Roundoff Noise Minimization of State-Space Digital Filters Using Separate and Joint Error Feedback/Coordinate Transformation Optimization

Takao Hinamoto, Fellow, IEEE, Hiroaki Ohnishi, and Wu-Sheng Lu, Fellow, IEEE

Abstract—This paper investigates the problem of minimizing roundoff noise under l_2 -norm dynamic-range scaling constraints in state-space digital filters by means of error feedback as well as joint error feedback/coordinate transformation optimization. First, several techniques for the determination of optimal full-scale, diagonal, and scalar error-feedback matrices for a given state-space digital filter are proposed, where three realization schemes, namely, the general state-space realization, input-balanced realization, and optimal realization in the sense of Hwang-Mullis-Roberts are examined. Furthermore, an iterative approach is developed for jointly optimizing a scalar error-feedback matrix and a coordinate transformation matrix so as to minimize the roundoff noise subject to the l_2 -norm dynamic-range scaling constraints. The proposed method may be regarded as an alternative, but much simpler and more general, approach to Hwang's method for synthesizing the optimal filter structure with minimum roundoff noise. A case study is included to illustrate the utility of the proposed techniques.

Index Terms—Optimal coordinate transformation, optimal error feedback, roundoff noise minimization, scaling constraints, state-space digital filters.

I. INTRODUCTION

T HE basic arithmetic operations involved in the imple-mentation of an infinite-impulse response (IIR) digital filter are multiplications of input/output samples by the filter coefficients and additions. For arithmetic operations involving fixed-point numbers, the result of a multiplication must be rounded or truncated. This quantization error generates roundoff noise at the filter output. In addition, because the result of an addition can exceed the finite register length, the dynamic range of the digital filter is always a concern in a fixed-point implementation. Error feedback has been known as an effective technique for reducing the roundoff noise at the filter output. This is achieved by extracting the quantization error after multiplication and addition, and then feeding the error signal back through simple circuits. It can be applied to digital filters that are described by either external or internal models for roundoff noise minimization without affecting the filter's input-output characteristics. Many techniques for error feedback have been presented in the past [1]-[10]. Another

T. Hinamoto and H. Ohnishi are with the Graduate School of Engineering, Hiroshima University, Higashi-Hiroshima 739-8527, Japan.

W.-S. Lu is with the Department of Electrial and Computer Engineering, University of Victoria, Victoria, BC V8W 3P6, Canada.

Digital Object Identifier 10.1109/TCSI.2002.807512

technique for reducing the roundoff noise is to synthesize an optimal state-space filter structure that minimizes the roundoff noise under l_2 -norm dynamic-range scaling constraints on the state-variables by means of coordinate transformation in the state space [11]–[14]. In addition, techniques for reducing the roundoff noise subject to the scaling constraints by combining the coordinate transformation with error feedback have been proposed [15], [16]. It has also been shown that the roundoff noise can be reduced by means of delta operator [16], [17] and that the digital filter in this case can be viewed as a special case of the filter with error feedback [16].

In this paper, several new algorithms for reducing the roundoff noise in state-space digital filters are proposed. First, the problem of roundoff noise reduction for several typical state-space realizations using error feedback is investigated. Closed-form formulas for evaluating the optimal full-scale, diagonal, and scalar error-feedback matrices for a given state-space digital filter are derived, where three realization schemes, namely, the general state-space realization, input-balanced realization (to be defined shortly), and optimal realization (in the sense of [12], [13]) are examined. Furthermore, an iterative noise reduction technique for state-space digital filters that jointly optimizes a scalar error-feedback matrix and a coordinate transformation matrix is proposed. A case study is presented to illustrate the algorithms proposed and to demonstrate their performance as compared with that of the existing methods [15], [16].

Throughout this paper, I_n stands for the identity matrix of dimension $n \times n$, the transpose (conjugate transpose) of a matrix A is indicated by A^T (A^*), and the trace, eigenvalue, and *i*th diagonal element of a square matrix A are denoted by tr[A], $\lambda(A)$, and (A)_{*ii*}, respectively.

II. STATE-SPACE DIGITAL FILTERS WITH ERROR FEEDBACK

Let $(A, b, c, d)_n$ be a state-space description of an *n*th-order IIR digital filter, i.e.

$$\boldsymbol{x}(k+1) = \boldsymbol{A}\boldsymbol{x}(k) + \boldsymbol{b}\boldsymbol{u}(k)$$
$$\boldsymbol{y}(k) = \boldsymbol{c}\boldsymbol{x}(k) + d\boldsymbol{u}(k) \tag{1}$$

where $\boldsymbol{x}(k)$ is an $n \times 1$ state-variable vector, u(k) is a scalar input, y(k) is a scalar output, and \boldsymbol{A} , \boldsymbol{b} , \boldsymbol{c} and d are real constant matrices of appropriate dimensions. The filter described in (1) is assumed to be stable, controllable and observable. Due to finite register sizes, finite-word-length (FWL) constraints are

Manuscript received February 20, 2002; revised September 9, 2002. This paper was recommended by Associate Editor W. P. Zhu.



Fig. 1. Error feedback in a state-space digital filter.

imposed on the state variables, input, output, and system parameters as shown in $(A, b, c, d)_n$.

By taking the quantization performed before matrix-vector multiplication into account, an FWL implementation of (1) can be expressed as

$$\tilde{\boldsymbol{x}}(k+1) = \boldsymbol{A}\boldsymbol{Q}\left[\tilde{\boldsymbol{x}}(k)\right] + \boldsymbol{b}\boldsymbol{u}(k)$$
$$\tilde{\boldsymbol{y}}(k) = \boldsymbol{c}\boldsymbol{Q}\left[\tilde{\boldsymbol{x}}(k)\right] + d\boldsymbol{u}(k)$$
(2)

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

The quantizer $Q[\cdot]$ in (2) rounds the *B*-b fraction $\tilde{x}(k)$ to $(B-B_c)$ b after completing the multiplications and additions, where the sign bit is not counted. In a fixed-point implementation, the quantization is usually carried out by two's complement truncation which discards the lower bits of a double-precision accumulator. Therefore, the quantization error

$$\boldsymbol{e}(k) = \tilde{\boldsymbol{x}}(k) - \boldsymbol{Q}\left[\tilde{\boldsymbol{x}}(k)\right]$$
(3)

is equal to the residue left in the lower part of $\tilde{\boldsymbol{x}}(k)$. The blockdiagram representation of a state-space digital filter with error feedback is shown in Fig. 1 where the roundoff error $\boldsymbol{e}(k)$ is modeled as a zero-mean noise process of covariance $\sigma^2 \boldsymbol{I}_n$ with

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

and the quantization error e(k) is fed back through an $n \times n$ constant matrix D in the FWL filter (2). From Fig. 1, it is obvious that the filter can be characterized by the state-space model

$$\tilde{\boldsymbol{x}}(k+1) = \boldsymbol{A}\boldsymbol{Q}\left[\tilde{\boldsymbol{x}}(k)\right] + \boldsymbol{b}\boldsymbol{u}(k) + \boldsymbol{D}\boldsymbol{e}(k)$$
$$\tilde{\boldsymbol{y}}(k) = \boldsymbol{c}\boldsymbol{Q}\left[\tilde{\boldsymbol{x}}(k)\right] + d\boldsymbol{u}(k)$$
(4)

where D is referred to as an *error-feedback matrix*. Subtracting (4) from (1) yields

$$\Delta \boldsymbol{x} (k+1) = \boldsymbol{A} \Delta \boldsymbol{x}(k) + (\boldsymbol{A} - \boldsymbol{D}) \boldsymbol{e}(k)$$

$$\Delta \boldsymbol{y}(k) = \boldsymbol{c} \Delta \boldsymbol{x}(k) + \boldsymbol{c} \boldsymbol{e}(k)$$
(5)

where

$$\Delta \boldsymbol{x}(k) = \boldsymbol{x}(k) - \tilde{\boldsymbol{x}}(k), \quad \Delta y(k) = y(k) - \tilde{y}(k).$$

Taking the z-transform on both sides of (5) and setting $\Delta x(0) = 0$, we have

$$\Delta Y(z) = \boldsymbol{G}_D(z)\boldsymbol{E}(z)$$

$$\boldsymbol{G}_D(z) = \boldsymbol{c}(z\boldsymbol{I}_n - \boldsymbol{A})^{-1}(\boldsymbol{A} - \boldsymbol{D}) + \boldsymbol{c}$$
(6)

where $\Delta Y(z)$ and E(z) represent the z-transform of $\Delta y(k)$ and e(k), respectively, and $G_D(z)$ is the transfer function from the quantization error, e(k), to the filter output $\Delta y(k)$.

