

# Roundoff Noise Minimization for State-Estimate Feedback Digital Controllers Using Joint Optimization of Error Feedback and Realization

Takao Hinamoto, Keijiro Kawai, Masayoshi Nakamoto and Wu-Sheng Lu  
*Name-of-the-University-Company*  
*Country*

## 1. INTRODUCTION

Due to the finite precision nature of computer arithmetic, the output roundoff noise of a fixed-point IIR digital filter usually arises. This noise is critically dependent on the internal structure of an IIR digital filter [1],[2]. Error feedback (EF) is known as an effective technique for reducing the output roundoff noise in an IIR digital filter [3]-[5]. Williamson [6] has reduced the output roundoff noise more effectively by choosing the filter structure and applying EF to the filter. Lu and Hinamoto [7] have developed a jointly optimized technique of EF and realization to minimize the effects of roundoff noise at the filter output subject to  $l_2$ -norm dynamic-range scaling constraints. Li and Gevers [8] have analyzed the output roundoff noise of the closed-loop system with a state-estimate feedback controller, and presented an algorithm for realizing the state-estimate feedback controller with minimum output roundoff noise under  $l_2$ -norm dynamic-range scaling constraints. Hinamoto and Yamamoto [9] have proposed a method for applying EF to a given closed-loop system with a state-estimate feedback controller.

This paper investigates the problem of jointly optimizing EF and realization for the closed-loop system with a state-estimate feedback controller so as to minimize the output roundoff noise subject to  $l_2$ -norm dynamic-range scaling constraints. To this end, the problem at hand is converted into an unconstrained optimization problem by using linear-algebraic techniques, and then an iterative technique which relies on a quasi-Newton algorithm [10] is developed. With a closed-form formula for gradient evaluation and an efficient quasi-Newton solver, the unconstrained optimization problem can be solved efficiently. Our computer simulation results demonstrate the validity and effectiveness of the proposed technique.

Throughout the 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 and  $i$ th diagonal element of a square matrix  $A$  are denoted by  $\text{tr}[A]$  and  $(A)_{ii}$ , respectively.

## 2. ROUND OFF NOISE ANALYSIS

Consider a stable, controllable and observable linear discrete-time system described by

$$\begin{aligned}x(k+1) &= \mathbf{A}_o x(k) + \mathbf{b}_o u(k) \\ y(k) &= \mathbf{c}_o x(k)\end{aligned}\tag{1}$$

where  $x(k)$  is an  $n \times 1$  state-variable vector,  $u(k)$  is a scalar input,  $y(k)$  is a scalar output, and  $A_o$ ,  $b_o$  and  $c_o$  are  $n \times n$ ,  $n \times 1$  and  $1 \times n$  real constant matrices, respectively. The transfer function of the linear system in (1) is given by

$$H_o(z) = c_o(zI_n - A_o)^{-1}b_o. \tag{2}$$

If a regulator is designed by using the full-order state observer, we obtain a state-estimate feedback controller as

$$\begin{aligned} \tilde{x}(k+1) &= F_o\tilde{x}(k) + b_o u(k) + g_o y(k) \\ &= R_o\tilde{x}(k) + b_o r(k) + g_o y(k) \\ u(k) &= -k_o\tilde{x}(k) + r(k) \end{aligned} \tag{3}$$

where  $\tilde{x}(k)$  is an  $n \times 1$  state-variable vector in the full-order state observer,  $g_o$  is an  $n \times 1$  gain vector chosen so that all the eigenvalues of  $F_o = A_o - g_o c_o$  are inside the unit circle in the complex plane,  $k_o$  is a  $1 \times n$  state-feedback gain vector chosen so that each of the eigenvalues of  $A_o - b_o k_o$  is at a desirable location within the unit circle,  $r(k)$  is a scalar reference signal, and  $R_o = F_o - b_o k_o$ . The closed-loop control system consisting of the linear system in (1) and the state-estimate feedback controller in (3) is illustrated in Fig. 1.

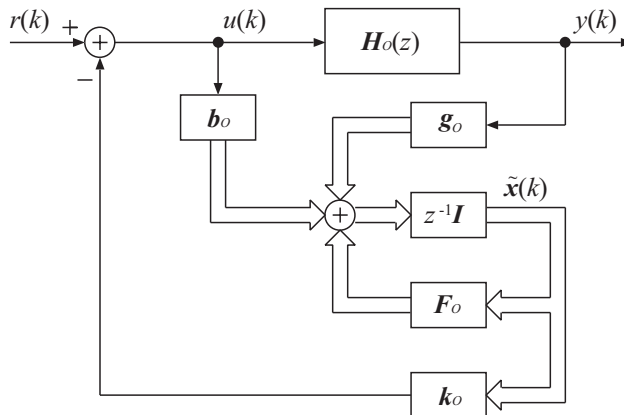


Fig. 1. The closed-loop control system with a state-estimate feedback controller.

