

Roundoff Noise Minimization for a Class of 2-D State-Space Digital Filters Using Separate/Joint Optimization of Error Feedback and Realization

Takao Hinamoto
Graduate School of Engineering
Hiroshima University
Higashi-Hiroshima 739-8527, Japan
Email: hinamoto@hiroshima-u.ac.jp

Toru Oumi
Graduate School of Engineering
Hiroshima University
Higashi-Hiroshima 739-8527, Japan
Email: oumi@hiroshima-u.ac.jp

Wu-Sheng Lu
Dept. of Elec. and Comp. Engineering
University of Victoria
Victoria, BC, Canada V8W 3P6
Email: wslu@ece.uvic.ca

Abstract—Techniques for the separate/joint optimization of error-feedback and realization are developed to minimize the roundoff noise subject to L_2 -norm dynamic-range scaling constraints for a class of 2-D state-space digital filters. In the joint optimization, the problem at hand is converted into an unconstrained optimization problem by using linear-algebraic techniques. The unconstrained problem obtained is then solved by applying an efficient quasi-Newton algorithm. A numerical example is presented to illustrate the utility of the proposed techniques.

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 proposed 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 EF and coordinate transformation for better performance in the RN reduction. In [8], separately/jointly-optimized iterative algorithms have been developed for 2-D filters with EF matrix.

This paper investigates the problems of separately/jointly optimizing EF and realization subject to L_2 -norm dynamic-range scaling constraints for 2-D state-space digital filters described by the Fornasini-Marchesini second model. The former is solved analytically and the latter iteratively by applying an efficient quasi-Newton algorithm [9]. Computer simulation results demonstrate the validity of the proposed techniques.

II. ROUND OFF NOISE ANALYSIS AND SCALING

Consider a 2-D IIR digital filter that is described by the Fornasini-Marchesini second local state-space (LSS) model $(\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}, d)_n$:

$$\begin{aligned} \mathbf{x}(i, j) &= \mathbf{A}_1 \mathbf{x}(i-1, j) + \mathbf{A}_2 \mathbf{x}(i, j-1) \\ &\quad + \mathbf{b}_1 u(i-1, j) + \mathbf{b}_2 u(i, j-1) \\ y(i, j) &= \mathbf{c} \mathbf{x}(i, j) + du(i, j) \end{aligned} \quad (1)$$

where $\mathbf{x}(i, j)$ is an $n \times 1$ local state vector, $u(i, j)$ is a scalar input, $y(i, j)$ is a scalar output, and $\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}$, and d are real matrices of appropriate dimensions. The LSS model in (1) is assumed to be stable, locally controllable and locally observable.

Due to finite register sizes, FWL constraints are imposed on the local state vector, input, output, and coefficients in the realization $(\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}, d)_n$. Assuming that the quantization is carried out before matrix-vector multiplication, the actual 2-D FWL filter of (1) with EF and error feedforward can be implemented as

$$\begin{aligned} \tilde{\mathbf{x}}(i, j) &= \mathbf{A}_1 \mathbf{Q}[\tilde{\mathbf{x}}(i-1, j)] + \mathbf{A}_2 \mathbf{Q}[\tilde{\mathbf{x}}(i, j-1)] \\ &\quad + \mathbf{b}_1 u(i-1, j) + \mathbf{b}_2 u(i, j-1) \\ &\quad + \mathbf{D}_1 \mathbf{e}(i-1, j) + \mathbf{D}_2 \mathbf{e}(i, j-1) \\ \tilde{y}(i, j) &= \mathbf{c} \mathbf{Q}[\tilde{\mathbf{x}}(i, j)] + du(i, j) + \mathbf{h} \mathbf{e}(i, j) \end{aligned} \quad (2)$$

where \mathbf{D}_1 and \mathbf{D}_2 are $n \times n$ EF matrices, \mathbf{h} is a $1 \times n$ error-feedforward vector,

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

and each component of matrices $\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}$ 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 $\sigma^2 \mathbf{I}_n$ with