The noise gain $I(\mathbf{D}) = \sigma_{\text{out}}^2 / \sigma^2$ is then defined by

$$I(D) = \frac{1}{2\pi j} \int_{|z|=1} G_D(z) G_D^*(z) \frac{dz}{z}$$
$$= \operatorname{tr}[W_D]$$
(7)

where

V

$$\boldsymbol{W}_D = \frac{1}{2\pi j} \int_{|z|=1} \boldsymbol{G}_D^*(z) \boldsymbol{G}_D(z) \frac{dz}{z}.$$

Utilizing the Cauchy integral theorem, the matrix W_D defined in (7) can be expressed in closed-form as

$$\boldsymbol{W}_{D} = (\boldsymbol{A} - \boldsymbol{D})^{T} \boldsymbol{W}_{o} (\boldsymbol{A} - \boldsymbol{D}) + \boldsymbol{c}^{T} \boldsymbol{c}$$
(8)

where W_o is the observability Gramian of the filter that can be obtained by solving the Lyapunov equation

$$\boldsymbol{W}_o = \boldsymbol{A}^T \boldsymbol{W}_o \boldsymbol{A} + \boldsymbol{c}^T \boldsymbol{c}. \tag{9}$$

Matrix W_o is referred to as the *unit noise matrix* for the filter, and W_D can be viewed as the unit noise matrix for the filter with error feedback specified by matrix D.

The l_2 -norm dynamic-range scaling constraints on the state variables involve the controllability Gramian of the filter, which is defined by

$$\boldsymbol{K}_{c} = \sum_{k=0}^{\infty} \boldsymbol{A}^{k} \boldsymbol{b} \boldsymbol{b}^{T} (\boldsymbol{A}^{T})^{k}$$
(10)

and can be computed by solving the Lyapunov equation

$$\boldsymbol{K}_c = \boldsymbol{A}\boldsymbol{K}_c \boldsymbol{A}^T + \boldsymbol{b}\boldsymbol{b}^T.$$
(11)

The problem considered here is to design the error-feedback matrix D so as to reduce the noise gain in the sense that

$$\operatorname{tr}[\boldsymbol{W}_D] < \operatorname{tr}[\boldsymbol{W}_o] \tag{12}$$

is satisfied subject to that all the diagonal elements of K_c equal unity. Such constraints on matrix K_c are known as the l_2 -norm dynamic-range scaling constraints [13].

A different yet equivalent state-space description of (1), $(\overline{A}, \overline{b}, \overline{c}, d)_n$, can be obtained via a coordinate transformation $\overline{x}(k) = T^{-1}x(k)$ where

$$\overline{A} = T^{-1}AT \quad \overline{b} = T^{-1}b \quad \overline{c} = cT.$$
(13)

Accordingly, the observability and controllability Gramians for $(\overline{A}, \overline{b}, \overline{c}, d)_n$ become

$$\overline{W}_o = T^T W_o T, \quad \overline{K}_c = T^{-1} K_c T^{-T}$$
(14)

respectively. If the l_2 -norm dynamic-range scaling constraints are imposed on the state-variable vector $\overline{\boldsymbol{x}}(k)$, i.e.

$$(\overline{K}_c)_{ii} = (T^{-1}K_cT^{-T})_{ii} = 1, \quad i = 1, 2, \cdots, n$$
 (15)

then, it can be shown that [12], [13]

$$\min_{\boldsymbol{T}} \operatorname{tr}[\boldsymbol{\overline{W}}_o] = \frac{1}{n} \left(\sum_{i=1}^n \sigma_i\right)^2 \tag{16}$$

where σ_i^2 for $i = 1, 2, \dots, n$ are the eigenvalues of matrix $K_c W_o$. The state-space realization satisfying (15) and (16) is called the *optimal realization* (which is sometimes also referred to as the *optimal filter structure*) with minimum roundoff noise. A method for constructing such a filter structure was proposed in [12], [13]. Obviously, if the realization in (1) is optimal with minimum roundoff noise, then the right-hand side of (12) is minimized.

In the next section, we will derive closed-form formulas for determining the optimal full-scale, diagonal, and scalar error-feedback matrix D that minimizes tr $[W_D]$ for a given state-space digital filter.

III. DETERMINATION OF OPTIMAL ERROR-FEEDBACK MATRICES

A. Case 1: D is a General Matrix

Substituting (8) into (7), we obtain

$$I(\boldsymbol{D}) = \operatorname{tr}[\boldsymbol{c}^{T}\boldsymbol{c} + (\boldsymbol{A} - \boldsymbol{D})^{T}\boldsymbol{W}_{o}(\boldsymbol{A} - \boldsymbol{D})]$$

= tr[\mathbf{W}_{o}] + tr[\mathbf{D}^{T}\boldsymbol{W}_{o}\boldsymbol{D}] - 2\operatorname{tr}[\mathbf{D}^{T}\boldsymbol{W}_{o}\boldsymbol{A}]. \quad (17)

Differentiating (17) with respect to the error-feedback matrix D yields

$$\frac{\partial I(\boldsymbol{D})}{\partial \boldsymbol{D}} = 2\boldsymbol{W}_o(\boldsymbol{D} - \boldsymbol{A}). \tag{18}$$

By choosing the error-feedback matrix as D = A, the noise gain I(D) in (17) achieves its minimum value

$$I_{\min}(\boldsymbol{D}) = \operatorname{tr}[\boldsymbol{W}_o] - \operatorname{tr}[\boldsymbol{A}^T \boldsymbol{W}_o \boldsymbol{A}] = \operatorname{tr}[\boldsymbol{c}^T \boldsymbol{c}].$$
(19)

B. Case 2: D is a Diagonal Matrix

In this case, matrix D assumes the form

$$\boldsymbol{D} = \operatorname{diag}\{\alpha_1, \alpha_2, \cdots, \alpha_n\}$$
(20)

which leads (17) to

$$I(\boldsymbol{D}) = \operatorname{tr}[\boldsymbol{W}_o] + \sum_{i=1}^{n} (\boldsymbol{W}_o)_{ii} \alpha_i^2 - 2 \sum_{i=1}^{n} (\boldsymbol{W}_o \boldsymbol{A})_{ii} \alpha_i. \quad (21)$$

From (21), it is clear that $I(D) = tr[W_D] < tr[W_o]$ holds if α_i 's satisfy

$$\alpha_i \left(\alpha_i - 2 \frac{(\boldsymbol{W}_o \boldsymbol{A})_{ii}}{(\boldsymbol{W}_o)_{ii}} \right) < 0, \qquad i = 1, 2, \cdots, n.$$
 (22)

If we compute the gradient of $I(\mathbf{D})$ in (21) and set it to zero, i.e.,

$$\frac{\partial I(\boldsymbol{D})}{\partial \alpha_i} = 2 \left[(\boldsymbol{W}_o)_{ii} \alpha_i - (\boldsymbol{W}_o \boldsymbol{A})_{ii} \right] = 0, \qquad i = 1, 2, \cdots, n$$
(23)

then, we obtain the minimizer $D = \text{diag}\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ with

$$\alpha_i = \frac{(\boldsymbol{W}_o \boldsymbol{A})_{ii}}{(\boldsymbol{W}_o)_{ii}}, \qquad i = 1, 2, \cdots, n$$
(24)

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

$$I_{\min}(D) = tr[W_o] - \sum_{i=1}^n \frac{(W_o A)_{ii}^2}{(W_o)_{ii}}.$$
 (25)

In the rest of the paper, the filter in (1) is said to be *input-balanced (internally balanced)* if $K_c = I_n$ and $W_o = \Sigma^2(K_c = W_o = \Sigma)$ where $\Sigma = \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_n\}$. The theorem below characterizes the diagonal error-feedback matrix that minimizes the noise gain I(D) for the input-balanced realization of an IIR digital filter.