When performing quantization before matrix-vector multiplication, we can express the finite-word-length (FWL) implementation of (3) with error feedback as

$$\begin{aligned} \hat{x}(k+1) &= R Q[\hat{x}(k)] + br(k) + gy(k) + De(k) \\ u(k) &= -k Q[\hat{x}(k)] + r(k) \end{aligned} \tag{4}$$

where

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

is an  $n \times 1$  roundoff error vector and  $D$  is an  $n \times n$  error feedback matrix. All coefficient matrices  $R$ ,  $b$ ,  $g$  and  $k$  are assumed to have an exact fractional  $B_c$  bit representation. The FWL

state-variable vector  $\hat{x}(k)$  and signal  $u(k)$  all have a  $B$  bit fractional representation, while the reference input  $r(k)$  is a  $(B - B_c)$  bit fraction. The vector quantizer  $Q[\cdot]$  in (4) rounds the  $B$  bit fraction  $\hat{x}(k)$  to  $(B - B_c)$  bits after completing the multiplications and additions, where the sign bit is not counted. It is assumed that the roundoff error vector  $e(k)$  can be modeled as a zero-mean noise process with covariance  $\sigma^2 I_n$  where

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

It is noted that if the  $i$ th element of the roundoff error vector  $e(k)$  is indicated by  $e_i(k)$  for  $i = 1, 2, \dots, n$  then the variable  $e_i(k)$  can be approximated by a white noise sequence uniformly distributed with the following probability density function:

$$p(e_i(k)) = \begin{cases} 2^{B-B_c} & \text{for } -\frac{1}{2}2^{-(B-B_c)} \leq e_i(k) \leq \frac{1}{2}2^{-(B-B_c)} \\ 0 & \text{otherwise} \end{cases}$$

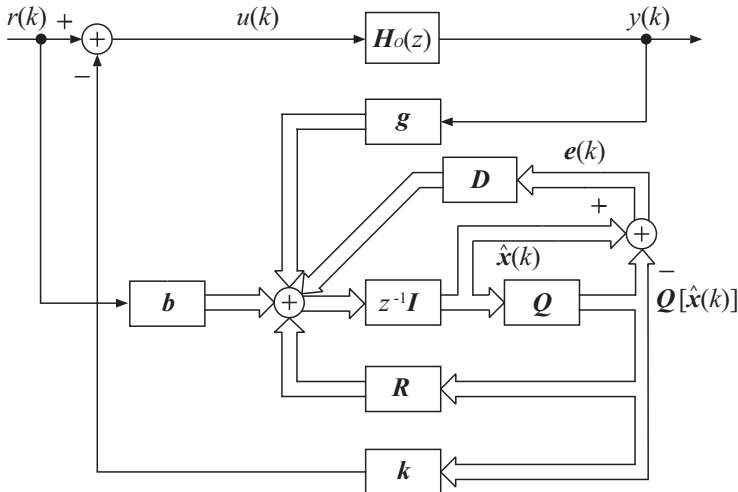


Fig. 2. A state-estimate feedback controller with error feedback.

The closed-loop system consisting of the linear system in (1) and the state-estimate feedback controller with error feedback in (4) is shown in Fig. 2, and is described by

$$\begin{bmatrix} x(k+1) \\ \hat{x}(k+1) \end{bmatrix} = \bar{A} \begin{bmatrix} x(k) \\ \hat{x}(k) \end{bmatrix} + \bar{b}r(k) + \bar{B}e(k) \tag{5}$$

$$y(k) = \bar{c} \begin{bmatrix} x(k) \\ \hat{x}(k) \end{bmatrix}$$

where

$$\bar{A} = \begin{bmatrix} A_o & -b_o k \\ g c_o & R \end{bmatrix}, \quad \bar{b} = \begin{bmatrix} b_o \\ b \end{bmatrix}$$

$$\bar{B} = \begin{bmatrix} b_o k \\ D - R \end{bmatrix}, \quad \bar{c} = [c_o \ 0].$$

