

# OPTIMAL SPARSE CONTROLLER STRUCTURE WITH MINIMUM ROUND-OFF NOISE GAIN

Jinxin Hao, Teck Chew Wee, Lucas S. Karatzas and Yew Fai Lee  
*School of Engineering, Temasek Polytechnic, 529757, Singapore*

**Keywords:** Roundoff noise gain, Sparse controller structure, Optimization, Direct-form II transposed (DFIIT) structure.

**Abstract:** This paper investigates the roundoff noise effect in the digital controller on the closed-loop output for a discrete-time feedback control system. Based on a polynomial parametrization approach, a sparse controller structure is derived. The performance of the proposed structure is analyzed by deriving the corresponding expression of closed-loop roundoff noise gain and the problem of finding optimized sparse structures is solved. A numerical example is presented to illustrate the design procedure and the performance of the proposed structure compared with those of some existing well-known structures.

## 1 INTRODUCTION

Finite word length (FWL) effects have been a well studied field in the design of digital filters for more than three decades (Mullis and Roberts, 1976), (Hwang, 1977), (Roberts and Mullis, 1987), (Gevers and Li, 1993). However, they have received less attention in the area of digital control. Nowadays, many researchers have recognized the importance of the numerical problems caused by FWL effects in digital controller implementation. The optimal FWL controller structure design (Fialho and Georgiou, 1994), (Li, 1998), (Wu et al., 2001), (Yu and Ko, 2003) has been considered as one of the most effective methods to minimize the effects of FWL errors on the performance of closed-loop control systems. The basic idea behind this approach is that for a given digital controller, there exist different structures which have different numerical properties, and the optimal structure problem is to identify those structures that optimize a certain FWL performance criterion.

Generally speaking, there are two types of FWL errors in the digital controller. The first one is the perturbation of the controller parameters implemented with FWL, and the second one is the rounding errors that occur in arithmetic operations, which are usually measured with the so-called roundoff noise gain. The effects of roundoff noise have been well studied in digital signal processing, particularly in digital filter implementation (Wong and Ng, 2000), (Wong and Ng, 2001). However, it was not un-

til the late 1980s that the problem of optimal controller realizations minimizing the roundoff noise gain was addressed. The roundoff noise gain was derived for a control system with a state-estimate feedback controller and the corresponding optimal realization problem was solved in (Li and Gevers, 1990), while the roundoff error effect on the linear quadratic regulation (LQG) performance was investigated in (Williamson and Kadiman, 1989) and the optimal solution was obtained by Liu *et al* (Liu et al., 1992). The problem of finding the optimum roundoff noise structures of digital controllers in a sampled-data system has been investigated in (Li et al., 2002).

It has been noted that the optimal controller realizations obtained with the above design methods are usually fully parametrized, which increase the complexity for real-time implementations. From a practical point of view, it is desired that the actually implemented controller have a nice performance against the FWL effects as well as a sparse structure that possesses many trivial parameters<sup>1</sup> which produce no FWL errors. As far as we know, a few results have been published on the sparseness issue for the controller structure design (Li, 1998), (Wu et al., 2003), however, it is noted that in these approaches, sophisticated numerical algorithms were utilized and the positions of trivial parameters were not predictable. In (Hao et al., 2006), we proposed two sparse structures

<sup>1</sup>By trivial parameters we mean those that are 0 and  $\pm 1$ , other parameters are, therefore, referred to as nontrivial parameters.

for digital controllers, which have some degrees of freedom that can be used to enhance the closed-loop stability robustness against the FWL effects.

In this paper, a new sparse controller structure is derived by adopting the polynomial parametrization approach in (Hao et al., 2006) and using the  $l_2$ -scaling scheme. This structure can be considered as a  $l_2$ -scaled generalized DFII (direct-form II transposed) structure. The expression of the roundoff noise gain is derived for a closed-loop feedback control system, in which the digital controller is implemented with the proposed structure. The problem of finding optimized sparse structures is solved by minimizing the corresponding closed-loop roundoff noise gain. A numerical example is given to illustrate the design procedure, which shows that the proposed structure beats the traditional DFII structures greatly in terms of roundoff noise performance, and furthermore, outperforms the fully parametrized optimal realization (Li et al., 2002) in terms of both roundoff noise gain and computation efficiency.