Theorem 1: If the filter in (1) is input-balanced, and the diagonal error-feedback matrix $D = \text{diag}\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ satisfies

$$\alpha_i(\alpha_i - 2a_{ii}) < 0, \qquad i = 1, 2, \cdots, n \tag{26}$$

then $I(\mathbf{D}) = tr[\mathbf{W}_D] < tr[\mathbf{W}_o]$ where a_{ij} denotes the (i, j)th element of matrix \mathbf{A} . Moreover, if

$$\alpha_i = a_{ii}, \qquad i = 1, 2, \cdots, n \tag{27}$$

then I(D) achieves the minimum value

$$I_{\min}(\mathbf{D}) = \sum_{i=1}^{n} (1 - a_{ii}^2) \sigma_i^2.$$
 (28)

Proof: $W_o = \Sigma^2$ implies that $w_{ii}^o = \sigma_i^2$ and $w_{ij}^o = 0$ for $i \neq j$. Under these circumstances, (22), (24) and (25) are reduced to (26), (27) and (28), respectively. This completes the proof.

C. Case 3: **D** is a Scalar Matrix αI_n

If $D = \alpha I_n$ with a scalar α , then (17) becomes

$$I(\boldsymbol{D}) = \operatorname{tr}[\boldsymbol{W}_o] + \alpha \left(\operatorname{tr}[\boldsymbol{W}_o]\alpha - 2\operatorname{tr}[\boldsymbol{W}_o\boldsymbol{A}] \right).$$
(29)

Hence, $I(\mathbf{D}) = tr[\mathbf{W}_D] < tr[\mathbf{W}_o]$ holds if α satisfies

$$\alpha \left(\alpha - 2 \frac{\operatorname{tr}[\boldsymbol{W}_{o}\boldsymbol{A}]}{\operatorname{tr}[\boldsymbol{W}_{o}]} \right) < 0.$$
(30)

Moreover, from $\partial I(\mathbf{D})/\partial \alpha = 0$ it follows that the value of α that minimizes $I(\mathbf{D})$ is given by

$$\alpha = \frac{\operatorname{tr}[\boldsymbol{W}_{o}\boldsymbol{A}]}{\operatorname{tr}[\boldsymbol{W}_{o}]} \tag{31}$$

which leads (29) to

$$I_{\min}(\boldsymbol{D}) = \operatorname{tr}[\boldsymbol{W}_o] \left[1 - \left(\frac{\operatorname{tr}[\boldsymbol{W}_o \boldsymbol{A}]}{\operatorname{tr}[\boldsymbol{W}_o]} \right)^2 \right].$$
(32)

Therem 2 concerns the optimal realization satisfying (15) and (16).

Theorem 2: If the state-space realization (1) has an optimal realization with minimum roundoff noise, then (30) and (31) can be expressed as

$$\alpha \left(\alpha - 2 \frac{a_{11}^{(b)} \sigma_1 + \dots + a_{nn}^{(b)} \sigma_n}{\sigma_1 + \sigma_2 + \dots + \sigma_n} \right) < 0$$
(33)

and

$$\alpha = \frac{a_{11}^{(b)}\sigma_1 + \dots + a_{nn}^{(b)}\sigma_n}{\sigma_1 + \sigma_2 + \dots + \sigma_n}$$
(34)

respectively, where $a_{ii}^{(b)}$ for $i = 1, 2, \dots, n$ are the diagonal elements of the matrix A in the input-balanced realization. Moreover, (32) becomes

$$I_{\min}(\boldsymbol{D}) = \frac{1}{n} \left[\left(\sum_{i=1}^{n} \sigma_i \right)^2 - \left(\sum_{i=1}^{n} a_{ii}^{(b)} \sigma_i \right)^2 \right].$$
(35)

Proof: See Appendix I.

(1)

Corollary 1: If the state-space realization (1) is input-balanced, then (30) and (31) can be expressed as

$$\alpha \left(\alpha - 2 \frac{a_{11}\sigma_1^2 + \dots + a_{nn}\sigma_n^2}{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2} \right) < 0 \tag{36}$$

and

$$\alpha = \frac{a_{11}\sigma_1^2 + \dots + a_{nn}\sigma_n^2}{\sigma_1^2 + \dots + \sigma_n^2} \tag{37}$$

respectively, and (32) becomes

$$I_{\min}(\mathbf{D}) = \left(\sum_{i=1}^{n} \sigma_i^2\right) - \frac{(a_{11}\sigma_1^2 + \dots + a_{nn}\sigma_n^2)^2}{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2}.$$
 (38)

Lemma 1: If the l_2 -norm dynamic-range scaling constraints $(\overline{K}_c)_{ii} = (T^{-1}K_cT^{-T})_{ii} = 1$ for $i = 1, 2, \cdots, n$ are imposed on the state variables, then

$$\min_{\boldsymbol{T}} \operatorname{tr}[\boldsymbol{W}_D] = \frac{1}{n} \left(\sum_{i=1}^n \nu_i \right)^2 \tag{39}$$

where $\overline{W}_D = T^T W_D T$ with $D = \alpha I_n$, and ν_i^2 for i = $1, 2, \cdots, n$ are the eigenvalues of $K_c W_D$.

The proof of Lemma 1 is similar to that in [13] and is therefore omitted.

Remark 1: Under the l_2 -norm dynamic-range scaling constraints, Williamson [15] obtained a suboptimal realization with an integer error-feedback matrix $D = I_n (D = -I_n)$ for low-pass (high-pass) narrow-band filters by choosing an appropriate coordinate transformation matrix T. A similar argument was also made in [16] by restricting D to I_n .

Theorem 3: Let the state-space realization in (1) be inputbalanced. If there exists a real number α such that

$$\alpha(\alpha - 2a_{ii}) < 0 \tag{40}$$

then

$$\sum_{i=1}^{n} \nu_i < \sum_{i=1}^{n} \sigma_i \tag{41}$$

or equivalently

$$\min_{\boldsymbol{T}} \operatorname{tr}[\boldsymbol{T}^T \boldsymbol{W}_D \boldsymbol{T}] < \min_{\boldsymbol{T}} \operatorname{tr}[\boldsymbol{T}^T \boldsymbol{W}_o \boldsymbol{T}].$$
(42)

Moreover, if the diagonal elements of matrix A, a_{ii} for i = $1, 2, \dots, n$, are equal, then the upper bound of $\sum_{i=1}^{n} \nu_i$ is minimized by choosing $\alpha = a_{ii}$.

Proof: See Appendix 2.

Remark 2: It is noted that if $\alpha = 1$ then (40) is changed to $a_{ii} > 1/2$ for all i which is identical to the condition shown in Theorem 4.1 in [16, p. 632].

Remark 3: As a simple example, we note that when a secondorder digital filter has complex conjugate poles, the diagonal elements of a $n \times n$ matrix **A** are equal, i.e.

$$\boldsymbol{A} = \begin{bmatrix} r\cos\theta & r\sin\theta\\ -r\sin\theta & r\cos\theta \end{bmatrix}$$

where the poles are given by $re^{\pm j\theta}$ with 0 < r < 1. In this case, we can choose $\alpha = r \cos \theta$.

It is known [16, p. 632] that for the filter in (1), there exists an internally balanced realization $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_n$ with

$$A = \begin{bmatrix} A_{11} & A_{12} \\ -A_{12}^T & A_{22} \end{bmatrix} \quad A_{11} = A_{11}^T \quad A_{22} = A_{22}^T.$$
(43)

The next theorem describes a counterpart of Theorem 3 for an internally balanced realization.

Theorem 4: Let the filter in (1) be realized by an internally balanced realization in (43), and let the union of $\lambda(A_{11})$ and $\lambda(\mathbf{A}_{22})$ be denoted by $\{\theta_i\}$. If there exists a scalar α satisfying the inequality

$$0 < \frac{\alpha}{2} \le \min_{i} \{\theta_i\} \quad \left(\max_{i} \{\theta_i\} \le \frac{\alpha}{2} < 0\right) \tag{44}$$

then (41) holds for such an α . *Proof:* Since $A_{ii} = A_{ii}^T$ for i = 1, 2 and since $\{\theta_i\} =$ $\lambda(\mathbf{A}_{11}) \cup \lambda(\mathbf{A}_{22})$, if

$$\min_{i} \{\theta_i\} > 0 \quad \left(\max_{i} \{\theta_i\} < 0\right) \tag{45}$$

then, matrices A_{11} and A_{22} are positive-definite (negative-definite). In general

$$\lambda_{\min}(\boldsymbol{A}) \le \boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x} \le \lambda_{\max}(\boldsymbol{A}) \tag{46}$$

holds for any symmetric matrix A where x is any real vector satisfying $||\mathbf{x}|| = 1$. This yields

$$\min_{i} \{\theta_i\} \le \min_{i} \{a_{ii}\} \quad \left(\max_{i} \{a_{ii}\} \le \max_{i} \{\theta_i\}\right).$$
(47)

From (44) and (47) it follows that

$$0 < \frac{\alpha}{2} < a_{ii} \quad \left(a_{ii} < \frac{\alpha}{2} < 0\right) \tag{48}$$

which satisfy the condition in (40). It then follows from Theorm 3 that (41) holds.

