# Jointly Optimized Error-Feedback and Realization for Roundoff Noise Minimization in State-Space Digital Filters

Wu-Sheng Lu Dept. of Elec. and Comp. Eng. University of Victoria, Canada E-mail: wslu@ece.uvic.ca

#### Abstract

Roundoff noise (RN) is known to exist in digital filters and systems under finite-precision operations and can become a critical factor for severe performance degradation in IIR filters and systems. Two classes of methods are available for RN reduction or minimization — one uses statespace coordinate transformation, the other uses error feedback of state variables. In this paper, we propose a method for the joint optimization of error feedback and state-space realization. It is shown that the problem at hand can be solved in an unconstrained optimization setting. With a closed-form formula for gradient evaluation and an efficient quasi-Newton solver, the unconstrained minimization problem can be solved efficiently.

## 1. Introduction

Since the work of [1][2], it has been well understood that the roundoff noise (RN) of infinite-impulse-response (IIR) digital filters under fixed-point arithmetic operations can be substantially reduced by using adequately chosen state-space realizations [3]-[5]. It has also been known that RN reduction can be accomplished by feeding the quantization error of state variables back to the filter's input through a memoryless (constant) error-feedback matrix without affecting the filter's input-output characteristics [6]-[8]. The success of these techniques leads to the consideration of a joint optimization of error feedback and state-space realization so as to achieve greater reduction in RN. It turns out that obtaining such a jointly optimized error feedback and realization requires the solution of a sophisticated constrained nonlinear minimization problem. In [8], an iterative algorithm for the above optimization problem was proposed for IIR filters with scalar error-feedback matrices, but it appears to be inherently difficult to extend the algorithm to the cases where the error-feedback matrices are diagonal or general.

In this paper, the problem of joint optimization of errorfeedback and state-space realization for RN minimization is investigated in a general nonlinear optimization frame-

Takao Hinamoto Graduate School of Engineering Hiroshima University, Japan E-mail: hinamoto@ecl.sys.hiroshima-u.ac.jp

> work where the error-feedback matrix can be a scalar, diagonal, or general matrix. Using linear-algebraic techniques, we convert the constrained optimization problem at hand into an unconstrained problem which can be solved using powerful quasi-Newton algorithms [10]. A nice feature of employing a general optimization setting for our problem is that both the realization optimization [1][2] and the errorfeedback-matrix optimization [8] become special cases of the proposed formulation that explains why digital filters with jointly optimal error feedback and realization always outperform previously reported systems.

### 2. Preliminaries

Let  $(A, b, c, d)_n$  be a minimal state-space realization of a stable IIR digital filter of order n. This realization can be expressed as

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

$$y(k) = cx(k) + du(k)$$
(1b)

where  $A \in \mathcal{R}^{n \times n}$ ,  $b \in \mathcal{R}^{n \times 1}$ ,  $c \in \mathcal{R}^{1 \times n}$ , and  $d \in \mathcal{R}$ . Now assume that the filter is implemented subject to finiteword-length (FWL) constraint and quantization takes place before matrix-vector multiplications, and an error-feedback for state variables is used for the sake of RN reduction, then the filter's model becomes [8]

$$\tilde{\boldsymbol{x}}(k+1) = \boldsymbol{A}Q[\tilde{\boldsymbol{x}}(k)] + \boldsymbol{b}u(k) + \boldsymbol{D}\boldsymbol{e}(k) \qquad (2a)$$

$$\tilde{y}(k) = cQ[\tilde{x}(k)] + du(k)$$
(2b)

where  $Q[\cdot]$  denotes the quantizer that rounds the fraction of each input component to a *b*-bit number, e(k) is the quantization error defined by

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

and D is referred to an error-feedback matrix. Fig. 1 shows a block diagram of the state-space filter described by (2).

From (1) and (2), the roundoff noise for the filter can be modeled as

$$\Delta \boldsymbol{x}(k+1) = \boldsymbol{A} \Delta \boldsymbol{x}(k) + (\boldsymbol{A} - \boldsymbol{D})\boldsymbol{e}(k) \quad (3a)$$

$$\Delta y(k) = c\Delta x(k) + ce(k)$$
(3b)