## 2 A SPARSE CONTROLLER STRUCTURE

Consider a discrete-time feedback control system depicted in Fig. 1, where  $P_d(z)$  is the discrete-time plant and  $C_d(z)$  is a well-designed digital controller. The controller can be represented by its transfer function which is parametrized with  $\{\xi_k, \zeta_k\}$  in the shift operator  $z$ :

$$C_d(z) = \frac{\sum_{k=0}^K \zeta_k z^{K-k}}{z^K + \sum_{k=1}^K \xi_k z^{K-k}}. \quad (1)$$

This controller can be implemented with many different structures, such as the direct forms or the following state-space equations:

$$\begin{cases} x(n+1) &= Ax(n) + Bu(n) \\ y(n) &= Cx(n) + du(n) \end{cases} \quad (2)$$

where  $x(n) \in \mathcal{R}^{K \times 1}$  is the state variable vector and  $u(n)$ ,  $y(n)$  are the input and output of the controller  $C_d(z)$ , respectively, while  $r(n)$  is the input signal of the closed-loop system.  $R \triangleq (A, B, C, d)$  is called a realization of  $C_d(z)$  with  $A \in \mathcal{R}^{K \times K}$ ,  $B \in \mathcal{R}^{K \times 1}$ ,  $C \in \mathcal{R}^{1 \times K}$  and  $d \in \mathcal{R}$ , satisfying

$$C_d(z) = d + C(zI - A)^{-1}B.$$

Denote  $S_C$  as the set of all the realizations:  $S_C \triangleq \{(A, B, C, d) : C_d(z) = d + C(zI - A)^{-1}B\}$ . Let  $R_0 \triangleq (A_0, B_0, C_0, d) \in S_C$  be an initial realization. It can be shown that  $S_C$  is characterized by

$$A = T^{-1}A_0T, \quad B = T^{-1}B_0, \quad C = C_0T \quad (3)$$

where  $T \in \mathcal{R}^{K \times K}$  is any nonsingular matrix. Such a matrix  $T$  is usually called a similarity transformation. Once an initial realization  $R_0$  is given, different controller realizations correspond to different similarity transformations  $T$ .

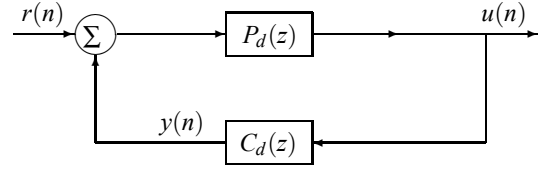


Figure 1: A discrete-time feedback control system.

### 2.1 A Generalized DFII Structure

Based on the approach in (Hao et al., 2006), we define

$$\rho_k(z) \triangleq \frac{z - \gamma_k}{\Delta_k}, \quad k = 1, 2, \dots, K, \quad (4)$$

where  $\{\gamma_k\}$  and  $\{\Delta_k > 0\}$  are two sets of constants to be discussed later. Let

$$\begin{aligned} p_k(z) &\triangleq \prod_{m=k+1}^K \rho_m(z), \quad \forall k \in \{0, 1, \dots, K-1\}, \\ p_K(z) &\triangleq 1. \end{aligned} \quad (5)$$

It can be shown that (1) can be rewritten as

$$C_d(z) = \frac{\beta_0 p_0(z) + \beta_1 p_1(z) + \dots + \beta_K p_K(z)}{p_0(z) + \alpha_1 p_1(z) + \dots + \alpha_K p_K(z)}, \quad (6)$$

where

$$\begin{aligned} \bar{\alpha} &\triangleq [1 \quad \alpha_1 \quad \dots \quad \alpha_K]^T \\ &= \kappa \bar{T}_p^{-T} [1 \quad \xi_1 \quad \dots \quad \xi_K]^T \\ \bar{\beta} &\triangleq [\beta_0 \quad \beta_1 \quad \dots \quad \beta_K]^T \\ &= \kappa \bar{T}_p^{-T} [\zeta_0 \quad \zeta_1 \quad \dots \quad \zeta_K]^T \end{aligned}$$

with  $\kappa = \prod_{k=1}^K \Delta_k^{-1}$  such that  $\bar{\alpha}(1) = 1$  and  $\bar{T}_p$  an upper triangular matrix whose  $k$ th row is formed with the coefficients of  $p_{k-1}(z)$  defined above. Equation (6) implies that the controller transfer function  $C_d(z)$  is reparametrized with  $\{\alpha_k\}$  and  $\{\beta_k\}$  in the new set of polynomial operators  $\{p_k(z)\}$ .