From (5), the transfer function from the roundoff error vector  $e(k)$  to the output  $y(k)$  is given by

$$G_D(z) = \bar{c} (zI_{2n} - \bar{A})^{-1} \bar{B}. \quad (6)$$

The output noise gain  $J(D) = \sigma_{out}^2 / \sigma^2$  is then computed as

$$J(D) = \text{tr}[W_D] \quad (7)$$

with

$$W_D = \frac{1}{2\pi j} \oint_{|z|=1} G_D^*(z) G_D(z) \frac{dz}{z} \quad (8)$$

where  $\sigma_{out}^2$  stands for the noise variance at the output. For tractability, we evaluate  $J(D)$  in (7) by replacing  $R, b, g$  and  $k$  by  $R_o, b_o, g_o$  and  $k_o$ , respectively. Defining

$$S = \begin{bmatrix} I_n & \mathbf{0} \\ I_n & -I_n \end{bmatrix}, \quad (9)$$

the transfer function in (6) can be expressed as

$$\begin{aligned} G_D(z) &= \bar{c} S (zI_{2n} - S^{-1} \bar{A} S)^{-1} S^{-1} \bar{B} \\ &= \bar{c} (zI_{2n} - \Phi)^{-1} \begin{bmatrix} b_o k_o \\ F_o - D \end{bmatrix} \\ &= c_o (zI_n - A_o + b_o k_o)^{-1} b_o k_o (zI_n - F_o)^{-1} \\ &\quad \cdot (zI_n - D) \\ &= \bar{c} (zI_{2n} - \Phi)^{-1} U (zI_n - D) \end{aligned} \quad (10)$$

where

$$\Phi = \begin{bmatrix} A_o - b_o k_o & b_o k_o \\ \mathbf{0} & F_o \end{bmatrix}$$

$$U = \begin{bmatrix} \mathbf{0} \\ I_n \end{bmatrix}.$$

It is noted that the stability of the closed-loop control system is determined by the eigenvalues of matrix  $\bar{A}$  in (5), or equivalently, those of matrix  $\Phi$  in (10). This means that neither of the roundoff error vector  $e(k)$  and the error-feedback matrix  $D$  affects the stability.

Substituting (10) into matrix  $W_D$  in (8) gives

$$\begin{aligned} W_D &= (b_o k_o)^T W_1 b_o k_o + (b_o k_o)^T W_2 (F_o - D) \\ &\quad + (F_o - D)^T W_3 b_o k_o \\ &\quad + (F_o - D)^T W_4 (F_o - D) \end{aligned} \quad (11)$$

where

$$W = \Phi^T W \Phi + \bar{c}^T \bar{c}$$

$$W = \begin{bmatrix} W_1 & W_2 \\ W_3 & W_4 \end{bmatrix}.$$

Since  $W$  is positive semidefinite, it can be shown that there exists an  $n \times n$  matrix  $P$  such that  $W_3 = W_4 P$ . In addition, (11) can be written by virtue of  $W_2 = W_3^T$  as

$$W_D = (F_0 + P b_0 k_0 - D)^T W_4 (F_0 + P b_0 k_0 - D) + (b_0 k_0)^T (W_1 - P^T W_4 P) b_0 k_0. \tag{12}$$

Alternatively, applying  $z$ -transform to the first equation in (5) under the assumption that  $e(k) = \mathbf{0}$ , we obtain

$$\begin{bmatrix} X(z) \\ \hat{X}(z) \end{bmatrix} = (zI - \bar{A})^{-1} \bar{b} R(z) \tag{13}$$

where  $X(z)$ ,  $\hat{X}(z)$  and  $R(z)$  represent the  $z$ -transforms of  $x(k)$ ,  $\hat{x}(k)$  and  $r(k)$ , respectively. Replacing  $R$ ,  $b$ ,  $k$  and  $g$  by  $R_o$ ,  $b_o$ ,  $k_o$  and  $g_o$ , respectively, and then using

$$S^{-1} \begin{bmatrix} X(z) \\ \hat{X}(z) \end{bmatrix} = (zI_{2n} - S^{-1} \bar{A} S)^{-1} S^{-1} \bar{b}$$

yields

$$\hat{X}(z) = X(z) = F(z) R(z) \tag{14}$$

where

$$F(z) = [zI_n - (A_o - b_o k_o)]^{-1} b_o.$$

The controllability Gramian  $K$  defined by

$$K = \frac{1}{2\pi j} \oint_{|z|=1} F(z) F^*(z) \frac{dz}{z} \tag{15}$$

can be obtained by solving the following Lyapunov equation:

$$K = (A_o - b_o k_o) K (A_o - b_o k_o)^T + b_o b_o^T. \tag{16}$$

### 3. ROUND OFF NOISE MINIMIZATION

Consider the system in (4) with  $D = \mathbf{0}$  and denote it by  $(R, b, g, k)_n$ . By applying a coordinate transformation  $\tilde{x}'(k) = T^{-1} \hat{x}(k)$  to the above system  $(R, b, g, k)_n$ , we obtain a new realization characterized by  $(\tilde{R}, \tilde{b}, \tilde{g}, \tilde{k})_n$  where

$$\begin{aligned} \tilde{R} &= T^{-1} R T, & \tilde{b} &= T^{-1} b \\ \tilde{g} &= T^{-1} g, & \tilde{k} &= k T. \end{aligned} \tag{17}$$

