}_{4c})^{\frac{1}{2}} \mathbf{c}^T \mathbf{c} (\mathbf{K}_{1c} \oplus \mathbf{K}_{4c})^{\frac{1}{2}} \\ \hat{\mathbf{W}}_o &= (\mathbf{K}_{1c} \oplus \mathbf{K}_{4c})^{\frac{1}{2}} \mathbf{W}_o (\mathbf{K}_{1c} \oplus \mathbf{K}_{4c})^{\frac{1}{2}}. \end{aligned}$$

Moreover, the objective function in (16) becomes

$$\begin{aligned} J(\mathbf{D}, \hat{\mathbf{T}}) &= (1 - \mu) \text{tr}[\hat{\mathbf{T}} (\hat{\mathbf{A}} - \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T})^T \\ &\quad \cdot \mathbf{W}_o (\hat{\mathbf{A}} - \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T}) \hat{\mathbf{T}}^T] \\ &\quad + (1 - \mu) \text{tr}[\hat{\mathbf{T}} \hat{\mathbf{C}} \hat{\mathbf{T}}^T] + \mu \text{tr}[\hat{\mathbf{T}} \hat{\mathbf{W}}_o \hat{\mathbf{T}}^T]. \end{aligned} \quad (22)$$

Therefore, the problem of obtaining matrices \mathbf{D} and $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$ that minimize (16) subject to the scaling constraints in (12) can be converted into an unconstrained optimization problem of obtaining matrices \mathbf{D} and $\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4$ that minimize (22).

C. Optimization Method

Let \mathbf{x} be the column vector that collects the variables in matrices \mathbf{D} and $\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4$. Then, $J(\mathbf{D}, \hat{\mathbf{T}})$ is a function of \mathbf{x} , denoted by $J(\mathbf{x})$. The algorithm starts with a trivial initial point \mathbf{x}_0 obtained from an initial assignment $\mathbf{D} = \hat{\mathbf{T}} = \mathbf{I}_{m+n}$. In the k th iteration, a quasi-Newton algorithm updates the most recent point \mathbf{x}_k to point \mathbf{x}_{k+1} as

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \quad (23)$$

where [9]

$$\begin{aligned} \mathbf{d}_k &= -\mathbf{S}_k \nabla J(\mathbf{x}_k), & \alpha_k &= \arg \min_{\alpha} J(\mathbf{x}_k + \alpha \mathbf{d}_k) \\ \mathbf{S}_{k+1} &= \mathbf{S}_k + \left(1 + \frac{\gamma_k^T \mathbf{S}_k \gamma_k}{\gamma_k^T \delta_k} \right) \frac{\delta_k \delta_k^T}{\gamma_k^T \delta_k} - \frac{\delta_k \gamma_k^T \mathbf{S}_k + \mathbf{S}_k \gamma_k \delta_k^T}{\gamma_k^T \delta_k} \\ \mathbf{S}_0 &= \mathbf{I}, & \delta_k &= \mathbf{x}_{k+1} - \mathbf{x}_k, & \gamma_k &= \nabla J(\mathbf{x}_{k+1}) - \nabla J(\mathbf{x}_k). \end{aligned}$$

Here, $\nabla J(\mathbf{x})$ is the gradients of $J(\mathbf{x})$ with respect to \mathbf{x} , and \mathbf{S}_k is a positive-definite approximation of the inverse Hessian matrix of $J(\mathbf{x})$. This iteration process continues until

$$|J(\mathbf{x}_{k+1}) - J(\mathbf{x}_k)| < \varepsilon \quad (24)$$

where $\varepsilon > 0$ is a prescribed tolerance. If the iteration is terminated at step k , \mathbf{x}_k is viewed as a solution point.