It follows from (5) and (6) that the output of the controller can be computed with the following equations

$$\begin{aligned} y(n) &= \beta_0 u(n) + w_1(n) \\ w_k(n) &= \rho_k^{-1} [\beta_k u(n) - \alpha_k y(n) + w_{k+1}(n)] \\ w_K(n) &= \rho_K^{-1} [\beta_K u(n) - \alpha_K y(n)] \end{aligned} \quad (7)$$

where  $w_k(n)$  is the output of  $\rho_k^{-1}(z)$  and can be computed with the structure depicted in Fig. 2. Fig. 3 shows the corresponding structure to (7). For convenience, a structure defined by Figs 2 and 3 is called a generalized DFII structure, denoted as  $\rho$ DFII. This structure possesses  $\{\alpha_k, \beta_k, \Delta_k\}$  and a set of free parameters  $\{\gamma_k\}$ . For a given digital controller  $C_d(z)$ , there exists a class of such structures, depending on the space within which  $\{\gamma_k\}$  take values. Clearly, when  $\gamma_k = 0, \Delta_k = 1, \forall k$ , Fig. 3 is the conventional direct-form II transposed (DFII) structure.

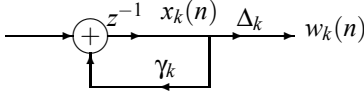


Figure 2: A realization of  $\rho_k^{-1}(z)$  defined in (4).

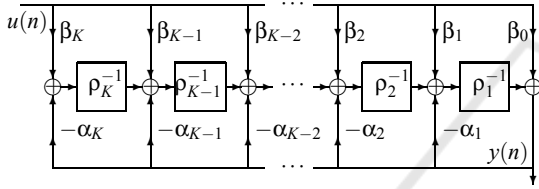


Figure 3: Block diagram of the  $\rho$ DFII structure.

With  $\{x_k(n)\}$  indicated in Fig. 2 as the state variables and  $x(n)$  denoting the state vector, one can find the equivalent state-space realization, denoted as  $(A_p, B_p, C_p, \beta_0)$ , of the proposed  $\rho$ DFII structure:

$$C_d(z) = \beta_0 + C_p(zI - A_p)^{-1}B_p \quad (8)$$

with  $B_p = \bar{V}_\beta - \beta_0 \bar{V}_\alpha$ , where  $\bar{V}_x \triangleq [x_1 \cdots x_k \cdots x_K]^T$  for  $x = \alpha, \beta$ ,  $C_p = [\Delta_1 \ 0 \ \cdots \ 0 \ 0]$ , and

$$A_p \triangleq \begin{bmatrix} a_{11} & \Delta_2 & 0 & \cdots & 0 & 0 \\ a_{21} & \gamma_2 & \Delta_3 & \cdots & 0 & 0 \\ & & & \ddots & & \\ a_{(K-1)1} & 0 & 0 & \cdots & \gamma_{K-1} & \Delta_K \\ a_{K1} & 0 & 0 & \cdots & 0 & \gamma_K \end{bmatrix}$$

with  $a_{11} = \gamma_1 - \Delta_1 \alpha_1$  and  $a_{k1} = -\Delta_1 \alpha_k, k \in \{2, 3, \dots, K\}$ .

## 2.2 Scaling Scheme

It is well known that in an implementation system, all the signals should be sustained within a certain dynamic range in order to avoid overflow. Under the assumption that the input  $r(n)$  and the output  $u(n)$  of the closed-loop system are properly pre-scaled, the only signals which may have overflow are the elements of the controller state vector  $x(n)$ , which, therefore, have to be scaled.