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

where  $\Delta \boldsymbol{x}(k) = \boldsymbol{x}(k) - \tilde{\boldsymbol{x}}(k)$  and  $\Delta y(k) = y(k) - \tilde{y}(k)$ . In the frequency domain, the noise process is modeled by

$$\Delta Y(z) = \boldsymbol{G}_D(z)\boldsymbol{E}(z) \tag{4a}$$

$$\boldsymbol{G}_D(z) = \boldsymbol{c}(z\boldsymbol{I} - \boldsymbol{A})^{-1}(\boldsymbol{A} - \boldsymbol{D}) + \boldsymbol{c} \qquad (4b)$$

where  $\Delta Y(z)$  and E(z) are the z-transforms of  $\Delta y(k)$  and e(k), respectively, and  $G_D(z)$  denotes the transfer function from the quantization error to output roundoff noise. Therefore, a *noise gain* due to quantization error can be defined as

where

(5a)

$$\boldsymbol{W}_{D} = \frac{1}{2\pi j} \oint_{|z|=1} \boldsymbol{G}_{D}(z) \boldsymbol{G}_{D}^{*}(z) \frac{dz}{z}$$
(5b)

It is known that the matrix  $W_D$  in (5b) can be expressed as [8]

$$W_D = (\boldsymbol{A} - \boldsymbol{D})^T W_o (\boldsymbol{A} - \boldsymbol{D}) + \boldsymbol{c}^T \boldsymbol{c}$$
 (6)

where  $W_o$  is the observability Gramian of the filter and can be computed by solving the Lyapunov equation [11]

$$\boldsymbol{W}_{o} - \boldsymbol{A}^{T} \boldsymbol{W}_{o} \boldsymbol{A} = \boldsymbol{c}^{T} \boldsymbol{c}$$
(7)

# 3. Joint Optimization of Error-Feedback and Realization

#### **3.1. An Optimization Formulation**

 $I(\boldsymbol{D}) = \operatorname{tr}(\boldsymbol{W}_{\mathcal{D}})$ 

As is well-known [11], the state-space realizations that are equivalent to a particular realization of a given digital filter, say  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_n$ , are characterized by  $(\overline{\mathbf{A}}, \overline{\mathbf{b}}, \overline{\mathbf{c}}, \overline{d})_n$  $= (\mathbf{T}^{-1}\mathbf{A}\mathbf{T}, \mathbf{T}^{-1}\mathbf{b}, \mathbf{c}\mathbf{T}, d)_n$  where  $\mathbf{T} \in \mathbb{R}^{n \times n}$  is a nonsingular coordinate transformation matrix. Under a transformation  $\mathbf{T}$ , the noise gain defined in (5) becomes

$$I(\boldsymbol{D}) = \operatorname{tr}(\overline{\boldsymbol{W}}_D) \tag{8a}$$

where

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

$$\overline{W}_o = T^T W_0 T \tag{8c}$$

A basic constraint imposed on RN minimization is the  $l_2$ norm dynamic range of the state variables [1][2]. Under a coordinate transformation, this constraint can be expressed as

$$(\overline{K}_c)_{ii} = 1$$
 for  $1 \le i \le n$  (9a)

where

$$\overline{\boldsymbol{K}}_c = \boldsymbol{T}^{-1} \boldsymbol{K}_c \boldsymbol{T}^{-T}$$
(9b)

and  $K_c$  is the controllability Gramian of the original realization and can be computed by solving the Lyapunov equation [11]

$$\boldsymbol{K}_c - \boldsymbol{A}\boldsymbol{K}_c \boldsymbol{A}^T = \boldsymbol{b}\boldsymbol{b}^T \tag{10}$$

Primarily we are concerned with the minimization of the roundoff noise gain I(D) in (8) subject to the constraints in (9). For the sake of reducing solution sensitivity (as will be detailed shortly), however, the "magnitude" of the observability Gramian,  $\overline{W}_o$ , needs to be controlled during the optimization process. To this end the objective function needs to be modified and the constrained minimization problem can be described as

$$\underset{\boldsymbol{D}, \boldsymbol{T}}{\text{minimize}} \quad J(\boldsymbol{D}, \boldsymbol{T}) = \text{tr}[(1-\mu)\overline{\boldsymbol{W}}_D + \mu \overline{\boldsymbol{W}}_o] \quad (11a)$$

subject to: 
$$(\boldsymbol{T}^{-1}\boldsymbol{K}_{c}\boldsymbol{T}^{-T})_{ii} = 1$$
 for  $1 \le i \le n$  (11b)

where  $\overline{W}_D$  and  $\overline{W}_o$  are given in (8), and  $0 \le \mu \le 1$  is a scalar that weighs the importance of reducing tr( $\overline{W}_o$ ) relative to reducing tr( $\overline{W}_D$ ).