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

From (22), the optimal choice of \mathbf{D} is given by

$$\mathbf{D} = \hat{\mathbf{T}}^{-T} \hat{\mathbf{A}} \hat{\mathbf{T}}^T \quad (25)$$

which leads to

$$J(\hat{\mathbf{T}}^{-T} \hat{\mathbf{A}} \hat{\mathbf{T}}^T, \hat{\mathbf{T}}) = \text{tr}[\hat{\mathbf{T}} \{ (1 - \mu) \hat{\mathbf{C}} + \mu \hat{\mathbf{W}}_o \} \hat{\mathbf{T}}^T]. \quad (26)$$

Then, the elements in vector \mathbf{x} consist of $\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4$ and the gradients of $J(\mathbf{x})$ are found to be

$$\begin{aligned} \frac{\partial J(\mathbf{x})}{\partial t_{ij}} &= \lim_{\Delta \rightarrow 0} \frac{J(\hat{\mathbf{T}}_{ij}) - J(\hat{\mathbf{T}})}{\Delta} \\ &= 2e_j^T \hat{\mathbf{T}} [(1 - \mu) \hat{\mathbf{C}} + \mu \hat{\mathbf{W}}_o] \hat{\mathbf{T}}^T \hat{\mathbf{T}} g_{ij} \\ &\quad (1 \leq i, j \leq m) \text{ or } (m+1 \leq i, j \leq m+n). \end{aligned} \quad (27)$$

where $\hat{\mathbf{T}}_{ij}$ is the matrix obtained from $\hat{\mathbf{T}}$ with a perturbed (i, j) th component, and is given by [10]

$$\hat{\mathbf{T}}_{ij} = \hat{\mathbf{T}} + \frac{\Delta \hat{\mathbf{T}} g_{ij} e_j^T \hat{\mathbf{T}}}{1 - \Delta e_j^T \hat{\mathbf{T}} g_{ij}}$$

and \mathbf{g}_{ij} is computed using

$$\mathbf{g}_{ij} = \partial \left\{ \frac{\mathbf{t}_j}{\|\mathbf{t}_j\|} \right\} / \partial t_{ij} = \frac{1}{\|\mathbf{t}_j\|^3} (t_{ij} \mathbf{t}_j - \|\mathbf{t}_j\|^2 \mathbf{e}_i).$$

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

$$\mathbf{D} = \mathbf{D}_1 \oplus \mathbf{D}_4 \quad (28)$$

where \mathbf{D}_1 and \mathbf{D}_4 are $m \times m$ and $n \times n$ matrices, respectively. The gradients of $J(\mathbf{x})$ can be derived as

$$\begin{aligned} \frac{\partial J(\mathbf{x})}{\partial t_{ij}} &= 2\beta_1 + (1 - \mu)(\beta_2 - \beta_3) \\ \frac{\partial J(\mathbf{x})}{\partial d_{ij}} &= 2\mathbf{e}_j^T (1 - \mu) \hat{\mathbf{T}} \hat{\mathbf{W}}_o (\hat{\mathbf{T}}^T \mathbf{D} - \hat{\mathbf{A}} \hat{\mathbf{T}}^T) \mathbf{e}_i \end{aligned} \quad (29)$$

where

$$\begin{aligned} \beta_1 &= \mathbf{e}_j^T \hat{\mathbf{T}} [(1 - \mu)(\hat{\mathbf{A}}^T \hat{\mathbf{W}}_o \hat{\mathbf{A}} + \hat{\mathbf{C}}) + \mu \hat{\mathbf{W}}_o] \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij} \\ \beta_2 &= \mathbf{e}_j^T \hat{\mathbf{T}} \hat{\mathbf{W}}_o \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{D} \mathbf{D}^T \hat{\mathbf{T}} \mathbf{g}_{ij} \\ \beta_3 &= \mathbf{e}_j^T \hat{\mathbf{T}} (\hat{\mathbf{A}}^T \hat{\mathbf{W}}_o \hat{\mathbf{T}}^T \mathbf{D} + \hat{\mathbf{W}}_o \hat{\mathbf{A}} \hat{\mathbf{T}}^T \mathbf{D}^T) \mathbf{g}_{ij} \end{aligned}$$