This completes the proof.

Remark 4: It is noted that the substitution of $\alpha = 1$ into (44) yields

$$\frac{1}{2} \le \min_i \{\theta_i\}$$

which corresponds to the condition stated in Theorem 4.2 in [16, p. 632].

Theorem 5: Let λ_i for $i = 1, 2, \dots, n$ be the eigenvalues of the matrix **A** in (1). If there exists an $\alpha > 0$ ($\alpha < 0$) satisfying

$$\bar{\lambda} \ge 1 - \frac{2 - \alpha}{2n} \quad \left(\bar{\lambda} \le -\left(1 - \frac{2 + \alpha}{2n}\right)\right)$$
(49)

then (41) holds for this α , where $\overline{\lambda}$ denotes the mean value of the eigenvalues, i.e., $\overline{\lambda} = (\lambda_1 + \lambda_2 + \dots + \lambda_n)/n$.

Proof: Without loss of generality, we assume that the filter in (1) is input-balanced. Then the diagonal element a_{ii} of matrix A satisfies

$$|a_{ii}| < 1, \qquad i = 1, 2, \cdots, n.$$
 (50)

Since

$$\sum_{i=1}^{n} a_{ii} = \sum_{i=1}^{n} \lambda_i \tag{51}$$

holds in general, it follows from (50) that

$$\min_{i} \{a_{ii}\} + n - 1 > \sum_{i=1}^{n} \lambda_{i} \quad \left(\max_{i} \{a_{ii}\} - (n-1) < \sum_{i=1}^{n} \lambda_{i}\right)$$
(52)

Alternatively, (49) can be written as

$$\sum_{i=1}^{n} \lambda_i \ge n - 1 + \frac{\alpha}{2} \quad \left(\sum_{i=1}^{n} \lambda_i \le -(n-1) + \frac{\alpha}{2}\right) \quad (53)$$

Together (52) and (53) imply that

$$\min_{i}\{a_{ii}\} > \frac{\alpha}{2} \quad \left(\max_{i}\{a_{ii}\} < \frac{\alpha}{2}\right) \tag{54}$$

which satisfy the condition in (40). It then follows from Theorem 3 that (41) holds. This completes the proof.

Remark 5: It is noted that the substitution of $\alpha = 1$ into (49) yields

$$\sum_{i=1}^{n} \lambda_i \ge n - \frac{1}{2}$$

which coincides with the condition stated in [16, p. 632, Th. 4.3].

IV. NOISE REDUCTION BY JOINT OPTIMIZATION OF ERROR FEEDBACK AND COORDINATE TRANSFORMATION

In this section, we consider the problem of joint optimization of a scalar error-feedback matrix αI_n and a coordinate transformation matrix T for roundoff noise minimization under l_2 -norm dynamic-range scaling constraints. The proposed joint optimization will be carried out in an iterative manner. First, a scalar α is obtained by modifying (31) under l_2 -norm dynamic-range scaling constraints. In what follows, the unit noise matrix \boldsymbol{W}_D in (8) with $\boldsymbol{D} = \alpha \boldsymbol{I}_n$ is denoted by \boldsymbol{W}_{α} . Under joint application of a scalar error-feedback and a coordinate transformation, the noise gain I(D) becomes tr[$T^T W_{\alpha} T$]. In order to minimize tr[$T^T W_{\alpha} T$] (with α fixed) over an $n \times n$ nonsingular matrix T subject to the constraints in (15), we define the Lagrange function

$$J(\alpha, \boldsymbol{P}, \lambda) = \operatorname{tr}[\boldsymbol{W}_{\alpha}\boldsymbol{P}] + \lambda(\operatorname{tr}[\boldsymbol{K}_{c}\boldsymbol{P}^{-1}] - n)$$
 (55)

where $P = TT^{T}$ and λ is a Lagrange multiplier. Using the formulas for evaluating matrix gradient [18, p. 275]

$$\frac{\frac{\partial \left(\operatorname{tr}[\boldsymbol{M}\boldsymbol{X}]\right)}{\partial \boldsymbol{X}} = \boldsymbol{M}^{T}}{\frac{\partial \left(\operatorname{tr}[\boldsymbol{M}\boldsymbol{X}^{-1}]\right)}{\partial \boldsymbol{X}}} = -\left[\boldsymbol{X}^{-1}\boldsymbol{M}\boldsymbol{X}^{-1}\right]^{T}$$
(56)

we co

$$\frac{\partial J(\alpha, \boldsymbol{P}, \lambda)}{\partial \boldsymbol{P}} = \boldsymbol{W}_{\alpha} - \lambda \boldsymbol{P}^{-1} \boldsymbol{K}_{c} \boldsymbol{P}^{-1}$$
$$\frac{\partial J(\alpha, \boldsymbol{P}, \lambda)}{\partial \lambda} = \operatorname{tr} \left[\boldsymbol{K}_{c} \boldsymbol{P}^{-1} \right] - n.$$
(57)

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

$$PW_{\alpha}P = \lambda K_c, \quad \text{tr}\left[K_cP^{-1}\right] = n.$$
 (58)

Note that if matrices W > 0 and $M \ge 0$ are symmetric, then the matrix equation PWP = M has the unique solution [19]

$$P = W^{-\frac{1}{2}} \left[W^{\frac{1}{2}} M W^{\frac{1}{2}} \right]^{\frac{1}{2}} W^{-\frac{1}{2}}.$$
 (59)

It follows from (58) that

$$P = \sqrt{\lambda} \boldsymbol{W}_{\alpha}^{-\frac{1}{2}} \left[\boldsymbol{W}_{\alpha}^{\frac{1}{2}} \boldsymbol{K}_{c} \boldsymbol{W}_{\alpha}^{\frac{1}{2}} \right]^{\frac{1}{2}} \boldsymbol{W}_{\alpha}^{-\frac{1}{2}}$$
$$\frac{1}{\sqrt{\lambda}} \operatorname{tr} \left[\boldsymbol{K}_{c} \boldsymbol{W}_{\alpha} \right]^{\frac{1}{2}} = \frac{1}{\sqrt{\lambda}} \left(\sum_{i=1}^{n} \theta_{i} \right) = n \tag{60}$$

where θ_i^2 for $i = 1, 2, \dots, n$ are the eigenvalues of $K_c W_{\alpha}$. Therefore, we obtain

$$\boldsymbol{P} = \frac{1}{n} \left(\sum_{i=1}^{n} \theta_i \right) \boldsymbol{W}_{\alpha}^{-\frac{1}{2}} \left[\boldsymbol{W}_{\alpha}^{\frac{1}{2}} \boldsymbol{K}_c \boldsymbol{W}_{\alpha}^{\frac{1}{2}} \right]^{\frac{1}{2}} \boldsymbol{W}_{\alpha}^{-\frac{1}{2}}.$$
 (61)

Substituting (61) into (55) yields the minimum value of $J(\alpha, \boldsymbol{P}, \lambda)$ as

$$\min_{\boldsymbol{P},\lambda} J(\alpha, \boldsymbol{P}, \lambda) = \frac{1}{n} \left(\sum_{i=1}^{n} \theta_i \right)^2 \tag{62}$$

for a given scalar α .