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

Subtracting (2) from (1) yields

$$\begin{aligned} \Delta \mathbf{x}(i, j) &= \mathbf{A}_1 \Delta \mathbf{x}(i-1, j) + \mathbf{A}_2 \Delta \mathbf{x}(i, j-1) \\ &\quad + (\mathbf{A}_1 - \mathbf{D}_1) \mathbf{e}(i-1, j) \\ &\quad + (\mathbf{A}_2 - \mathbf{D}_2) \mathbf{e}(i, j-1) \\ \Delta y(i, j) &= \mathbf{c} \Delta \mathbf{x}(i, j) + (\mathbf{c} - \mathbf{h}) \mathbf{e}(i, j) \end{aligned} \quad (3)$$

where

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

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

$$\begin{aligned} G(z_1, z_2) &= \mathbf{c} (\mathbf{I}_n - z_1^{-1} \mathbf{A}_1 - z_2^{-1} \mathbf{A}_2)^{-1} \\ &\quad \cdot [z_1^{-1} (\mathbf{A}_1 - \mathbf{D}_1) + z_2^{-1} (\mathbf{A}_2 - \mathbf{D}_2)] + \mathbf{c} - \mathbf{h}. \end{aligned} \quad (4)$$

For the 2-D filter in (2), the noise gain $I(\mathbf{D}_1, \mathbf{D}_2, \mathbf{h}) = \sigma_{out}^2 / \sigma^2$ can be evaluated by

$$I(\mathbf{D}_1, \mathbf{D}_2, \mathbf{h}) = \text{tr} \left[\frac{1}{(2\pi j)^2} \oint_{\Gamma_1} \oint_{\Gamma_2} \mathbf{G}^*(z_1, z_2) \mathbf{G}(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \right] \quad (5)$$

where σ_{out}^2 denotes noise variance at the output, and $\Gamma_i = \{z_i : |z_i| = 1\}$ for $i = 1, 2$.

Let the transition matrix $\mathbf{A}^{(i, j)}$ be defined by

$$(\mathbf{I}_n - z_1^{-1} \mathbf{A}_1 - z_2^{-1} \mathbf{A}_2)^{-1} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{A}^{(i, j)} z_1^{-i} z_2^{-j} \quad (6)$$

where $i, j \geq 0$. Then the following properties holds:

$$\begin{aligned} \mathbf{A}^{(0, 0)} &= \mathbf{I}_n, \quad \mathbf{A}^{(i, j)} = \mathbf{0} \quad \text{for } i < 0 \text{ or } j < 0 \\ \mathbf{A}^{(i, j)} &= \mathbf{A}_1 \mathbf{A}^{(i-1, j)} + \mathbf{A}_2 \mathbf{A}^{(i, j-1)} \\ &= \mathbf{A}^{(i-1, j)} \mathbf{A}_1 + \mathbf{A}^{(i, j-1)} \mathbf{A}_2 \quad \text{for } i, j > 0. \end{aligned} \quad (7)$$

Substituting (6) into (4) yields

$$\begin{aligned} \mathbf{G}(z_1, z_2) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} [\mathbf{c} \mathbf{A}^{(i-1, j)} (\mathbf{A}_1 - \mathbf{D}_1) \\ &\quad + \mathbf{c} \mathbf{A}^{(i, j-1)} (\mathbf{A}_2 - \mathbf{D}_2)] z_1^{-i} z_2^{-j} \\ &\quad + \mathbf{c} - \mathbf{h}. \end{aligned} \quad (8)$$

Substituting (8) into (5), it follows that

$$I(\mathbf{D}_1, \mathbf{D}_2, \mathbf{h}) = J_1(\mathbf{D}_1, \mathbf{D}_2) + \text{tr}[(\mathbf{c} - \mathbf{h})^T (\mathbf{c} - \mathbf{h})] \quad (9)$$