The problem formulation in (11) is rather general. As a matter of fact, it includes the following two special cases: (i) if error feedback is not used, then D is set to zero, which in conjunction with (8b) and (7) implies that

$$J(\mathbf{0}, \mathbf{T}) = \operatorname{tr}(\overline{\mathbf{W}}_o) \tag{12}$$

Several methods of minimizing J(0, T) in (12) subject to (11b) were investigated in [1][2]. (ii) to minimize the objective function J(D, T) for a fixed T (i.e., for a given state-space realization). This problem has been addressed in [8]. Consequently, the solution of the general optimization problem in (11) that finds the jointly optimized errorfeedback matrix D and coordinate transformation matrix Tis expect to be superior to the solutions obtained from these two special cases.

#### 3.2. An Equivalent Unconstrained Problem

Since the IIR filter at hand is assumed to be stable, controllable and observable, the controllability matrix  $K_c$  is positive definite [11]. Let  $K_c^{1/2}$  denote the symmetric square root of  $K_c$ , i.e.,  $K_c^{1/2}$  is a symmetric matrix satisfying  $K_c^{1/2}K_c^{1/2} = K_c$ , then  $K_c^{1/2}$  is also positive define and we can define

$$\hat{\boldsymbol{T}} = \boldsymbol{T}^T \boldsymbol{K}_c^{-1/2} \tag{13}$$

which implies that  $T^{-1} = \hat{T}^{-T} K_c^{-1/2}$  and the constraints in (11b) become

$$(\hat{\boldsymbol{T}}^{-T}\hat{\boldsymbol{T}}^{-1})_{ii} = 1$$
 for  $1 \le i \le n$  (14)

The constraints in (14) simply mean that each cloumn in  $\hat{T}^{-1}$  must be a unity vector. This can be satisfied if  $\hat{T}^{-1}$  assumes the form

$$\hat{\boldsymbol{T}}^{-1} = \begin{bmatrix} \boldsymbol{t}_1 & \boldsymbol{t}_2 \\ \|\boldsymbol{t}_1\| & \|\boldsymbol{t}_2\| & \cdots & \|\boldsymbol{t}_n\| \end{bmatrix}$$
 (15)

with  $t_i \in \mathcal{R}^{n \times 1}$ . To complete our problem conversion, we need to re-write the objective function in (11a) in terms of D and  $\hat{T}$ , and this can be done as follows.

From (8) and (13), we can write

$$\overline{W}_{D} = (T^{-1}AT - D)^{T}T^{T}W_{o}T(T^{-1}AT - D) + T^{T}c^{T}cT \\
= \hat{T}[(\hat{A} - \hat{T}^{T}D\hat{T}^{-T})^{T}\hat{W}_{o}(\hat{A} - \hat{T}^{T}D\hat{T}^{-T}) + \hat{C}]\hat{T}^{T} (16a)$$

where

$$\hat{A} = K_c^{-1/2} A K_c^{1/2}$$
 (16b)

$$\hat{C} = K_c^{1/2} c^T c K_c^{1/2}$$
 (16c)

$$\hat{W}_o = K_c^{1/2} W_o K_c^{1/2}$$
 (16d)

and the objective function in (11a) becomes

$$J(\boldsymbol{D}, \hat{\boldsymbol{T}}) = \operatorname{tr}\{\hat{\boldsymbol{T}}[(1-\mu)(\hat{\boldsymbol{A}} - \hat{\boldsymbol{T}}^{T}\boldsymbol{D}\hat{\boldsymbol{T}}^{-T})^{T}\hat{\boldsymbol{W}}_{o} \\ (\hat{\boldsymbol{A}} - \hat{\boldsymbol{T}}^{T}\boldsymbol{D}\hat{\boldsymbol{T}}^{-T}) + (1-\mu)\hat{\boldsymbol{C}} + \mu\hat{\boldsymbol{W}}_{o}]\hat{\boldsymbol{T}}^{T}\}$$
(17)

which can also be expressed as

$$J(\boldsymbol{D}, \hat{\boldsymbol{T}}) = (1 - \mu)(J_1 + J_2) + \mu J_3$$
(18a)

with  

$$J_{1} = \text{tr}[\hat{T}(\hat{A} - \hat{T}^{T}D\hat{T}^{-T})^{T}\hat{W}_{o}(\hat{A} - \hat{T}^{T}D\hat{T}^{-T})\hat{T}^{T}](18\text{b})$$

$$J_2 = \operatorname{tr}(\hat{\boldsymbol{T}}\hat{\boldsymbol{C}}\hat{\boldsymbol{T}}^{T}) \tag{18c}$$

$$J_3 = \operatorname{tr}(\hat{\boldsymbol{T}}\hat{\boldsymbol{W}}_o\hat{\boldsymbol{T}}^T) \tag{18d}$$