From (61), the optimal coordinate transformation matrix Tthat minimizes (55) can be obtained in closed form as

$$\boldsymbol{T} = \frac{1}{\sqrt{n}} \left(\sum_{i=1}^{n} \theta_i \right)^{\frac{1}{2}} \boldsymbol{W}_{\alpha}^{-\frac{1}{2}} \left[\boldsymbol{W}_{\alpha}^{\frac{1}{2}} \boldsymbol{K}_c \boldsymbol{W}_{\alpha}^{\frac{1}{2}} \right]^{\frac{1}{4}} \boldsymbol{U}$$
(63)

where U is an arbitrary $n \times n$ orthogonal matrix. From (63) it follows that

$$\overline{\boldsymbol{K}}_{c} = \boldsymbol{T}^{-1} \boldsymbol{K}_{c} \boldsymbol{T}^{-T}$$
$$= n \left(\sum_{i=1}^{n} \theta_{i} \right)^{-1} \boldsymbol{U}^{T} \left[\boldsymbol{W}_{\alpha}^{\frac{1}{2}} \boldsymbol{K}_{c} \boldsymbol{W}_{\alpha}^{\frac{1}{2}} \right]^{\frac{1}{2}} \boldsymbol{U}.$$
(64)

Next, we choose the $n \times n$ orthogonal matrix U such that the matrix \overline{K} in (64) satisfies the l_2 -norm dynamic-range scaling constraints on the state-variables in (15). To this end, we perform the eigenvalue-eigenvector decomposition

$$\left[\boldsymbol{W}_{\alpha}^{\frac{1}{2}}\boldsymbol{K}_{c}\boldsymbol{W}_{\alpha}^{\frac{1}{2}}\right]^{\frac{1}{2}} = \boldsymbol{R}\boldsymbol{\Theta}\boldsymbol{R}^{T}$$
(65)

where

$$\boldsymbol{\Theta} = \operatorname{diag}\{\theta_1, \theta_2, \cdots, \theta_n\}$$
$$\boldsymbol{R}\boldsymbol{R}^T = \boldsymbol{I}_n.$$

Consequently

$$n\left(\sum_{i=1}^{n}\theta_{i}\right)^{-1}\left[\boldsymbol{W}_{\alpha}^{\frac{1}{2}}\boldsymbol{K}_{c}\boldsymbol{W}_{\alpha}^{\frac{1}{2}}\right]^{\frac{1}{2}}=\boldsymbol{R}\boldsymbol{\Lambda}^{-2}\boldsymbol{R}^{T} \qquad (66)$$

where

$$\Lambda = \operatorname{diag}\{\lambda_1, \lambda_2, \cdots, \lambda_n\}$$
$$\lambda_i = \left(\frac{\theta_1 + \theta_2 + \cdots + \theta_n}{n\theta_i}\right)^{\frac{1}{2}},$$
$$i = 1, 2, \cdots, n.$$

Now an $n \times n$ orthogonal matrix S such that

$$\boldsymbol{S}\boldsymbol{\Lambda}^{-2}\boldsymbol{S}^{T} = \begin{bmatrix} 1 & * & \cdots & * \\ * & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & * \\ * & \cdots & * & 1 \end{bmatrix}$$
(67)

can be obtained by numerical manipulations [13, p. 278]. By choosing $U = RS^T$ in (63), the optimal coordinate transformation matrix T satisfying (15) and (62) simultaneously can now be constructed as

$$\boldsymbol{T} = \frac{1}{\sqrt{n}} \left(\sum_{i=1}^{n} \theta_i \right)^{\frac{1}{2}} \boldsymbol{W}_{\alpha}^{-\frac{1}{2}} \left[\boldsymbol{W}_{\alpha}^{\frac{1}{2}} \boldsymbol{K}_c \boldsymbol{W}_{\alpha}^{\frac{1}{2}} \right]^{\frac{1}{4}} \boldsymbol{R} \boldsymbol{S}^T.$$
(68)

This coordinate transformation matrix T is used to minimize $tr[\overline{W}_{\alpha}] = tr[T^T W_{\alpha}T]$ subject to the constraints in (15). This completes the first round of iteration and, if necessary, this process may continue until both T and α converge. Having obtained transformation matrix T, an improved value of scalar α can be obtained using

$$\alpha = \frac{\operatorname{tr} \left[\boldsymbol{T}^{T} \boldsymbol{W}_{o} \boldsymbol{A} \boldsymbol{T} \right]}{\operatorname{tr} \left[\boldsymbol{T}^{T} \boldsymbol{W}_{o} \boldsymbol{T} \right]}$$
$$= \frac{\operatorname{tr} \left[\boldsymbol{W}_{o} \boldsymbol{A} \boldsymbol{P} \right]}{\operatorname{tr} \left[\boldsymbol{W}_{o} \boldsymbol{P} \right]}.$$
(69)

This iterative procedure for minimizing the roundoff noise under l_2 -norm scaling constraints with respect to a scalar parameter α as well as an $n \times n$ symmetric positive-definite matrix P can be summarized as follows:

1) Set
$$i = 1$$
 and
 $P(0) = \text{diag}\{(K_c)_{11}^{-1}, (K_c)_{22}^{-1}, \cdots, (K_c)_{nn}^{-1}\}$

2) Compute a scalar $\alpha(i)$ using

$$\alpha(i) = \frac{\operatorname{tr} \left[\boldsymbol{W}_o \boldsymbol{A} \boldsymbol{P}(i-1) \right]}{\operatorname{tr} \left[\boldsymbol{W}_o \boldsymbol{P}(i-1) \right]}.$$

3) Compute
$$I_{\min}(\alpha(i)I_n) = (1 - \alpha(i)^2) \operatorname{tr}[W_o P(i-1)].$$

4) Replace W_{α} by $W_{\alpha(i)}$ computed using

$$\boldsymbol{W}_{\alpha(i)} = \left(1 + \alpha(i)^2\right) \boldsymbol{Wo} - \alpha(i) \left(\boldsymbol{A}^T \boldsymbol{W}_o + \boldsymbol{W}_o \boldsymbol{A}\right)$$

- 5) Derive matrix P from (61), and take the resulting P as P(i).
- 6) Compute tr[$\boldsymbol{W}_{\alpha(i)}\boldsymbol{P}(i)$].
- 7) Update i to i + 1.
- 8) Repeat from Step 2) until the change in either $I_{\min}[\alpha(i)I_n]$ or tr[$W_{\alpha(i)}P(i)$] becomes negligible.
- 9) Obtain matrix T by using (68).

We have applied the proposed algorithm to quite a number of IIR filters, in all the cases we have tried so far, the noise gain is monotonically decreasing and the algorithm converges to the unique optimal solution P regardless of the initial point chosen. In all the cases that we tried, the number of iterations required has been fairly small. For example, simulation results for a ninth-order IIR low-pass digital filter showed that the algorithm converges after nine iterations. Although this may change for high-order filters, it should not be a concern as the order of IIR filters is usually not very high, and the algorithm is typically applied offline.

Suppose the above algorithm converges after N iterations and the optimal coordinate transformation matrix T(N) has been computed from (65)–(68). Then, according to (13), (14), and (21)–(24), the diagonal error-feedback matrix $D = \text{diag}\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ that minimizes

$$I(\boldsymbol{D}) = \operatorname{tr} \left[\boldsymbol{T}^{T}(N) \boldsymbol{W}_{o} \boldsymbol{T}(N) \right] + \operatorname{tr} \left[\boldsymbol{T}^{T}(N) \boldsymbol{W}_{o} \boldsymbol{T}(N) \boldsymbol{D}^{2} \right] -2 \operatorname{tr} \left[\boldsymbol{T}^{T}(N) \boldsymbol{A}^{T} \boldsymbol{W}_{o} \boldsymbol{T}(N) \boldsymbol{D} \right]$$
(70)

is given by

$$\alpha_i = \frac{\left(\boldsymbol{T}^T(N)\boldsymbol{W}_o\boldsymbol{A}\boldsymbol{T}(N)\right)_{ii}}{\left(\boldsymbol{T}^T(N)\boldsymbol{W}_o\boldsymbol{T}(N)\right)_{ii}}, \quad i = 1, 2, \cdots, n.$$
(71)

This diagonal error-feedback matrix D makes it possible to produce more reduction of the noise gain, i.e.

$$I_{\min}(D) < I_{\min}[\alpha(N)I_n].$$
(72)

V. CASE STUDY

In this section, we present a case study to illustrate the roundoff noise reduction methods proposed in the preceding sections.

A. Low-Pass Digital Filter

The system considered here is a third-order stable low-pass IIR digital filter realized as $(A_o, b_o, c_o, d)_3$ in controllable canonical form

$$\boldsymbol{A}_{o} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0.339\,377 & -1.152\,652 & 1.520\,167 \end{bmatrix}$$
$$\boldsymbol{b}_{o} = \begin{bmatrix} 0 & 0 & 0.437\,881 \end{bmatrix}^{T}$$
$$\boldsymbol{c}_{o} = \begin{bmatrix} 0.212\,964 & 0.293\,733 & 0.718\,718 \end{bmatrix}$$
$$\boldsymbol{d} = 6.595\,92 \times 10^{-2}$$



Fig. 2. Magnitude response of the low-pass digital filter.

whose controllability and observability Gramians are given by

$$\boldsymbol{K}_{c} = \begin{bmatrix} 1.0 & 0.741\,988 & 0.227\,107 \\ 0.741\,988 & 1.0 & 0.741\,988 \\ 0.227\,107 & 0.741\,988 & 1.0 \end{bmatrix}$$
$$\boldsymbol{W}_{o} = \begin{bmatrix} 0.720\,426 & -1.635\,144 & 1.753\,511 \\ -1.635\,144 & 4.551\,538 & -4.194\,133 \\ 1.753\,511 & -4.194\,133 & 5.861\,185 \end{bmatrix}.$$

The normalized cutoff frequency of the filter is about 0.32 and its magnitude response is shown in Fig. 2.