where

$$\begin{aligned} J_1(\mathbf{D}_1, \mathbf{D}_2) &= \text{tr} \left[[(\mathbf{A}_1 - \mathbf{D}_1)^T, (\mathbf{A}_2 - \mathbf{D}_2)^T] \mathbf{W}' \begin{bmatrix} \mathbf{A}_1 - \mathbf{D}_1 \\ \mathbf{A}_2 - \mathbf{D}_2 \end{bmatrix} \right] \\ &= \text{tr} \left[(\mathbf{A}_1 - \mathbf{D}_1)^T \mathbf{W}_o (\mathbf{A}_1 - \mathbf{D}_1) \right. \\ &\quad + (\mathbf{A}_2 - \mathbf{D}_2)^T \mathbf{W}^T (\mathbf{A}_1 - \mathbf{D}_1) \\ &\quad + (\mathbf{A}_1 - \mathbf{D}_1)^T \mathbf{W} (\mathbf{A}_2 - \mathbf{D}_2) \\ &\quad \left. + (\mathbf{A}_2 - \mathbf{D}_2)^T \mathbf{W}_o (\mathbf{A}_2 - \mathbf{D}_2) \right]. \end{aligned}$$

Here, the $2n \times 2n$ matrix \mathbf{W}' is defined by

$$\begin{aligned} \mathbf{W}' &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} (\mathbf{c} \mathbf{A}^{(i-1, j)})^T \\ (\mathbf{c} \mathbf{A}^{(i, j-1)})^T \end{bmatrix} \begin{bmatrix} \mathbf{c} \mathbf{A}^{(i-1, j)} & \mathbf{c} \mathbf{A}^{(i, j-1)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{W}_o & \mathbf{W} \\ \mathbf{W}^T & \mathbf{W}_o \end{bmatrix} \end{aligned} \quad (10)$$

where matrix \mathbf{W}_o is the local observability Gramian of the LSS model in (1). In the case when there is no EF but error feedforward exists, it follows from (9) that

$$\begin{aligned} I(\mathbf{0}, \mathbf{0}, \mathbf{c}) &= \text{tr} \left[[\mathbf{A}_1^T, \mathbf{A}_2^T] \mathbf{W}' \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix} \right] \\ &= \text{tr} \left[\mathbf{A}_1^T \mathbf{W}_o \mathbf{A}_1 + \mathbf{A}_2^T \mathbf{W}^T \mathbf{A}_1 \right. \\ &\quad \left. + \mathbf{A}_1^T \mathbf{W} \mathbf{A}_2 + \mathbf{A}_2^T \mathbf{W}_o \mathbf{A}_2 \right]. \end{aligned} \quad (11)$$

The local controllability Gramian \mathbf{K}_c is defined by

$$\mathbf{K}_c = \sum_{k=1}^{\infty} \sum_{i=0}^k \mathbf{f}(i, k-i) \mathbf{f}^T(i, k-i) \quad (12)$$

where

$$\mathbf{f}(i, j) = \mathbf{A}^{(i-1, j)} \mathbf{b}_1 + \mathbf{A}^{(i, j-1)} \mathbf{b}_2.$$

The LSS model in (1) is said to satisfy L_2 -norm dynamic-range scaling constraints provided that

$$(\mathbf{K}_c)_{ii} = 1 \quad \text{for } i = 1, 2, \dots, n \quad (13)$$

where $(\mathbf{K}_c)_{ii}$ denotes the i th entry of matrix \mathbf{K}_c .

III. SEPARATE OPTIMIZATION OF REALIZATION AND ERROR FEEDBACK

Applying a coordinate transformation defined by

$$\bar{\mathbf{x}}(i, j) = \mathbf{T}^{-1} \mathbf{x}(i, j) \quad (14)$$

to the LSS model $(\mathbf{A}_1, \mathbf{A}_2, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}, d)_n$ in (1), we obtain a new realization $(\bar{\mathbf{A}}_1, \bar{\mathbf{A}}_2, \bar{\mathbf{b}}_1, \bar{\mathbf{b}}_2, \bar{\mathbf{c}}, d)_n$ characterized by

$$\begin{aligned} \bar{\mathbf{A}}_1 &= \mathbf{T}^{-1} \mathbf{A}_1 \mathbf{T}, & \bar{\mathbf{A}}_2 &= \mathbf{T}^{-1} \mathbf{A}_2 \mathbf{T} \\ \bar{\mathbf{b}}_1 &= \mathbf{T}^{-1} \mathbf{b}_1, & \bar{\mathbf{b}}_2 &= \mathbf{T}^{-1} \mathbf{b}_2, & \bar{\mathbf{c}} &= \mathbf{c} \mathbf{T} \end{aligned} \quad (15)$$