For the system described by (17), the counterparts of  $W_i$  for  $i = 1, 2, 3, 4$  are given by

$$\tilde{W}_i = T^T W_i T \tag{18}$$

and the corresponding output noise gain is given by

$$J(\mathbf{D}, \mathbf{T}) = \text{tr}[\tilde{\mathbf{W}}_D] \tag{19}$$

where  $\tilde{\mathbf{W}}_D$  can be obtained referring to (11) as

$$\begin{aligned} \tilde{\mathbf{W}}_D &= \left[ \mathbf{T}^{-1}(\mathbf{F}_0 + \mathbf{P}\mathbf{b}_0\mathbf{k}_0)\mathbf{T} - \mathbf{D} \right]^T \\ &\quad \cdot \mathbf{T}^T \mathbf{W}_4 \mathbf{T} \left[ \mathbf{T}^{-1}(\mathbf{F}_0 + \mathbf{P}\mathbf{b}_0\mathbf{k}_0)\mathbf{T} - \mathbf{D} \right] \\ &\quad + \mathbf{T}^T (\mathbf{b}_0\mathbf{k}_0)^T (\mathbf{W}_1 - \mathbf{P}^T \mathbf{W}_4 \mathbf{P}) \mathbf{b}_0\mathbf{k}_0 \mathbf{T}. \end{aligned}$$

In addition, (15) can be written as

$$\begin{aligned} \tilde{\mathbf{K}} &= \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{T}^{-1} \mathbf{F}(z) \mathbf{F}^*(z) \mathbf{T}^{-T} \frac{dz}{z} \\ &= \mathbf{T}^{-1} \mathbf{K} \mathbf{T}^{-T}. \end{aligned} \tag{20}$$

As a result, the output roundoff noise minimization problem amounts to obtaining matrices  $\mathbf{D}$  and  $\mathbf{T}$  which jointly minimize  $J(\mathbf{D}, \mathbf{T})$  in (19) subject to the  $l_2$ -norm dynamic-range scaling constraints specified by

$$(\tilde{\mathbf{K}})_{ii} = (\mathbf{T}^{-1} \mathbf{K} \mathbf{T}^{-T})_{ii} = 1, \quad i = 1, 2, \dots, n. \tag{21}$$

To deal with (21), we define

$$\hat{\mathbf{T}} = \mathbf{T}^T \mathbf{K}^{-\frac{1}{2}}. \tag{22}$$

Then the  $l_2$ -norm dynamic-range scaling constraints in (21) can be written as

$$(\hat{\mathbf{T}}^{-T} \hat{\mathbf{T}}^{-1})_{ii} = 1, \quad i = 1, 2, \dots, n. \tag{23}$$

These constraints are always satisfied if  $\hat{\mathbf{T}}^{-1}$  assumes the form

$$\hat{\mathbf{T}}^{-1} = \left[ \frac{\mathbf{t}_1}{\|\mathbf{t}_1\|}, \frac{\mathbf{t}_2}{\|\mathbf{t}_2\|}, \dots, \frac{\mathbf{t}_n}{\|\mathbf{t}_n\|} \right]. \tag{24}$$

Substituting (22) into (19), we obtain

$$\begin{aligned} J(\mathbf{D}, \hat{\mathbf{T}}) &= \text{tr} \left[ \hat{\mathbf{T}} (\hat{\mathbf{A}} - \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T})^T \hat{\mathbf{W}}_4 \right. \\ &\quad \left. \cdot (\hat{\mathbf{A}} - \hat{\mathbf{T}}^T \mathbf{D} \hat{\mathbf{T}}^{-T}) \hat{\mathbf{T}}^T + \hat{\mathbf{T}} \hat{\mathbf{C}} \hat{\mathbf{T}}^T \right] \end{aligned} \tag{25}$$