In the case study, three realization schemes of the above statespace digital filter subject to the l_2 -norm dynamic-range scaling constraints will be considered.

• Input-balanced realization $(\boldsymbol{A}_{ib}, \boldsymbol{b}_{ib}, \boldsymbol{c}_{ib}, d)_3$ where

$$\boldsymbol{A}_{ib} = \begin{bmatrix} 0.718\,259 & -0.367\,681 & -0.004\,649 \\ 0.680\,718 & 0.402\,944 & 0.208\,077 \\ -0.032\,955 & -0.796\,728 & 0.398\,964 \end{bmatrix}$$
$$\boldsymbol{b}_{ib} = \begin{bmatrix} -0.590\,672 & 0.575\,294 & 0.452\,733 \end{bmatrix}^{T}$$
$$\boldsymbol{c}_{ib} = \begin{bmatrix} -0.933\,934 & -0.491\,318 & 0.100\,978 \end{bmatrix}$$

with the controllability and observability Gramians

$$K_c^{ib} = \text{diag}\{1.0, 1.0, 1.0\}$$

 $W_a^{ib} = \text{diag}\{2.499\,998, 0.729\,367, 0.049\,748\}$

• Optimal realization $(A_{opt}, b_{opt}, c_{opt}, d)_3$ where

$$\boldsymbol{A}_{\text{opt}} = \begin{bmatrix} 0.476\ 474 & 0.568\ 865 & -0.071\ 856 \\ -0.715\ 925 & 0.476\ 474 & 0.154\ 270 \\ -0.154\ 270 & 0.071\ 856 & 0.567\ 219 \end{bmatrix}$$
$$\boldsymbol{b}_{\text{opt}} = \begin{bmatrix} 0.686\ 292 & 0.112\ 451 & -0.713\ 811 \end{bmatrix}^T$$
$$\boldsymbol{c}_{\text{opt}} = \begin{bmatrix} -0.099\ 639 & -0.608\ 103 & -0.632\ 487 \end{bmatrix}$$

with the controllability and observability Gramians

$$\begin{split} \boldsymbol{K}_{c}^{\mathrm{opt}} &= \begin{bmatrix} 1.0 & -0.036\,160 & -0.541\,747 \\ -0.036\,160 & 1.0 & 0.541\,747 \\ -0.541\,747 & 0.541\,747 & 1.0 \end{bmatrix} \\ \boldsymbol{W}_{o}^{\mathrm{opt}} &= \begin{bmatrix} 0.785\,120 & -0.028\,390 & -0.425\,337 \\ -0.028\,390 & 0.785\,120 & 0.425\,337 \\ -0.425\,337 & 0.425\,337 & 0.785\,120 \end{bmatrix}. \end{split}$$

Without using error feedback, the noise gain I(D) in (17) becomes $I(0) = tr[W_o]$. For comparison purposes, the noise

 TABLE I
 I

 ROUNDOFF NOISE GAIN OF THREE REALIZATIONS WITHOUT ERROR FEEDBACK

Realization	Controllable Canonical	Input-Balanced	Optimal
$I(0) = \operatorname{tr}[\boldsymbol{W}_o]$	11.133150	3.279113	2.355360

TABLE II Roundoff Noise Gain With Optimal Diagonal Error-Feedback Matrices

Realization	Controllable Canonical	Input-Balanced	Optimal
Infinite Precision	4.598541	1.863033	1.520074
3-Bit Quantization	4.627836	1.866150	1.524145
Integer Quantization	5.245937	2.187819	1.922510

gain without using error feedback is computed for each of the three realizations and listed in Table I.

Case A. Noise Reduction Using a Diagonal Error-Feedback Matrix: If a diagonal error-feedback matrix is calculated using (24) for each of the above three realizations, then we obtain

 $D = \text{diag}\{0.826\,041, \, 0.702\,890, \, 0.804\,589\}$

for the controllable canonical realization

 $D = diag\{0.718259, 0.402944, 0.398964\}$

for the input-balanced realization, and

 $D = diag\{0.585\,937, 0.494\,832, 0.689\,722\}$

for the optimal realization. The noise reduction performance for these three realizations with optimal diagonal error-feedback is given in Table II, where the noise gains with the diagonal elements quantized to 3-bit numbers or integers are also included.

On comparing the above results with that in Table I, it is observed that error feedback with an appropriate diagonal matrix can reduce the roundoff noise considerably even with 3-bit or integer quantization.

Case B. Noise Reduction Using a Scalar Error-Feedback Matrix: We now consider the case where $D = \alpha I_3$ with a scalar α . Applying (31) to the three realizations, we obtain

	0.764 400	for controllable canonical realization
$\alpha = \langle$	0.643280	for input-balanced realization
	0.590 163	for optimal realization.

The noise gain performance with the optimal scalar error feedback (with in finite precision or 3-bit or integer rounding) for the three realizations are shown in Table III.

On comparing the results obtained with that in Tables I and II we see that the performance with an optimal scalar error-feedback degrades only slightly in relative to what an optimal diagonal error-feedback matrix can achieve and remains significantly better than their no-error-feedback counterparts.

Case C. Noise Reduction by Joint Optimization of Scalar Error Feedback and Coordinate Transformation: In this case, we apply the iterative optimization procedure described in Section IV to the three realizations. For each realization, the

TABLE III
ROUNDOFF NOISE GAIN WITH OPTIMAL SCALAR ERROR-FEEDBACK MATRICES

Realization	Controllable Canonical	Input-Balanced	Optimal
Infinite Precision	4.627966	1.922186	1.535005
3-Bit Quantization	4.630275	1.923282	1.537864
Integer Quantization	5.245937	2.339451	1.930626

algorithm converges after seven iterations to the same matrix P° and a scalar $\alpha = 0.647686$. In the case of controllable canonical realization, the coordinate transformation matrix T° is given by

$$\boldsymbol{T}^{\circ} = \begin{bmatrix} -1.973\,853 & -0.153\,371 & -2.328\,357 \\ -0.063\,334 & -1.398\,294 & -1.260\,527 \\ 1.402\,772 & -0.676\,604 & -0.969\,851 \end{bmatrix}.$$

As is expected, the noise gain resulted from the algorithm is identical regardless of the realization to which it applies, and is given by $I_{\min}(\mathbf{D}) = 1.450049$. If $\alpha = 0.647686$ is rounded to power-of-two representation with 3 b after binary point, then, the noise gain is founded to be $I(\mathbf{D}) = 1.451335$.

Next, a refined solution which offers further reduced noise gain is deduced by applying an optimal diagonal error-feedback matrix to the optimized realization, i.e., $(T^{o-1}AT^o, T^{o-1}b, cT^o, d)_3$. The optimal diagonal error-feedback matrix obtained by using (71) is given by

$$D = diag\{0.705402, 0.510713, 0.683277\}$$

which yields $I_{\min}(D) = 1.433755$. The above diagonal errorfeedback matrix after 3-b quantization (power-of-two representation with 3 b after binary point) gives $I_{\min}(D) = 1.438801$, which is less than $I_{\min}(D) = 1.450049$ in the optimal scalar error feedback.

To compare the proposed methods with that reported in [15], [16], we choose the error-feedback matrix $D = I_n$ as in [15], [16], and minimize tr[$T^T W_D T$] with respect to a coordinate transformation matrix T. This gives

$$\min_{\boldsymbol{T}} \operatorname{tr} \left[\boldsymbol{T}^T \mathbf{W}_D \boldsymbol{T} \right] = 1.752\,546$$

which is considerably larger than our results described above.

Case D. Noise Reduction Using a General Error-Feedback Matrix: Finally, we examine the case where D = A or matrix D is given by its approximation in each realization scheme. The performance of optimized general error-feedback matrices after different rounding polices in three state-space realizations as compared with their infinite-precision error-feedback counterparts are summarized in Table IV. More specifically, the table includes (i) the case with infinite precision error feedback; (ii) the case of rounding the optimal infinite-precision coefficients to integers; and (iii) the case of rounding the coefficients to power-of-two representations with 3 b after the binary point.

From the simulations, it is observed that the controllable canonical realization with an infinite-precision full-scale error-feedback matrix D = A yields the smallest noise gain.