Because  $\hat{W}_o$  is positive definite and  $\hat{C}$  is positive semidefinite, we have  $J_i \ge 0$  for i = 1, 2, 3. It follows immediately that if the error-feedback matrix D is allowed to be a general matrix, then the optimal choice of D is

$$\boldsymbol{D} = \hat{\boldsymbol{T}}^{-T} \hat{\boldsymbol{A}} \hat{\boldsymbol{T}}^{T}$$
(19)

as it leads to  $J_1 = 0$ . In other words, in the case of D being a general matrix, the objective function is simplified to

$$J(\hat{\boldsymbol{T}}^{-T}\hat{\boldsymbol{A}}\hat{\boldsymbol{T}}^{T},\hat{\boldsymbol{T}}) = \operatorname{tr}\{\hat{\boldsymbol{T}}[(1-\mu)\hat{\boldsymbol{C}}+\mu\hat{\boldsymbol{W}}_{o}]\hat{\boldsymbol{T}}^{T}\} \quad (20)$$

Another case that entails a simplified  $J(D, \hat{T})$  is when D is a scalar matrix, i.e.,  $D = \alpha I$ , which leads  $J(D, \hat{T})$  to

$$J(\alpha \boldsymbol{I}, \hat{\boldsymbol{T}}) = \operatorname{tr}\{\hat{\boldsymbol{T}}[(1-\mu)(\hat{\boldsymbol{A}}-\alpha \boldsymbol{I})^T \hat{\boldsymbol{W}}_o(\hat{\boldsymbol{A}}-\alpha \boldsymbol{I}) + (1-\mu)\hat{\boldsymbol{C}}+\mu \hat{\boldsymbol{W}}_o]\hat{\boldsymbol{T}}^T\}$$
(21)

In summary, the joint optimization problem in (11) is now re-formulated as

$$\underset{\boldsymbol{D},\,\hat{\boldsymbol{T}}}{\text{minimize}} \ J(\boldsymbol{D},\hat{\boldsymbol{T}}) \tag{22}$$

where  $J(D, \hat{T})$  assumes the form in (20) if D is a general matrix, the form in (21) if D is a scalar matrix, and otherwise the form in (17); and  $\hat{T}$  assumes the form in (15). From above discussion, it is quite clear that the variables in problem (22) consist of the vectors  $t_1, t_2, \dots, t_n$  (see (15)) plus certain entries of D: for example it would include  $\alpha$  if  $D = \alpha I$ , the *n* diagonal elements of *D* if *D* is a diagonal matrix, but no entries of D need to be included if D is a general matrix. It should be emphasized that although the vectors  $\{t_i, 1 \leq i \leq n\}$  have to be such that  $\hat{T}$  is nonsingular, this type of "constraint" needs not to be imposed explicitly because a near singular T would make the value of  $J(\boldsymbol{D}, \boldsymbol{T})$  very large, hence the process of minimizing  $J(\boldsymbol{D}, \hat{\boldsymbol{T}})$  automatically avoids considering ill-conditioned  $\hat{T}$ . Consequently, the problem in (22) is practically an *un*constrained minimization problem.

Finally, a remark on the term  $\mu J_3$  in (18): from (8) it is clear that if  $D \neq \overline{A}$  then the contribution from term  $(\overline{A} - D)^T \overline{W}_o (\overline{A} - D)$  to the noise gain is critically depending on the "magnitude" of  $\overline{W}_o$ . Since in a practical implementation D is of finite precision,  $\overline{A} - D$  is always a nonzero matrix and therefore the magnitude of  $\overline{W}_o$  (in terms of its norm, for example) should be controlled. In particular, when D is allowed to be a general matrix, the objective function then becomes the form in (20) and  $\mu J_3$  is the only term there to control the value of  $J_3$  (which turns out to be the Frobenius norm of  $\overline{W}_o^{1/2}$ ). This necessitates the use of a nonzero  $\mu$  for the objective function in (20). In other cases such as those in (17) and (21), the term equivalent to  $(\overline{A} - D)^T \overline{W}_o(\overline{A} - D)$  always presents, and the use of  $\mu = 0$  would not in general lead to ill-conditioned results.