where

$$\begin{aligned} \hat{\mathbf{A}} &= \mathbf{K}^{-\frac{1}{2}} (\mathbf{F}_0 + \mathbf{P}\mathbf{b}_0\mathbf{k}_0) \mathbf{K}^{\frac{1}{2}}, \quad \hat{\mathbf{W}}_4 = \mathbf{K}^{\frac{1}{2}} \mathbf{W}_4 \mathbf{K}^{\frac{1}{2}} \\ \hat{\mathbf{C}} &= \mathbf{K}^{\frac{1}{2}} (\mathbf{b}_0\mathbf{k}_0)^T (\mathbf{W}_1 - \mathbf{P}^T \mathbf{W}_4 \mathbf{P}) \mathbf{b}_0\mathbf{k}_0 \mathbf{K}^{\frac{1}{2}}. \end{aligned}$$

From the foregoing arguments, the problem of obtaining matrices  $\mathbf{D}$  and  $\mathbf{T}$  that minimize (19) subject to the scaling constraints in (21) is now converted into an unconstrained optimization problem of obtaining  $\mathbf{D}$  and  $\hat{\mathbf{T}}$  that jointly minimize  $J(\mathbf{D}, \hat{\mathbf{T}})$  in (25).

Let  $x$  be the column vector that collects the variables in matrix  $D$  and matrix  $[t_1, t_2, \dots, t_n]$ . Then  $J(D, \hat{T})$  is a function of  $x$ , denoted by  $J(x)$ . The proposed algorithm starts with an initial point  $x_0$  obtained from an initial assignment  $D = \hat{T} = I_n$ . In the  $k$ th iteration, a quasi-Newton algorithm updates the most recent point  $x_k$  to point  $x_{k+1}$  as [10]

$$x_{k+1} = x_k + \alpha_k d_k \tag{26}$$

where

$$\begin{aligned} d_k &= -S_k \nabla J(x_k) \\ \alpha_k &= \arg \left[ \min_{\alpha} J(x_k + \alpha d_k) \right] \\ S_{k+1} &= S_k + \left( 1 + \frac{\gamma_k^T S_k \gamma_k}{\gamma_k^T \delta_k} \right) \frac{\delta_k \delta_k^T}{\gamma_k^T \delta_k} - \frac{\delta_k \gamma_k^T S_k + S_k \gamma_k \delta_k^T}{\gamma_k^T \delta_k} \\ S_0 &= I, \quad \delta_k = x_{k+1} - x_k, \quad \gamma_k = \nabla J(x_{k+1}) - \nabla J(x_k). \end{aligned}$$

Here,  $\nabla J(x)$  is the gradient of  $J(x)$  with respect to  $x$ , and  $S_k$  is a positive-definite approximation of the inverse Hessian matrix of  $J(x_k)$ . This iteration process continues until

$$|J(x_{k+1}) - J(x_k)| < \varepsilon \tag{27}$$

is satisfied where  $\varepsilon > 0$  is a prescribed tolerance.

In what follows, we derive closed-form expressions of  $\nabla J(x)$  for the cases where  $D$  assumes the form of a general, diagonal, or scalar matrix.

1) *Case 1: D Is a General Matrix:* From (25), the optimal choice of  $D$  is given by

$$D = \hat{T}^{-T} \hat{A} \hat{T}^T, \tag{28}$$

which leads to

$$J(\hat{T}^{-T} \hat{A} \hat{T}^T, \hat{T}) = \text{tr} [\hat{T} \hat{C} \hat{T}^T]. \tag{29}$$

In this case, the number of elements in vector  $x$  consisting of  $\hat{T}$  is equal to  $n^2$  and the gradient of  $J(x)$  is found to be

$$\begin{aligned} \frac{\partial J(x)}{\partial t_{ij}} &= \lim_{\Delta \rightarrow 0} \frac{J(\hat{T}_{ij}) - J(\hat{T})}{\Delta} \\ &= 2e_j^T \hat{T} \hat{C} \hat{T}^T \hat{T} g_{ij}, \quad i, j = 1, 2, \dots, n \end{aligned} \tag{30}$$

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

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

and  $g_{ij}$  is computed using

$$g_{ij} = \partial \left\{ \frac{t_j}{\|t_j\|} \right\} / \partial t_{ij} = \frac{1}{\|t_j\|^3} (t_{ij} t_j - \|t_j\|^2 e_i).$$