TABLE IV ROUNDOFF NOISE GAIN USING GENERAL ERROR-FEEDBACK MATRICES

Realization	Controllable Canonical	Input-Balanced	Optimal
Infinite Precision	0.648188	1.123823	0.779757
3-Bit Quantization	0.662491	1.130837	0.790560
Integer Quantization	2.809319	1.894676	1.540156



Fig. 3. Magnitude response of the high-pass digital filter.

However, as can be expected, the noise gain with this realization scheme is quite sensitive to the variations in the optimized D. On the other hand, the noise gain for the input-balanced and optimal realizations are relatively insensitive to the quantization of D and the best overall performance is obviously offered by the optimal realization scheme.

B. A High-Pass Digital Filter

Now, let us consider a fourth-order stable high-pass filter realized as $(A_o, b_o, c_o, d)_4$ in controllable canonical form

$$\boldsymbol{A}_{o}^{T} = \begin{bmatrix} 0 & 0 & 0 & -0.386\,952 \\ 1 & 0 & 0 & -1.467\,381 \\ 0 & 1 & 0 & -2.496\,662 \\ 0 & 0 & 1 & -2.225\,780 \end{bmatrix}$$
$$\boldsymbol{b}_{o} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.240\,444 \end{bmatrix}$$
$$\boldsymbol{c}_{o}^{T} = \begin{bmatrix} 0.018\,325 \\ -0.219\,056 \\ 0.141\,882 \\ -0.271\,413 \end{bmatrix}.$$

The normalized cutoff frequency of the filter is about 0.66 and its magnitude response is drawn in Fig. 3.

Tables I–IV in the previous example are changed to Tables V–VIII in this example, respectively.

When the iterative optimization procedure described in Section IV is applied to the three realizations, the algorithm converges after seven iterations to a scalar $\alpha = -0.667273$ and $I_{\min}(\mathbf{D}) = 0.472563$ for each realization. If $\alpha = -0.667273$ is rounded to power-of-two representation with 3 b after binary point, then the noise gain is founded to be $I(\mathbf{D}) = 0.474085$.

 TABLE
 V

 ROUNDOFF NOISE GAIN OF THREE REALIZATIONS WITHOUT ERROR FEEDBACK

Realization	Controllable Canonical	Input-Balanced	Optimal
$I(0) = \operatorname{tr}[\boldsymbol{W}_o]$	22.685416	1.125488	0.796551

TABLE VI Roundoff Noise Gain With Optimal Diagonal Error-Feedback Matrices

Realization	Controllable Canonical	Input-Balanced	Optimal
Infinite Precision	8.803179	0.616628	0.469765
3-Bit Quantization	8.837581	0.617109	0.470395
Integer Quantization	9.892990	0.747868	0.535587

TABLE VII ROUNDOFF NOISE GAIN WITH OPTIMAL SCALAR ERROR-FEEDBACK MATRICES

Realization	Controllable Canonical	Input-Balanced	Optimal
Infinite Precision	8.81442	0.641995	0.501540
3-Bit Quantization	8.837581	0.643037	0.501755
Integer Quantization	9.892990	0.775624	0.623585

 TABLE
 VIII

 ROUNDOFF NOISE GAIN USING GENERAL ERROR-FEEDBACK MATRICES

Realization	Controllable Canonical	Input-Balanced	Optimal
Infinite Precision	0.142117	0.268941	0.190752
3-Bit Quantization	0.150512	0.273714	0.195471
Integer Quantization	3.239478	0.636372	0.526276

Next, a refined solution which offers further reduced noise gain is deduced by applying an optimal diagonal error-feedback matrix to the resulting optimized realization. The optimal diagonal error-feedback matrix obtained by using (71) is given by

$$D = diag\{-0.712866, -0.750118, -0.434626, -0.661687\}$$

which yields $I_{\min}(\mathbf{D}) = 0.462\,844$. The above diagonal error-feedback matrix after 3-bit quantization (power-of-two representation with 3 b after binary point) gives $I_{\min}(\mathbf{D}) = 0.463\,942$, which is less than $I_{\min}(\mathbf{D}) = 0.472\,563$ in the optimal scalar error feedback.

The proposed methods can be compared with that reported in [15], [16] by choosing matrix $D = I_n$ and then minimizing tr[$T^T W_D T$] with respect to matrix T. Namely,

$$\min_{\boldsymbol{T}} \operatorname{tr} \left[\boldsymbol{T}^T \boldsymbol{W}_D \boldsymbol{T} \right] = 2.546\,985$$

which is much larger than our results described above.

VI. CONCLUSION

The roundoff noise minimization in state-space digital filters in several typical realization schemes by means of error feedback has been investigated, and general, diagonal, and scalar error-feedback matrices that minimize the normalized noise gain in a given state-space digital filter have been derived. Moreover, the noise minimization problem has also been addressed in scenario where a coordinate transformation and a scalar error-feedback matrix is jointly optimized subject to the usual l_2 -norm dynamic-range scaling constraints. Simulation results in the form of a case study have been presented to illustrate and support our theoretical analysis and proposed algorithms. The proposed method may be regarded as an alternative, but much simpler and more general, approach to Hwang's method for synthesizing the optimal filter structure with minimum roundoff noise.

The extension of the results obtained in this paper to multidimensional case will appear elsewhere.

APPENDIX I PROOF OF THEOREM 2

The optimal coordinate transformation matrix T described by (68) can be expressed in the form

$$T = T_b T_o \tag{A1}$$

where α is set to zero and

$$T_{b} = W_{o}^{-\frac{1}{2}} \left[W_{o}^{\frac{1}{2}} K_{c} W_{o}^{\frac{1}{2}} \right]^{\frac{1}{2}} R$$
$$T_{o} = \frac{1}{\sqrt{n}} \left(\sum_{i=1}^{n} \sigma_{i} \right)^{\frac{1}{2}} R^{T} \left[W_{o}^{\frac{1}{2}} K_{c} W_{o}^{\frac{1}{2}} \right]^{-\frac{1}{4}} RS^{T}.$$

It follows that

$$\boldsymbol{T}_{b}^{-1}\boldsymbol{K}_{c}\boldsymbol{T}_{b}^{-T} = \boldsymbol{I}_{n}, \quad \boldsymbol{T}_{b}^{T}\boldsymbol{W}_{o}\boldsymbol{T}_{b} = \boldsymbol{\Sigma}^{2}.$$
(A2)

Equation (A2) shows that T_b is the coordinate transformation matrix that converts realization $(A, b, c, d)_n$ to an input-balanced realization. If $\alpha = 0$, then, (13), (14), (65), (A1), and (A2) imply that

$$\operatorname{tr}[\overline{\boldsymbol{W}}_{o}\overline{\boldsymbol{A}}] = \operatorname{tr}\left[\boldsymbol{T}_{o}\boldsymbol{T}_{o}^{T}\boldsymbol{\Sigma}^{2}\boldsymbol{T}_{b}^{-1}\boldsymbol{A}\boldsymbol{T}_{b}\right]$$
$$= \frac{1}{n}\left(\sum_{i=1}^{n}\sigma_{i}\right)\operatorname{tr}\left[\boldsymbol{\Sigma}\boldsymbol{T}_{b}^{-1}\boldsymbol{A}\boldsymbol{T}_{b}\right]$$
$$= \frac{1}{n}\left(\sum_{i=1}^{n}\sigma_{i}\right)\left[\sum_{i=1}^{n}\left(\boldsymbol{T}_{b}^{-1}\boldsymbol{A}\boldsymbol{T}_{b}\right)_{ii}\sigma_{i}\right] \quad (A3)$$

since $\Theta = \Sigma$ in this case, where $T_b^{-1}AT_b$ is the system matrix in the input-balanced realization whose *i*th diagonal element $(T_b^{-1}AT_b)_{ii}$ is replaced by $a_{ii}^{(b)}$ for $i = 1, 2, \dots, n$ in this theorem. By replacing W_o and A by \overline{W}_o and \overline{A} , respectively, and then substituting (16) and (A3) into (30)–(32), we obtain the results in (33)–(35). This completes the proof.

APPENDIX II PROOF OF THEOREM 3

Lemma 2: Let $\{\rho_i^2, \rho_i \ge 0\}$ and $\{\delta_i^2, \delta_i \ge 0\}$ be the sets of diagonal elements and eigenvalues of a positive semidefinite symmetric matrix, respectively. Then

$$\sum_{i=1}^{n} \rho_i \ge \sum_{i=1}^{n} \delta_i \tag{A4}$$

where equality holds if and only if matrix \boldsymbol{M} is diagonal, i.e., $\rho_i = \delta_i$ for any *i* [13].