where \mathbf{T} is an $n \times n$ nonsingular matrix. The Gramians $\bar{\mathbf{K}}_c$, $\bar{\mathbf{W}}_o$ and $\bar{\mathbf{W}}$ in the new realization can then be written as

$$\bar{\mathbf{K}}_c = \mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T}, \quad \bar{\mathbf{W}}_o = \mathbf{T}^T \mathbf{W}_o \mathbf{T}, \quad \bar{\mathbf{W}} = \mathbf{T}^T \mathbf{W} \mathbf{T} \quad (16)$$

respectively. If the L_2 -norm dynamic-range scaling constraints are imposed on the new realization, then we have

$$(\bar{\mathbf{K}}_c)_{ii} = (\mathbf{T}^{-1} \mathbf{K}_c \mathbf{T}^{-T})_{ii} = 1 \quad \text{for } i = 1, 2, \dots, n. \quad (17)$$

For the new realization with no EF and error feedforward, we consider the problem of minimizing a measure

$$M(\mathbf{P}, \lambda) = \text{tr}[\mathbf{V} \mathbf{P}] + \lambda (\text{tr}[\mathbf{K}_c \mathbf{P}^{-1}] - n) \quad (18)$$

with respect to $n \times n$ nonsingular matrix \mathbf{P} and scalar λ where $\mathbf{P} = \mathbf{T} \mathbf{T}^T$, λ is a Lagrange multiplier, and

$$\mathbf{V} = \mathbf{A}_1^T \mathbf{W}_o \mathbf{A}_1 + \mathbf{A}_2^T \mathbf{W}^T \mathbf{A}_1 + \mathbf{A}_1^T \mathbf{W} \mathbf{A}_2 + \mathbf{A}_2^T \mathbf{W}_o \mathbf{A}_2.$$

We compute

$$\begin{aligned} \frac{\partial J(\mathbf{P}, \lambda)}{\partial \mathbf{P}} &= \mathbf{V} - \lambda \mathbf{P}^{-1} \mathbf{K}_c \mathbf{P}^{-1} \\ \frac{\partial J(\mathbf{P}, \lambda)}{\partial \lambda} &= \text{tr}[\mathbf{K}_c \mathbf{P}^{-1}] - n. \end{aligned} \quad (19)$$

If we let $\partial J(\mathbf{P}, \lambda)/\partial \mathbf{P} = \mathbf{0}$ and $\partial J(\mathbf{P}, \lambda)/\partial \lambda = 0$, then

$$\mathbf{P} \mathbf{V} \mathbf{P} = \lambda \mathbf{K}_c, \quad \text{tr}[\mathbf{K}_c \mathbf{P}^{-1}] = n. \quad (20)$$

It follows from (20) that

$$\begin{aligned} \mathbf{P} &= \sqrt{\lambda} \mathbf{V}^{-\frac{1}{2}} [\mathbf{V}^{\frac{1}{2}} \mathbf{K}_c \mathbf{V}^{\frac{1}{2}}]^{\frac{1}{2}} \mathbf{V}^{-\frac{1}{2}} \\ \frac{1}{\sqrt{\lambda}} \text{tr}[\mathbf{K}_c \mathbf{V}]^{\frac{1}{2}} &= \frac{1}{\sqrt{\lambda}} \left(\sum_{i=1}^n \theta_i \right) = n \end{aligned} \quad (21)$$

where θ_i^2 for $i = 1, 2, \dots, n$ are the eigenvalues of $\mathbf{K}_c \mathbf{V}$. Therefore, we obtain

$$\mathbf{P} = \frac{1}{n} \left(\sum_{i=1}^n \theta_i \right) \mathbf{V}^{-\frac{1}{2}} [\mathbf{V}^{\frac{1}{2}} \mathbf{K}_c \mathbf{V}^{\frac{1}{2}}]^{\frac{1}{2}} \mathbf{V}^{-\frac{1}{2}}. \quad (22)$$

Substituting (22) into (18) yields the minimum value of $M(\mathbf{P}, \lambda)$ as

$$\min_{\mathbf{P}, \lambda} M(\mathbf{P}, \lambda) = \frac{1}{n} \left(\sum_{i=1}^n \theta_i \right)^2. \quad (23)$$

Note that matrix \mathbf{T} assumes the form

$$\mathbf{T} = \mathbf{P}^{\frac{1}{2}} \mathbf{U} \quad (24)$$

where $\mathbf{P}^{1/2}$ is the square root of \mathbf{P} obtained above, and \mathbf{U} is the $n \times n$ orthogonal matrix such that (17) is satisfied.