2) *Case 2: D Is a Diagonal Matrix:* Here, matrix  $D$  assumes the form

$$D = \text{diag}\{d_1, d_2, \dots, d_n\}. \tag{31}$$

In this case, (25) becomes

$$J(D, \hat{T}) = \text{tr} [\hat{T} M_d \hat{T}^T] \tag{32}$$

where

$$M_d = \hat{C} + \hat{A}^T \hat{W}_4 \hat{A} + \hat{W}_4 \hat{T}^T D^2 \hat{T}^{-T} - \hat{A}^T \hat{W}_4 \hat{T}^T D \hat{T}^{-T} - \hat{W}_4 \hat{A} \hat{T}^T D \hat{T}^{-T}.$$

It follows that

$$\begin{aligned} \frac{\partial J(x)}{\partial t_{ij}} &= 2e_j^T \hat{T} M_d \hat{T}^T \hat{T} g_{ij}, \quad i, j = 1, 2, \dots, n \\ \frac{\partial J(x)}{\partial d_i} &= 2e_i^T (D \hat{T} - \hat{T} \hat{A}^T) \hat{W}_4 \hat{T}^T e_i, \quad i = 1, 2, \dots, n. \end{aligned} \tag{33}$$

3) *Case 3: D Is a Scalar Matrix:* It is assumed here that  $D = \alpha I_n$  with a scalar  $\alpha$ . The gradient of  $J(x)$  can then be calculated as

$$\begin{aligned} \frac{\partial J(x)}{\partial t_{ij}} &= 2e_j^T \hat{T} M_s \hat{T}^T \hat{T} g_{ij}, \quad i, j = 1, 2, \dots, n \\ \frac{\partial J(x)}{\partial \alpha} &= \text{tr} [\hat{T} (2\alpha \hat{W}_4 - \hat{A}^T \hat{W}_4 - \hat{W}_4 \hat{A}) \hat{T}^T] \end{aligned} \tag{34}$$

where

$$M_s = (\hat{A} - \alpha I_n)^T \hat{W}_4 (\hat{A} - \alpha I_n) + \hat{C}.$$

### 4. A NUMERICAL EXAMPLE

In this section we illustrate the proposed method by considering a linear discrete-time system specified by

$$\begin{aligned} A_o &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0.339377 & -1.152652 & 1.520167 \end{bmatrix}, \quad b_o = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ c_o &= [ 0.093253 \quad 0.128620 \quad 0.314713 ]. \end{aligned}$$

Suppose that the poles of the observer and regulator in the system are required to be located at  $z = 0.1532, 0.2861, 0.1137$ , and  $z = 0.5067, 0.6023, 0.4331$ , respectively. This can be achieved by choosing

$$\begin{aligned} k_o &= [ 0.471552 \quad -0.367158 \quad 3.062267 ] \\ g_o &= [ -0.006436 \quad 3.683651 \quad 5.083920 ]^T. \end{aligned}$$

Performing the  $l_2$ -norm dynamic-range scaling to the state-estimate feedback controller, we obtain  $J(0) = 686.4121$  in (7) where  $D = 0$ . Next, the controller is transformed into the optimal realization that minimizes  $J(0)$  in (7) under the  $l_2$ -norm dynamic-range scaling constraints. This leads to  $J_{min}(0) = 28.6187$ . Finally, EF and state-variable coordinate transformation are applied to the above optimal realization so as to jointly minimize the output roundoff noise. The profiles of  $J(x)$  during the first 20 iteration for the cases of  $D$  being a general, diagonal, and scalar matrix are depicted in Fig. 3.



1) *Case 1: D Is a General Matrix:* The quasi-Newton algorithm was applied to minimize (25). It took the algorithm 20 iterations to converge to the solution

$$D = \begin{bmatrix} 0.211191 & -3.078211 & -3.344596 \\ -1.321589 & 1.897308 & 3.243515 \\ 1.917916 & -1.890027 & -3.807473 \end{bmatrix}$$

$$T = \begin{bmatrix} -11.039974 & -43.683697 & -30.131793 \\ -3.231505 & 8.919473 & 9.118205 \\ 2.620911 & 6.462685 & 7.032260 \end{bmatrix}$$

and the minimized noise gain was found to be  $J(D, \hat{T}) = 4.8823$ . Next, the above optimal EF matrix  $D$  was rounded to a power-of-two representation with 3 bits after the binary point, which resulted in

$$D_{3bit} = \begin{bmatrix} 0.250 & -3.125 & -3.375 \\ -1.375 & 1.875 & 3.250 \\ 1.875 & -1.875 & -3.750 \end{bmatrix}$$

and a noise gain  $J(D_{3bit}, \hat{T}) = 23.4873$ . Furthermore, when the optimal EF matrix  $D$  was rounded to the integer representation

$$D_{int} = \begin{bmatrix} 0 & -3 & -3 \\ -1 & 2 & 3 \\ 2 & -2 & -4 \end{bmatrix},$$

the noise gain was found to be  $J(D_{int}, \hat{T}) = 293.0187$ .