From (8) it is derived that

$$\boldsymbol{W}_{D} = \boldsymbol{D}^{T} \boldsymbol{W}_{o} \boldsymbol{D} - \boldsymbol{A}^{T} \boldsymbol{W}_{o} \boldsymbol{D} - \boldsymbol{D}^{T} \boldsymbol{W}_{o} \boldsymbol{A} + \boldsymbol{W}_{o}.$$
(A5)

By substituting $D = \alpha I_n$ and $W_o = \Sigma^2$ into (A5), the *i*th diagonal element of matrix W_D, w_{ii} for $i = 1, 2, \dots, n$, can be expressed as

$$w_{ii} = (\alpha^2 - 2a_{ii}\alpha + 1)\sigma_i^2 = [(\alpha - a_{ii})^2 + 1 - a_{ii}^2]\sigma_i^2.$$
(A6)

Now, noting that $K_c = I_n$, Lemma 2 and (A6) can be applied to yield

$$\sum_{i=1}^{n} \nu_i \le \sum_{i=1}^{n} w_{ii}^{\frac{1}{2}} = \sum_{i=1}^{n} \sqrt{\alpha^2 - 2a_{ii}\alpha + 1} \,\sigma_i.$$
(A7)

Since the filter in (1) is input-balanced, the Lyapunov equation in (11) becomes

$$\boldsymbol{I}_n = \boldsymbol{A}\boldsymbol{A}^T + \boldsymbol{b}\boldsymbol{b}^T. \tag{A8}$$

This implies that each diagonal element a_{ii} of a stable matrix A satisfies

$$|a_{ii}| < 1, \quad i = 1, 2, \cdots, n.$$
 (A9)

Therefore, for a scalar α satisfying (40), we have $0 < \alpha^2 - 2a_{ii}\alpha + 1 < 1$. This proves the inequality in (41).

In case the diagonal elements of A are equal, the choice of $\alpha = a_{ii}$ for all i yields

$$\sum_{i=1}^{n} \nu_i \le \sqrt{1 - \alpha^2} \sum_{i=1}^{n} \sigma_i \tag{A10}$$

and in this case, the upper bound of $\sum_{i=1}^{n} \nu_i$ in (A10) is the tightest. This completes the proof.

REFERENCES

- H. A. Spang III and P. M. Shultheiss, "Reduction of quantizing noise by use of feedback," *IRE Trans. Commun. Syst.*, vol. 10, pp. 373–380, Dec. 1962.
- [2] T. Thong and B. Liu, "Error spectrum shaping in narrowband recursive digital filters," *IEEE Trans. Acoust. Speech, Signal Processing*, vol. ASSP-25, pp. 200–203, Apr. 1977.
- [3] T. L. Chang and S. A. White, "An error cancellation digital filter structure and its distributed-arithmetic implementation," *IEEE Trans. Circuits Syst.*, vol. CAS-28, pp. 339–342, Apr. 1981.
- [4] D. C. Munson and D. Liu, "Narrowband recursive filters with error spectrum shaping," *IEEE Trans. Circuits Syst.*, vol. CAS-28, pp. 160–163, Feb. 1981.
- [5] 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. ASSP-30, pp. 963–973, Dec. 1982.
- [6] M. Renfors, "Roundoff noise in error-feedback state-space filters," in Proc. Int. Conf. Acoustics, Speech, Signal Processing (ICASSP'83), Apr. 1983, pp. 619–622.

- [7] W. E. Higgins and D. C. Munson, "Optimal and suboptimal error-spectrum shaping for cascade-form digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-31, pp. 429–437, May 1984.
- [8] 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.
- [9] P. P. Vaidyanathan, "On error-spectrum shaping in state-space digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-32, pp. 88–92, Jan. 1985.
- [10] 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.
- [11] S. Y. Hwang, "Roundoff noise in state-space digital filtering: a general analysis," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-24, pp. 256–262, June 1976.
- [12] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed-point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, pp. 551–562, Sept. 1976.
- [13] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-25, pp. 273–281, Aug. 1977.
- [14] 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.
- [15] D. Williamson, "Roundoff noise minimization and pole-zero sensitivity in fixed-point digital filters using residue feedback," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-34, pp. 1210–1220, Oct. 1986.
- [16] G. Li and M. Gevers, "Roundoff noise minimization using delta-operator realizations," *IEEE Trans. Signal Processing*, vol. 41, pp. 629–637, Feb. 1993.
- [17] D. Williamson, "Delay replacement in direct form structures," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 36, pp. 453–460, Apr. 1988.
- [18] L. L. Scharf, *Statistical Signal Processing*. Reading, MA: Addison-Wesley, 1991.
- [19] G. Li, B. D. O. Anderson, M. Gevers, and J. E. Perkins, "Optimal FWL design of state-space digital systems with weighted sensitivity minimization and sparseness considerations," *IEEE Trans. Circuits Syst. I*, vol. 39, pp. 365–377, May 1992.



Takao Hinamoto (M'77–SM'84–F'01) received the B.E. degree from Okayama University, Okayama, Japan, in 1969, the M.E. degree from Kobe University, Kobe, Japan, in 1971, and the Dr. Eng. degree from Osaka University, Osaka, Japan, in 1977, all in electrical engineering.

From 1972 to 1988, he was with the Faculty of Engineering, Kobe University. From 1979 to 1981, he was a Visiting Member of Staff in the Department of Electrical Engineering, Queen's University, Kingston, ON, Canada, on leave from Kobe

University. During 1988–1991, he was a Professor of electronic circuits in the Faculty of Engineering, Tottori University, Tottori, Japan. Since January 1992, he has been a Professor of electronic control in the Department of Electrical Engineering, Hiroshima University, Hiroshima, Japan. His research interests include digital signal processing, system theory, and control engineering. He has more than 270 published papers in these areas and is the co-editor and co-author of *Two-Dimensional Signal and Image Processing* (Tokyo, Japan: SICE, 1996). He served as an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II from 1993 to 1995, and presently serves as an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS I. He was the Guest Editor of the special section of DSP in the August 1998 issue of the *IEICE Transactions on Fundamentals*.

He also served as Chair of the 12th Digital Signal Processing (DSP) Symposium held in Hiroshima in November 1997, sponsored by the DSP Technical Committee of IEICE. Since 1995, he has been a member of the Steering Committee of the IEEE Midwest Symposium on Circuits and Systems, and since 1998, a member of the Digital Signal Processing Technical Committee in the IEEE Circuits and Systems Society. He served as a member of the Technical Program Committee for ISCAS'99. From 1993 to 2000, he served as a senator or member of the Board of Directors in the Society of Instrument and Control Engineers (SICE), and from 1999 to 2001, he was Chair of the Chugoku Chapter of SICE. He played a leading role in establishing the Hiroshima Section of IEEE, and served as the Interim Chair of the section. He is a recipient of the IEEE Third Millennium Medal.



Hiroaki Ohnishi received the B.E. degree in electrical engineering from Hiroshima University, Hiroshima, Japan, in 2001, where he is presently working toward the M.E. degree in electrical engineering.

His research interest is in the fields of digital signal processing.



Wu-Sheng Lu (S'81–M'85–SM'90–F'99) received the undergraduate degree in mathematics from Fudan University, Shanghai, China, in 1964, and the M.S. degree in electrical engineering, and Ph.D. degree in control science from University of Minnesota, Minneapolis, in 1983, and 1984, respectively.

He was a Post-Doctoral Fellow at the University of Victoria, Victoria, BC, Canada, in 1985, and a Visiting Assistant Professor at University of Minnesota in 1986. Since 1987, he has been with University of Victoria where he is currently a

Professor. His teaching and research interests are in the areas of digital signal processing and application of optimization methods. He is the co-author, with A. Antoniou, of *Two-Dimensional Digital Filters* (New York, Marcel Dekker, 1992). He was an Associate Editor of the *Canadian Journal of Electrical and Computer Engineering* in 1989, and its Editor from 1990 to 1992. He served as an Associate Editor for IEEE TRANS. ON CIRCUITS AND SYSTEMS II from 1993 to 1995, and for IEEE TRANS. ON CIRCUITS AND SYSTEMS II from 1999 to 2001. He is presently an Associate Editor for the *International Journal of Multidimensional Systems and Signal Processing*.

Dr. Lu is a Fellow of the Engineering Institute of Canada.