with \mathbf{g}_{ij} defined in (27). Here, $d_{ij} \in \mathbf{D}_1 \oplus \mathbf{D}_4$ such that $d_{ij} \in \mathbf{D}_1$ for $(1, 1) \leq (i, j) \leq (m, m)$ and $d_{ij} \in \mathbf{D}_4$ for $(m+1, m+1) \leq (i, j) \leq (m+n, m+n)$.

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

$$\mathbf{D} = \text{diag}\{d_{11}, d_{22}, \dots, d_{m+n, m+n}\} \quad (30)$$

which leads to

$$\frac{\partial J(\mathbf{x})}{\partial d_{ii}} = 2\mathbf{e}_i^T (1 - \mu) \hat{\mathbf{T}} \hat{\mathbf{W}}_o (\hat{\mathbf{T}}^T \mathbf{D} - \hat{\mathbf{A}} \hat{\mathbf{T}}^T) \mathbf{e}_i \quad (31)$$

where $1 \leq i \leq m+n$. In this case, $\partial J(\mathbf{x}) / \partial t_{ij}$ is the same as in (29).

Case 4: $\mathbf{D}_1 = \alpha \mathbf{I}_m$ and $\mathbf{D}_4 = \beta \mathbf{I}_n$ with scalars α, β

The gradients of $J(\mathbf{x})$ can be calculated using

$$\begin{aligned} \frac{\partial J(\mathbf{x})}{\partial \alpha} &= 2\mathbf{e}_1^T (1 - \mu) \hat{\mathbf{T}} \hat{\mathbf{W}}_o (\hat{\mathbf{T}}^T \mathbf{D} - \hat{\mathbf{A}} \hat{\mathbf{T}}^T) \mathbf{e}_1 \\ \frac{\partial J(\mathbf{x})}{\partial \beta} &= 2\mathbf{e}_{m+1}^T (1 - \mu) \hat{\mathbf{T}} \hat{\mathbf{W}}_o (\hat{\mathbf{T}}^T \mathbf{D} - \hat{\mathbf{A}} \hat{\mathbf{T}}^T) \mathbf{e}_{m+1} \end{aligned} \quad (32)$$

and $\partial J(\mathbf{x}) / \partial t_{ij}$ is computed using (29).

IV. CONCLUSION

The joint optimization of a error feedback matrix and a coordinate-transformation matrix in 2-D state-space digital filters for roundoff noise minimization subject to L_2 -norm dynamic-range scaling constraints has been investigated. It has been shown that the problem at

hand can be converted into an unconstrained optimization problem by using linear algebraic techniques. An efficient quasi-Newton algorithm has been employed to solve the unconstrained optimization problem iteratively. It has been clarified that the proposed technique can be applied to the cases where the error feedback matrix is a scalar, diagonal, block-diagonal, or general matrix. Our computer simulation results have demonstrated the effectiveness of the proposed technique compared with the existing method.

1. REFERENCES

- [1] 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.
- [2] 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.
- [3] 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.
- [4] T. Hinamoto, N. Kuroda and T. Kuma, "Error feedback for noise reduction in 2-D digital filters with quadrantally symmetric or antisymmetric coefficients," *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS'97)*, vol. 4, pp. 2461-2464, June 1997.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] T. Hinamoto, K. Higashi and W.-S. Lu, "Separate/joint optimization of error feedback and coordinate transformation for roundoff noise minimization in two-dimensional state-space digital filters," *IEEE Trans. Signal Processing*, vol. 51, 2003 (to appear).
- [9] R. Fletcher, *Practical Methods of Optimization*, 2nd ed. Wiley, New York, 1987.
- [10] T. Kailath, *Linear Systems*, Prentice Hall, 1980.