2) *Case 2: D Is a Diagonal Matrix:* Again, the quasi-Newton algorithm was applied to minimize  $J(D, \hat{T})$  in (25) for a diagonal EF matrix  $D$ . It took the algorithm 20 iterations to converge to the solution

$$D = \text{diag}\{0.050638, -0.608845, -0.951572\}$$

$$T = \begin{bmatrix} 3.588878 & 0.735966 & 0.010417 \\ -2.457241 & 0.728171 & 0.556762 \\ 1.514232 & -2.058856 & 0.142204 \end{bmatrix}$$

and the minimized noise gain was found to be  $J(D, \hat{T}) = 12.7097$ . Next, the above optimal diagonal EF matrix  $D$  was rounded to a power-of-two representation with 3 bits after the binary point to yield  $D_{3bit} = \text{diag}\{0.000, -0.625, -1.000\}$ , which leads to a noise gain  $J(D_{3bit}, \hat{T}) = 12.7722$ . Furthermore, when the optimized diagonal EF matrix  $D$  was rounded to the integer representation  $D_{int} = \text{diag}\{0, -1, -1\}$ , the noise gain was found to be  $J(D_{int}, \hat{T}) = 13.7535$ .

3) *Case 3: D Is a Scalar Matrix:* In this case, the quasi-Newton algorithm was applied to minimize (25) for  $D = \alpha I_3$  with a scalar  $\alpha$ . The algorithm converges after 20 iterations to converge to the solution

$$D = -0.779678 I_3$$

$$T = \begin{bmatrix} 3.252790 & -0.081745 & -0.198376 \\ -1.717225 & 1.220068 & -0.792487 \\ 0.546599 & -0.854316 & 2.295944 \end{bmatrix}$$

and the minimized noise gain was found to be  $J(D, \hat{T}) = 16.2006$ . Next, the EF matrix  $D = \alpha I_3$  was rounded to a power-of-two representation with 3 bits after the binary point as well as

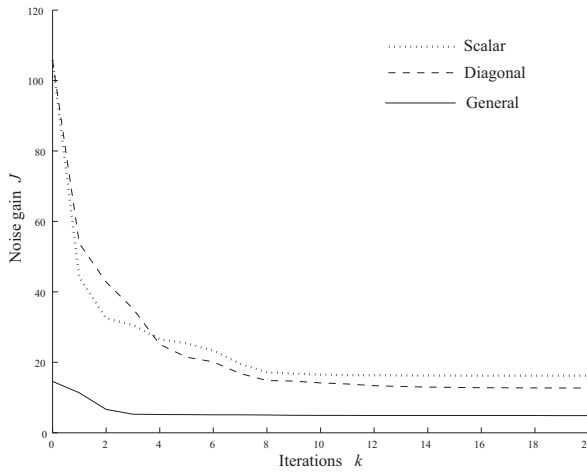


Fig. 3. Profiles of iterative noise gain minimization.

an integer representation. It was found that these representations were given by  $D_{3bit} = \text{diag}\{0.750, 0.750, 0.750\}$  and  $D_{int} = \text{diag}\{1, 1, 1\}$ , respectively. The corresponding noise gains were obtained as  $J(D_{3bit}, \hat{T}) = 16.2370$  and  $J(D_{int}, \hat{T}) = 18.2063$ , respectively.

The above simulation results in terms of noise gain  $J(D, \hat{T})$  in (25) are summarized in Table 1. For comparison purpose, their counterparts obtained using the method in [9] are also included in the table, where the minimization of the roundoff noise was carried out using EF and state-variable coordinate transformation, but in a separate manner. From the table, it is observed that the proposed joint optimization offers improved reduction in roundoff noise gain for the cases of a scalar EF matrix and a diagonal EF matrix when compared with those obtained by using *separate* optimization. However, in the case of a general EF matrix, the optimal solution with infinite precision appears to be quite sensitive to the parameter perturbations.

Error-Feedback Scheme	Accuracy of $D$		
	Infinite Precision	3 Bit Quantization	Integer Quantization
$D = 0$	28.6187		
Separate	28.6187		
Scalar Separate [9]	20.1235	20.1810	26.0527
Scalar Joint	16.2006	16.2370	18.2063
Diagonal Separate [9]	16.4104	16.4547	17.4039
Diagonal Joint	12.7097	12.7722	13.7535
General Separate [9]	11.6352	11.7054	16.5814
General Joint	4.8823	23.4873	293.0187

Table 1. Noise gain  $J(D, \hat{T})$  for different EF schemes.