#### 3.3. A Quasi-Newton Algorithm for Problem (22)

Let x be the column vector that collects the variables in D and  $\hat{T}$ , thus  $J(D, \hat{T})$  is a function of x, denoted by J(x). The algorithm starts with a trivial initial point  $x_0$ obtained by letting D = I and  $\hat{T} = I$ . Now suppose we are in the *k*th iteration to update the most recent point  $x_k$ . A quasi-Newton algorithm updates  $x_k$  to  $x_{k+1}$  as

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \lambda_k \boldsymbol{d}_k \tag{23}$$

in two steps: (i) Determine a search direction  $d_k = -S_k g_k$ where  $g_k = \nabla J(x)$  is the gradient of the objective function and  $S_k$  is a positive-definite approximation of the inverse Hessian matrix of J(x). A popular quasi-Newton algorithm is the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm [10] which updates  $S_k$  through the recursive relation

$$\boldsymbol{S}_{k+1} = \boldsymbol{S}_{k} + \left(1 + \frac{\boldsymbol{\gamma}_{k}^{T} \boldsymbol{S}_{k} \boldsymbol{\gamma}_{k}}{\boldsymbol{\gamma}_{k}^{T} \boldsymbol{\delta}_{k}}\right) \frac{\boldsymbol{\delta}_{k} \boldsymbol{\delta}_{k}^{T}}{\boldsymbol{\gamma}_{k}^{T} \boldsymbol{\delta}_{k}} - \frac{\left(\boldsymbol{\delta}_{k} \boldsymbol{\gamma}_{k}^{T} \boldsymbol{S}_{k} + \boldsymbol{S}_{k} \boldsymbol{\gamma}_{k} \boldsymbol{\delta}_{k}^{T}\right)}{\boldsymbol{\gamma}_{k}^{T} \boldsymbol{\delta}_{k}}$$
(24)

where  $S_0 = I$ ,  $\delta_k = x_{k+1} - x_k$ , and  $\gamma_k = g_{k+1} - g_k$ , (ii) Once the search direction  $d_k$  is computed, the onedimensional optimization (often called line search)

$$\lambda_k = \arg\min_{\lambda} J(\boldsymbol{x}_k + \lambda \boldsymbol{d}_k)$$
 (25)

is carried out to determine the value of  $\lambda_k$ . If the iteration progress measured by  $\|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\|$ , is greater than a prescribed tolerance  $\varepsilon$ , then set k := k + 1 and repeat from Step (i), otherwise the iteration is terminated and  $\boldsymbol{x}_{k+1}$  is claimed to be a solution point.

The implementation of (24) requires the gradient of J(x). Closed form expressions for J(x) with scalar, diagonal, and general error-feedback matrix D are given in Appendix A.

#### 3.4. Examples

We present two examples to illustrate the proposed optimization method. The first example concerns a 3rd-order lowpass IIR digital filter which was also used in [8]. The second example is about a 9th-order lowpass IIR filter which is used to demonstrate the ability of the proposed algorithm to deal with relatively large number of variables.

*Example 1* Consider a 3rd-order stable IIR lowpass digital filter whose controllable canonical realization is denoted by  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_3$  with

$$\boldsymbol{A} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0.339377 & -1.152652 & 1.52016 \end{bmatrix}$$
(26a)

$$\boldsymbol{b} = \begin{bmatrix} 0 & 0 & 0.4437881 \end{bmatrix}^T$$
(26b)

$$\boldsymbol{c} = \begin{bmatrix} 0.212964 & 0.293733 & 0.718718 \end{bmatrix}$$
(26c)

$$d = 6.59592 \times 10^{-2} \tag{26d}$$

The controllability Gramian  $K_c$  of the above filter has been normalized to satisfy the constraints  $(K_c)_{ii} = 1$  for i =1, 2, 3. Without error feedback (i.e., D = 0), the noise gain of filter (26) was found to be [8]

$$tr(W_o) = 11.1332$$

If one applies the method of [1][2] to filter (26) to obtain a realization  $(\overline{A}, \overline{b}, \overline{c}, d)_3$  for roundoff noise minimization (without error feedback), then the noise gain was reduced to [8]

$$\operatorname{tr}(\overline{W}_o) = 2.3554$$

Next, we compute jointly optimal error-feedback and statespace realization, with D being scalar, diagonal, and general matrices, by solving the respective minimization problem (22) using BFGS updates. The total number of variables involved in the optimization are 10 (for a scalar D), 12 (for a diagonal D), and 9 (for a general D). For the cases where D is either a scalar or diagonal matrix,  $\mu = 0$  was used. For the case of D being a general matrix,  $\mu = 0.01$ was assumed. In all three cases the initial point used corresponds to D = I and  $\hat{T} = I$ , and with  $\varepsilon = 10^{-8}$ the algorithm converges with less than 20 iterations. The minimized noise gains obtained are given in Table 1. For comparison purposes, Table 1 also includes the noise gain values obtained in [8] where the error-feedback matrix is optimized for a fixed state-space realization that is optimal (without error feedback) for roundoff noise minimization.

These values are listed in Table 1 in the lines where "Separate" is indicated for the column "Joint/Separate". Form the simulation results, it is evident that the proposed joint optimization offers improved performance for RN reduction for all three types of D.