If the coordinate transformation in (14) is applied to the LSS model in (1), then (9) is changed to

$$\bar{I}(\mathbf{D}_1, \mathbf{D}_2, \mathbf{h}, \mathbf{T}) = J_2(\mathbf{D}_1, \mathbf{D}_2, \mathbf{T}) + \text{tr}[(\bar{\mathbf{c}} - \mathbf{h})^T (\bar{\mathbf{c}} - \mathbf{h})] \quad (25)$$

where

$$\begin{aligned} J_2(\mathbf{D}_1, \mathbf{D}_2, \mathbf{T}) &= \text{tr}[(\bar{\mathbf{A}}_1 - \mathbf{D}_1)^T \bar{\mathbf{W}}_o (\bar{\mathbf{A}}_1 - \mathbf{D}_1) \\ &\quad + (\bar{\mathbf{A}}_2 - \mathbf{D}_2)^T \bar{\mathbf{W}}^T (\bar{\mathbf{A}}_1 - \mathbf{D}_1) \\ &\quad + (\bar{\mathbf{A}}_1 - \mathbf{D}_1)^T \bar{\mathbf{W}} (\bar{\mathbf{A}}_2 - \mathbf{D}_2) \\ &\quad + (\bar{\mathbf{A}}_2 - \mathbf{D}_2)^T \bar{\mathbf{W}}_o (\bar{\mathbf{A}}_2 - \mathbf{D}_2)]. \end{aligned}$$

Case 1: \mathbf{D}_1 and \mathbf{D}_2 are general matrices

In this case, we select the matrices \mathbf{D}_1 and \mathbf{D}_2 as

$$\mathbf{D}_1 = \bar{\mathbf{A}}_1, \quad \mathbf{D}_2 = \bar{\mathbf{A}}_2. \quad (26)$$

Case 2: \mathbf{D}_1 and \mathbf{D}_2 are diagonal matrices

We define

$$\begin{aligned} \mathbf{D}_1 &= \text{diag}\{\alpha_1, \alpha_2, \dots, \alpha_n\} \\ \mathbf{D}_2 &= \text{diag}\{\beta_1, \beta_2, \dots, \beta_n\}. \end{aligned} \quad (27)$$

From $\partial J_2(\mathbf{D}_1, \mathbf{D}_2, \mathbf{T})/\partial \alpha_i = 0$ and $\partial J_2(\mathbf{D}_1, \mathbf{D}_2, \mathbf{T})/\partial \beta_i = 0$, it is derived that for $i = 1, 2, \dots, n$

$$\begin{aligned} \alpha_i &= \frac{\bar{\mathbf{W}}_o(i, i) \mathbf{M}_1(i, i) - \bar{\mathbf{W}}(i, i) \mathbf{M}_2(i, i)}{\bar{\mathbf{W}}_o(i, i)^2 - \bar{\mathbf{W}}(i, i)^2} \\ \beta_i &= \frac{\bar{\mathbf{W}}_o(i, i) \mathbf{M}_2(i, i) - \bar{\mathbf{W}}(i, i) \mathbf{M}_1(i, i)}{\bar{\mathbf{W}}_o(i, i)^2 - \bar{\mathbf{W}}(i, i)^2} \end{aligned} \quad (28)$$

where $X(i, j)$ denotes the ij th element of matrix \mathbf{X} and

$$\begin{aligned} \mathbf{M}_1 &= \bar{\mathbf{W}}_o \bar{\mathbf{A}}_1 + \bar{\mathbf{W}} \bar{\mathbf{A}}_2 \\ \mathbf{M}_2 &= \bar{\mathbf{W}}_o \bar{\mathbf{A}}_2 + \bar{\mathbf{W}}^T \bar{\mathbf{A}}_1. \end{aligned}$$