More reduction of the noise gain might be possible by re-designing the coordinate transformation matrix  $T$  for the optimally quantized  $D$ .

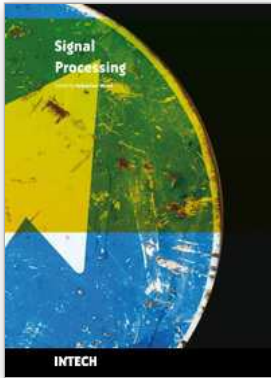
## 5. CONCLUSION

The joint optimization problem of EF and realization to minimize the effects of roundoff noise of the closed-loop system with a state-estimate feedback controller subject to  $l_2$ -norm dynamic-range scaling constraints has been investigated. The problem at hand has been converted into an unconstrained optimization problem by using linear algebraic techniques. An efficient quasi-Newton algorithm has been employed to solve the unconstrained optimization problem. The proposed technique has been applied to the cases where EF matrix is a general, diagonal, or scalar matrix. The effectiveness for the cases of a scalar EF matrix and a diagonal EF matrix compared with the existing method [9] has been illustrated by a numerical example.

## 6. References

- 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.
- 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.
- 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.
- 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.
- P. P. Vaidyanathan, "On error-spectrum shaping in state-space digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-32, pp. 88-92, Jan. 1985.
- 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.
- W.-S. Lu and T. Hinamoto, "Jointly optimized error-feedback and realization for roundoff noise minimization in state-space digital filters," *IEEE Trans. Signal Processing*, vol. 53, pp. 2135-2145, June 2005.
- G. Li and M. Gevers, "Optimal finite precision implementation of a state-estimate feedback controller," *IEEE Trans. Circuits Syst.*, vol. CAS-37, pp. 1487-1498, Dec. 1990.
- T. Hinamoto and S. Yamamoto, "Error spectrum shaping in closed-loop systems with state-estimate feedback controller," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS'02)*, May 2002, vol. 1, pp. 289-292.
- R. Fletcher, *Practical Methods of Optimization*, 2nd ed. New York: Wiley, 1987.





## Signal Processing

Edited by Sebastian Miron

ISBN 978-953-7619-91-6

Hard cover, 528 pages

**Publisher** InTech

**Published online** 01, March, 2010

**Published in print edition** March, 2010

This book intends to provide highlights of the current research in signal processing area and to offer a snapshot of the recent advances in this field. This work is mainly destined to researchers in the signal processing related areas but it is also accessible to anyone with a scientific background desiring to have an up-to-date overview of this domain. The twenty-five chapters present methodological advances and recent applications of signal processing algorithms in various domains as telecommunications, array processing, biology, cryptography, image and speech processing. The methodologies illustrated in this book, such as sparse signal recovery, are hot topics in the signal processing community at this moment. The editor would like to thank all the authors for their excellent contributions in different areas of signal processing and hopes that this book will be of valuable help to the readers.

### How to reference

In order to correctly reference this scholarly work, feel free to copy and paste the following:

Takao Hinamoto, Keijiro Kawai, Masayoshi Nakamoto and Wu-Sheng Lu (2010). Roundoff Noise Minimization for State-Estimate Feedback Digital Controllers Using Joint Optimization of Error Feedback and Realization, Signal Processing, Sebastian Miron (Ed.), ISBN: 978-953-7619-91-6, InTech, Available from: <http://www.intechopen.com/books/signal-processing/roundoff-noise-minimization-for-state-estimate-feedback-digital-controllers-using-joint-optimization>

# INTECH

open science | open minds

### InTech Europe

University Campus STeP Ri  
Slavka Krautzeka 83/A  
51000 Rijeka, Croatia  
Phone: +385 (51) 770 447  
Fax: +385 (51) 686 166  
[www.intechopen.com](http://www.intechopen.com)

### InTech China

Unit 405, Office Block, Hotel Equatorial Shanghai  
No.65, Yan An Road (West), Shanghai, 200040, China  
中国上海市延安西路65号上海国际贵都大饭店办公楼405单元  
Phone: +86-21-62489820  
Fax: +86-21-62489821

© 2010 The Author(s). Licensee IntechOpen. This chapter is distributed under the terms of the [Creative Commons Attribution-NonCommercial-ShareAlike-3.0 License](#), which permits use, distribution and reproduction for non-commercial purposes, provided the original is properly cited and derivative works building on this content are distributed under the same license.