Table 1. Performance Comparison for Example 1

Realization	Error	Joint/Separate	Noise Gain
	Feedback	Optimization	
Canonical	0		11.1332
Optimal	0	_	2.3554
Optimal	Scalar	Separate	1.5350
Optimal	Scalar	Joint	1.4500
Optimal	Diagonal	Separate	1.4338
Optimal	Diagonal	Joint	1.3090
Optimal	General	Separate	0.7798
Optimal	General	Joint	0.4208

*Example 2* Now we consider a 9th-order stable IIR lowpass filter whose transfer function is denoted by

$$H(z) = \frac{b_1 z^9 + b_2 z^8 + \dots + b_9 z + b_{10}}{a_1 z^9 + a_2 z^8 + \dots + a_9 z + a_{10}}$$

where the coefficients are given in Table 2.

Table 2. Coefficients of H(z)i $a_i$  $b_i$ 1 1 0.002198 2 3.640015 -0.0079933 7.148374 0.010478 4 9.521133 0.007251 5 9.297299 0.011297 6 -6.830931 -0.0105823.754299 -0.017728 7 8 -1.4851090.023996 9 0.384233 0.006712 10 -0.049851 0.046116

Next we obtain the controllable canonical realization of the filter and normalize its controllability matrix by scaling so as to satisfy the constraints  $(\mathbf{K}_c)_{ii} = 1$  for  $i = 1, 2, \ldots, 9$ . The state-space realization obtained is denoted by  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_9$  and its noise gain (without error feedback) was found to be

$$\operatorname{tr}(\boldsymbol{W}_o) = 3.1354 \times 10^3$$

Without using error feedback, the method of [1][2] was applied to obtain a realization  $(\overline{A}, \overline{b}, \overline{c}, \overline{d})_9$  for RN minimization, whose noise gain was reduced to

$$\operatorname{tr}(\overline{W}_o) = 2.5315$$

The proposed joint optimization method was then applied to the controllable canonical realization with error-feedback matrix D being scalar, diagonal, and general matrices. For the cases of scalar and diagonal D, the weighting factor  $\mu$ was set to zero; while for the case of a general D,  $\mu = 0.03$ was used. As in Example 1, for all three cases the same initial point, which corresponds to the choice of D = I,  $\hat{T} = \hat{I}$ , was used. With  $\hat{\varepsilon} = 10^{-4}$ , it took the algorithm 35 (for a scalar D), 196 (for a diagonal D), and 54 (for a general D) iterations to converge. The minimized noise gains are given in Table 3. Again, for comparison purposes, Table 3 also lists the noise gains obtained using separate realization and error-feedback matrix optimization proposed in [8]. From the table, it is observed that the performance improvement provided by using the joint optimization appears to be more pronounced. Based on this and a large number of simulations conducted so far, we conclude that the proposed joint optimization can offer improved performance gain for IIR state-space digital filters of relatively high order.

Table 3. Performance Comparison for Example 2

Realization	Error	Joint/Separate	Noise Gain
	Feedback	Optimization	
Canonical	0		$3.1354 \times 10^{3}$
Optimal	0	_	2.5315
Optimal	Scalar	Separate	1.3622
Optimal	Scalar	Joint	1.1628
Optimal	Diagonal	Separate	1.2995
Optimal	Diagonal	Joint	0.9866
Optimal	General	Separate	0.2776
Optimal	General	Joint	0.0868
Optimal Optimal	General General	Separate Joint	0.2776

### **Appendix A** Evaluation of $\nabla J(\boldsymbol{x})$

Depending on the type of matrix D, the objective function J(x) may assume one of the three expressions in (17), (20), and (21). In what follows, the derivation of  $\nabla J(x)$  is carried out for two separate cases.

#### A.1 If D is a general or a scalar matrix

The objective function in both (20) and (21) has the form  $J(\boldsymbol{x}) = \operatorname{tr}(\hat{\boldsymbol{T}}\boldsymbol{M}\hat{\boldsymbol{T}}^T)$  which, in the light of (15), can be expressed as

$$J(\boldsymbol{x}) = \operatorname{tr} \left\{ \begin{bmatrix} \boldsymbol{t}_1 & \cdots & \boldsymbol{t}_n \\ \|\boldsymbol{t}_1\| & \cdots & \|\boldsymbol{t}_n\| \end{bmatrix}^{-1} M \begin{bmatrix} \boldsymbol{t}_1 & \cdots & \boldsymbol{t}_n \\ \|\boldsymbol{t}_1\| & \cdots & \|\boldsymbol{t}_n\| \end{bmatrix}^{-T} \right\}$$
(A1)