Case 3: \mathbf{D}_1 and \mathbf{D}_2 are scalar matrices

With scalars α and β , we define

$$\mathbf{D}_1 = \alpha \mathbf{I}_n, \quad \mathbf{D}_2 = \beta \mathbf{I}_n. \quad (29)$$

From $\partial J_2(\mathbf{D}_1, \mathbf{D}_2, \mathbf{T})/\partial \alpha = 0$ and $\partial J_2(\mathbf{D}_1, \mathbf{D}_2, \mathbf{T})/\partial \beta = 0$, we obtain

$$\begin{aligned} \alpha &= \frac{\text{tr}[\bar{\mathbf{W}}_o] \text{tr}[\mathbf{M}_1] - \text{tr}[\bar{\mathbf{W}}] \text{tr}[\mathbf{M}_2]}{(\text{tr}[\bar{\mathbf{W}}_o])^2 - (\text{tr}[\bar{\mathbf{W}}])^2} \\ \beta &= \frac{\text{tr}[\bar{\mathbf{W}}_o] \text{tr}[\mathbf{M}_2] - \text{tr}[\bar{\mathbf{W}}] \text{tr}[\mathbf{M}_1]}{(\text{tr}[\bar{\mathbf{W}}_o])^2 - (\text{tr}[\bar{\mathbf{W}}])^2}. \end{aligned} \quad (30)$$

IV. JOINT OPTIMIZATION OF ERROR FEEDBACK AND REALIZATION

Define

$$\hat{\mathbf{T}} = \mathbf{T}^T \mathbf{K}_c^{-\frac{1}{2}}. \quad (31)$$

Then (17) can be written as

$$(\hat{\mathbf{T}}^{-T} \hat{\mathbf{T}}^{-1})_{ii} = 1 \quad \text{for } i = 1, 2, \dots, n. \quad (32)$$

The constraints in (32) simply state that each column in matrix $\hat{\mathbf{T}}^{-1}$ must be a unity vector. These are satisfied if $\hat{\mathbf{T}}^{-1}$ assumes the form

$$\hat{\mathbf{T}}^{-1} = \left[\frac{\mathbf{t}_1}{\|\mathbf{t}_1\|}, \frac{\mathbf{t}_2}{\|\mathbf{t}_2\|}, \dots, \frac{\mathbf{t}_n}{\|\mathbf{t}_n\|} \right] \quad (33)$$

where \mathbf{t}_i 's for $i = 1, 2, \dots, n$ are $n \times 1$ real vectors. In such a case, (25) can be expressed as

$$\hat{I}(\mathbf{D}_1, \mathbf{D}_2, \mathbf{h}, \hat{\mathbf{T}}) = J_3(\mathbf{D}_1, \mathbf{D}_2, \hat{\mathbf{T}}) + \text{tr}[(\hat{\mathbf{c}} - \mathbf{h})^T (\hat{\mathbf{c}} - \mathbf{h})] \quad (34)$$

where

$$\begin{aligned} J_3(\mathbf{D}_1, \mathbf{D}_2, \hat{\mathbf{T}}) &= \text{tr}[(\hat{\mathbf{A}}_1 - \mathbf{D}_1)^T \hat{\mathbf{W}}_o (\hat{\mathbf{A}}_1 - \mathbf{D}_1) \\ &\quad + (\hat{\mathbf{A}}_2 - \mathbf{D}_2)^T \hat{\mathbf{W}}^T (\hat{\mathbf{A}}_1 - \mathbf{D}_1) \\ &\quad + (\hat{\mathbf{A}}_1 - \mathbf{D}_1)^T \hat{\mathbf{W}} (\hat{\mathbf{A}}_2 - \mathbf{D}_2) \\ &\quad + (\hat{\mathbf{A}}_2 - \mathbf{D}_2)^T \hat{\mathbf{W}}_o (\hat{\mathbf{A}}_2 - \mathbf{D}_2)] \end{aligned}$$