In the case of *D* being a general matrix, *M* is a constant matrix (see (20)) and *x* contains a total of n<sup>2</sup> variables, i.e., *t*<sub>1</sub>, *t*<sub>2</sub>, ..., *t<sub>n</sub>*. To compute ∂*J*(*x*)/∂*t<sub>ij</sub>*, we perturb the *i*th component of vector *t<sub>j</sub>* by a small amount, say δ, and keep the rest of *T̂* unchanged. If we denote the perturbed *j*th column of *T̂*<sup>-1</sup> by *t̃*<sub>j</sub>/|*t̃*<sub>j</sub>||, then we can write a linear approximation of *t̃*<sub>j</sub>/|*t̃*<sub>j</sub>|| as

$$rac{ ilde{oldsymbol{t}}_j}{\| ilde{oldsymbol{t}}_j\|}pprox rac{oldsymbol{t}_j}{\|oldsymbol{t}_j\|}-\deltaoldsymbol{g}_{ij}$$

where  $g_{ij}$  is a vector given by

$$g_{ij} = \frac{1}{\|t_j\|^3} (t_{ij}t_j - \|t_j\|^2 e_i)$$
 (A2)

and  $e_i$  is the *i*th coordinate vector. Now let  $\hat{T}_{ij}$  be the matrix obtained from  $\hat{T}$  with a perturbed (i, j)th component, then up to the first order the matrix inversion formula [11, p. 655] gives

$$\hat{\boldsymbol{T}}_{ij} = \hat{\boldsymbol{T}} + rac{\delta(\hat{\boldsymbol{T}} \boldsymbol{g}_{ij})(\boldsymbol{e}_j^T \hat{\boldsymbol{T}})}{1 - \delta \boldsymbol{e}_j^T \hat{\boldsymbol{T}} \boldsymbol{g}_{ij}}$$

Consequently, we have

$$\frac{\partial J(\boldsymbol{x})}{\partial t_{ij}} = \lim_{\delta \to 0} [\operatorname{tr}(\hat{\boldsymbol{T}}_{ij}\boldsymbol{M}\hat{\boldsymbol{T}}_{ij}^{T}) - \operatorname{tr}(\hat{\boldsymbol{T}}\boldsymbol{M}\hat{\boldsymbol{T}}^{T})]/\delta$$
$$= 2\operatorname{tr}[(\hat{\boldsymbol{T}}\boldsymbol{g}_{ij})(\boldsymbol{e}_{j}^{T}\hat{\boldsymbol{T}})\boldsymbol{M}\hat{\boldsymbol{T}}^{T}]$$
$$= 2\boldsymbol{e}_{j}^{T}(\hat{\boldsymbol{T}}\boldsymbol{M}\hat{\boldsymbol{T}}^{T}\hat{\boldsymbol{T}})\boldsymbol{g}_{ij} \text{ for } 1 \leq i, j \leq n \quad (A3)$$

• If D is a scalar matrix,  $D = \alpha I$ , then vector x contains a total of  $n^2 + 1$  variables, i.e.,  $\alpha$ ,  $t_1$ ,  $t_2$ , ...,  $t_n$ . In this case  $\partial J(x)/\partial t_{ij}$  can also be computed using (A3), and it follows from (21) that

$$\frac{\partial J(\boldsymbol{x})}{\partial \alpha} = (1 - \mu) \operatorname{tr}[\hat{\boldsymbol{T}}(2\alpha \hat{\boldsymbol{W}}_o - \hat{\boldsymbol{A}}^T \hat{\boldsymbol{W}}_o - \hat{\boldsymbol{W}}_o \hat{\boldsymbol{A}}) \hat{\boldsymbol{T}}^T]$$
(A4)

# A.2 If matrix *D* contains certain number of zero components in fixed places

The type of D matrices we deal with here obviously includes the case of D being a diagonal matrix. To evaluate the gradient of  $J(D, \hat{T})$  in (17), we write it as

$$J(\boldsymbol{x}) = \operatorname{tr}(\hat{\boldsymbol{T}}\boldsymbol{M}\hat{\boldsymbol{T}}^{T}) + (1-\mu)[\operatorname{tr}(\boldsymbol{D}^{T}\hat{\boldsymbol{T}}\hat{\boldsymbol{W}}_{o}\hat{\boldsymbol{T}}^{T}\boldsymbol{D}) - 2\operatorname{tr}(\hat{\boldsymbol{T}}\hat{\boldsymbol{A}}^{T}\hat{\boldsymbol{W}}_{o}\hat{\boldsymbol{T}}^{T}\boldsymbol{D})]$$
(A5)