with

$$\begin{aligned} \hat{\mathbf{A}}_1 &= \hat{\mathbf{T}}^{-T} (\mathbf{K}_c^{-\frac{1}{2}} \mathbf{A}_1 \mathbf{K}_c^{\frac{1}{2}}) \hat{\mathbf{T}}^T \\ \hat{\mathbf{A}}_2 &= \hat{\mathbf{T}}^{-T} (\mathbf{K}_c^{-\frac{1}{2}} \mathbf{A}_2 \mathbf{K}_c^{\frac{1}{2}}) \hat{\mathbf{T}}^T, \quad \hat{\mathbf{c}} = (\mathbf{c} \mathbf{K}_c^{\frac{1}{2}}) \hat{\mathbf{T}}^T \\ \hat{\mathbf{W}}_o &= \hat{\mathbf{T}} (\mathbf{K}_c^{\frac{1}{2}} \mathbf{W}_o \mathbf{K}_c^{\frac{1}{2}}) \hat{\mathbf{T}}^T, \quad \hat{\mathbf{W}} = \hat{\mathbf{T}} (\mathbf{K}_c^{-\frac{1}{2}} \mathbf{W} \mathbf{K}_c^{\frac{1}{2}}) \hat{\mathbf{T}}^T. \end{aligned}$$

When selecting vector \mathbf{h} as $\mathbf{h} = \hat{\mathbf{c}}$, the problem of obtaining matrices \mathbf{D}_1 , \mathbf{D}_2 and \mathbf{T} that minimize $J_2(\mathbf{D}_1, \mathbf{D}_2, \mathbf{T})$ in (25) subject to the scaling constraints in (17) can be converted into an unconstrained optimization problem of obtaining matrices \mathbf{D}_1 , \mathbf{D}_2 and $\hat{\mathbf{T}}$ that minimize $J_3(\mathbf{D}_1, \mathbf{D}_2, \hat{\mathbf{T}})$ in (34).

Let \mathbf{x} be the column vector that collects the variables in matrices \mathbf{D}_1 , \mathbf{D}_2 and $\hat{\mathbf{T}}$. Then, $J_3(\mathbf{D}_1, \mathbf{D}_2, \hat{\mathbf{T}})$ is a function

of \mathbf{x} , denoted by $J_3(\mathbf{x})$. The algorithm starts with a trivial initial point \mathbf{x}_0 obtained from an initial assignment $\mathbf{D}_1 = \mathbf{D}_2 = \hat{\mathbf{T}} = \mathbf{I}_n$. In the k th iteration, a quasi-Newton algorithm updates the most recent point \mathbf{x}_k to point \mathbf{x}_{k+1} as [9]

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \quad (35)$$

where

$$\begin{aligned} \mathbf{d}_k &= -\mathbf{S}_k \nabla J_3(\mathbf{x}_k), \quad \alpha_k = \arg \min_{\alpha} J_3(\mathbf{x}_k + \alpha \mathbf{d}_k) \\ \mathbf{S}_{k+1} &= \mathbf{S}_k + \left(1 + \frac{\gamma_k^T \mathbf{S}_k \gamma_k}{\gamma_k^T \delta_k} \right) \frac{\delta_k \delta_k^T}{\gamma_k^T \delta_k} - \frac{\delta_k \gamma_k^T \mathbf{S}_k + \mathbf{S}_k \gamma_k \delta_k^T}{\gamma_k^T \delta_k} \\ \mathbf{S}_0 &= \mathbf{I}, \quad \delta_k = \mathbf{x}_{k+1} - \mathbf{x}_k, \quad \gamma_k = \nabla J_3(\mathbf{x}_{k+1}) - \nabla J_3(\mathbf{x}_k). \end{aligned}$$

Here, $\nabla J_3(\mathbf{x})$ is the gradients of $J_3(\mathbf{x})$ with respect to \mathbf{x} , and \mathbf{S}_k is a positive-definite approximation of the inverse Hessian matrix of $J_3(\mathbf{x})$. This iteration process continues until

$$|J_3(\mathbf{x}_{k+1}) - J_3(\mathbf{x}_k)| < \varepsilon \quad (36)$$

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

When \mathbf{D}_1 and \mathbf{D}_2 are general matrices, vector \mathbf{x} consists of matrix $\hat{\mathbf{T}}$ only. After obtaining matrix $\hat{\mathbf{T}}$, we select matrices \mathbf{D}_1 and \mathbf{D}_2 as

$$\mathbf{D}_1 = \hat{\mathbf{A}}_1, \quad \mathbf{D}_2 = \hat{\mathbf{A}}_2. \quad (37)$$

The gradient of $J(\mathbf{x})$ with respect to the ij th element of $\hat{\mathbf{T}}$ is found to be

$$\frac{\partial J(\mathbf{x})}{\partial t_{ij}} = \lim_{\Delta \rightarrow 0} \frac{J(\hat{\mathbf{T}}_{ij}) - J(\hat{\mathbf{T}})}{\Delta} \quad (38)$$

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]

$$\hat{\mathbf{T}}_{ij} = \hat{\mathbf{T}} + \frac{\Delta \hat{\mathbf{T}} \mathbf{g}_{ij} \mathbf{e}_j^T \hat{\mathbf{T}}}{1 - \Delta \mathbf{e}_j^T \hat{\mathbf{T}} \mathbf{g}_{ij}}$$

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

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

V. A NUMERICAL EXAMPLE

Consider a stable, locally controllable, and locally observable 2-D state-space digital filter with order $n = 4$ specified by

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 0 & 0 & -0.00411 \\ 1 & 0 & 0 & 0.08007 \\ 0 & 1 & 0 & -0.42458 \\ 0 & 0 & 1 & 1.04460 \end{bmatrix}, \quad \mathbf{b}_1 = \begin{bmatrix} -0.01452 \\ 0.01234 \\ 0.02054 \\ 0.04762 \end{bmatrix}$$

$$\mathbf{A}_2 = \begin{bmatrix} -0.22608 & -0.40594 & -0.30955 & -0.14469 \\ 1.61428 & 1.61040 & 1.02336 & 0.43872 \\ 0.10054 & -0.60615 & -0.45322 & -0.31019 \\ -0.00723 & 0.24580 & 0.38668 & 0.56289 \end{bmatrix}$$

$$\mathbf{b}_2 = [0.01189 \quad 0.02351 \quad -0.00637 \quad 0.02094]^T$$

$$\mathbf{c} = [0 \quad 0 \quad 0 \quad 1], \quad d = 0.00943.$$

After carrying out the L_2 -scaling for the above LSS model with a diagonal coordinate transformation matrix, the noise gain of the scaled LSS model with error feedforward was computed as $I(\mathbf{0}, \mathbf{0}, \mathbf{c}) = 76.641884$. Next, the matrix \mathbf{P} was derived from (22) and substituting it into (18) produced $M(\mathbf{P}, \lambda) = 3.230958$. The other results obtained by applying the proposed technique are summarized in Tables I and II.

VI. CONCLUSION

The separate/joint optimization of EF and realization has been investigated to minimize RN subject to L_2 -scaling constraints for a class of 2-D state-space digital filters. It has been shown that the problem in the joint optimization can be converted into an unconstrained optimization problem by using linear algebraic techniques. An efficient quasi-Newton algorithm has then been employed to solve the unconstrained optimization problem iteratively. Our computer simulation results have demonstrated the effectiveness of the proposed techniques.

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 quadrantly 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, pp. 2436-2445, Sept. 2003.
- [9] R. Fletcher, *Practical Methods of Optimization*, 2nd ed. Wiley, New York, 1987.
- [10] T. Kailath, *Linear Systems*, Prentice Hall, 1980.

TABLE I
ROUND OFF NOISE GAIN IN SEPARATE OPTIMIZATION

Error Feedback Matrices	General	Diagonal	Scalar
Infinite Precision	0	0.454620	0.469953
3-Bit Quantization	0.049936	0.461004	0.481157
Integer Quantization	3.100355	2.217879	1.586120

TABLE II
ROUND OFF NOISE GAIN IN JOINT OPTIMIZATION

Error Feedback Matrices	General	Diagonal	Scalar
Infinite Precision	0	0.186190	0.233562
3-Bit Quantization	0.049936	0.216724	0.243361
Integer Quantization	3.100355	1.736732	1.808645