where

$$\boldsymbol{M} = (1-\mu)(\hat{\boldsymbol{A}}^T \hat{\boldsymbol{W}}_o \hat{\boldsymbol{A}} + \hat{\boldsymbol{C}}) + \mu \hat{\boldsymbol{W}}_o$$

and x in this case contains the nonzero entries of D plus vectors  $t_1, t_2, \ldots, t_n$ . To compute  $\partial J(x)/\partial t_{ij}$ , we treat all the quantities other than  $t_{ij}$  in (A5) including D as constant terms. It then follows from Sec. A.1 that

$$\frac{\partial J(\boldsymbol{x})}{\partial t_{ij}} = 2\beta_1 + 2(1-\mu)(\beta_2 - \beta_3) \qquad \text{for } 1 \le i, j \le n$$
(A6)

with

$$\beta_1 = \boldsymbol{e}_j^T (\hat{\boldsymbol{T}} \boldsymbol{M} \hat{\boldsymbol{T}}^T \hat{\boldsymbol{T}}) \boldsymbol{g}_{ij} \tag{A7}$$

$$\beta_2 = \boldsymbol{e}_j^T (\hat{\boldsymbol{T}} \hat{\boldsymbol{W}}_o \hat{\boldsymbol{T}}^T \boldsymbol{D} \boldsymbol{D}^T \hat{\boldsymbol{T}}) \boldsymbol{g}_{ij}$$
(A8)

$$\beta_3 = \boldsymbol{e}_j^T \hat{\boldsymbol{T}} (\hat{\boldsymbol{A}}^T \hat{\boldsymbol{W}}_o \hat{\boldsymbol{T}}^T \boldsymbol{D} + \hat{\boldsymbol{W}}_o \hat{\boldsymbol{A}} \hat{\boldsymbol{T}}^T \boldsymbol{D}^T) \hat{\boldsymbol{T}} \boldsymbol{g}_{ij} (A9)$$

Finally, using (A5) we compute for  $\boldsymbol{D} = \{d_{ij}\}$  the derivative

$$\frac{\partial J(\boldsymbol{x})}{\partial d_{ij}} = 2(1-\mu)\boldsymbol{e}_j^T (\boldsymbol{D}^T \hat{\boldsymbol{T}} - \hat{\boldsymbol{T}} \hat{\boldsymbol{A}}^T) \hat{\boldsymbol{W}}_o \hat{\boldsymbol{T}}^T \boldsymbol{g}_{ij} \quad (A10)$$

In particular, if D is a diagonal matrix, i.e.,  $D = diag\{d_1, d_2, \ldots, d_n\}$ , then

$$\frac{\partial J(\boldsymbol{x})}{\partial d_i} = 2(1-\mu)\boldsymbol{e}_i^T (\boldsymbol{D}^T \hat{\boldsymbol{T}} - \hat{\boldsymbol{T}} \hat{\boldsymbol{A}}^T) \hat{\boldsymbol{W}}_o \hat{\boldsymbol{T}}^T \boldsymbol{g}_{ii} \quad (A11)$$

# References

- S. Y. Hwang, "Roundoff noise in state-space digital filtering: A general analysis," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 24, pp. 254-262, June 1976.
- [2] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed-point digital filters," *IEEE Trans. Circuits syst.*, vol. 23, pp. 551-562, Sept. 1976.

- [3] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 25, pp. 273-281, Aug. 1977.
- [4] 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. 26, pp. 149-153, Mar. 1979.
- [5] B. W. Bomar, "State-space structure for the realization of low roundoff noise digital filters," *Dissertation*, University of Tennessee, 1983.
- [6] D. Williamson, "Roundoff noise minimization and pole-zero sensitivity in fixed-point digital filter using residue feedback," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 34, pp. 1210-1220, Oct. 1986.
- [7] G. Li and M. Gevers, "Roundoff noise minimization using delta-operator realizations," *IEEE Trans. Signal Processing*, vol. 41, pp. 629-637, Feb. 1993.
- [8] T. Hinamoto, H. Ohnishi, and W.-S. Lu, "Roundoff noise minimization of state-space digital using separate and joint error feedback/coordinate transformation optimization," *IEEE Trans. Circuits Syst.*, *I*, 2003. (to appear)
- [9] T. Hinamoto, S. Karino, N. Kuroda and T. Kuma, "Error spectrum shaping in two-dimensional recursive digital filters," *IEEE Trans. Circuits Syst.*, vol. 46, pp. 1203-1215, Oct. 1999.
- [10] R. Fletcher, Practical Methods of Optimization, 2nd ed., Wiley, New York, 1987.
- [11] T. Kailath, Linear Systems, Prentice Hall, 1980.