There exist different scaling schemes for preventing variables from overflow. The popularly used ones are the  $l_2$ - and  $l_\infty$ -scalings. In what follows, we will concentrate on the  $l_2$ -scaling scheme. The  $l_2$ -scaling means that each element of the controller state vector  $x(n)$  should have a unit variance when the input  $r(n)$  is a white noise with a unit variance. This can be achieved if

$$\bar{\mathcal{X}}(l, l) = 1, \quad l = N+1, N+2, \dots, N+K \quad (9)$$

where  $\bar{\mathcal{X}}$  is the controllability Gramian of the closed-loop system of order  $N+K$ . Assuming that  $P_d(z)$  is strictly proper and has a realization  $(A_z, B_z, C_z, 0)$ , let  $(A_{cl}, B_{cl}, C_{cl}, 0)$  be the closed-loop realization, where

$$\begin{aligned} A_{cl} &= \begin{bmatrix} A_z + dB_z C_z & B_z C \\ BC_z & A \end{bmatrix} \\ B_{cl} &= \begin{bmatrix} B_z \\ \mathbf{0} \end{bmatrix} \\ C_{cl} &= [C_z \ \mathbf{0}] \end{aligned} \quad (10)$$

with  $\mathbf{0}$  denoting the zero vector of appropriate dimension. Then  $\bar{\mathcal{X}}$  is given by

$$\bar{\mathcal{X}} = \sum_{k=0}^{+\infty} A_{cl}^k B_{cl} B_{cl}^T (A_{cl}^T)^k \quad (11)$$

satisfying

$$\bar{\mathcal{X}} = A_{cl} \bar{\mathcal{X}} A_{cl}^T + B_{cl} B_{cl}^T.$$

Let  $(A_{cl}, B_{cl}, C_{cl})$  and  $(A_{cl}^0, B_{cl}^0, C_{cl}^0)$  be two realizations of the closed-loop system with  $A_{cl}, B_{cl}$  and  $C_{cl}$  defined in (10), corresponding to the two digital controller realizations  $R \triangleq (A, B, C, d)$  and  $R_0 \triangleq (A_0, B_0, C_0, d)$  which are related with (3), respectively. It can be shown that

$$\begin{aligned} A_{cl} &= \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & T \end{bmatrix}^{-1} A_{cl}^0 \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & T \end{bmatrix} \\ B_{cl} &= \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & T \end{bmatrix}^{-1} B_{cl}^0 \\ C_{cl} &= C_{cl}^0 \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & T \end{bmatrix}. \end{aligned} \quad (12)$$

It then follows from (12) that

$$\bar{\mathcal{X}} = \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & T \end{bmatrix}^{-1} \bar{\mathcal{X}}^0 \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & T \end{bmatrix}^{-T}$$

where  $\bar{\mathcal{X}}^0$  is the closed-loop controllability Gramian corresponding to  $R_0$ . Let

$$\bar{\mathcal{X}} \triangleq \begin{bmatrix} \mathcal{X}_{11} & \mathcal{X}_{12} \\ \mathcal{X}_{21} & \mathcal{X} \end{bmatrix}, \quad \bar{\mathcal{X}}^0 \triangleq \begin{bmatrix} \mathcal{X}_{11}^0 & \mathcal{X}_{12}^0 \\ \mathcal{X}_{21}^0 & \mathcal{X}_0 \end{bmatrix} \quad (13)$$

have the same partition as  $\begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & T \end{bmatrix}$ , then

$$\mathcal{X} = T^{-1} \mathcal{X}_0 T^{-T} \quad (14)$$

where  $\mathcal{X}_0$  is a positive-definite matrix independent of  $T$ .

It is easy to see from above equations that the  $l_2$ -scaling constraint (9) can be satisfied if the diagonal elements of  $\mathcal{X}$  are all equal to one, that is

$$\mathcal{X}(k, k) = 1, \forall k. \quad (15)$$

When the pDFIIt structure is used to implement a digital controller, it has to be  $l_2$ -scaled in order to prevent the signals in the controller from overflow, which can be achieved by choosing  $\{\Delta_k\}$  properly. It is interesting to note that

$$p_k(z) = \left[ \prod_{l=k+1}^K \Delta_l^{-1} \right] \bar{p}_k(z), \quad \forall k \quad (16)$$

where all  $\bar{p}_k(z)$  are obtained using (5) with  $\Delta_k = 1, \forall k$ .

Let  $(A_\rho^0, B_\rho^0, C_\rho^0, \beta_0)$  be the equivalent state-space realization corresponding to  $\Delta_k = 1, \forall k$ . With (16), it can be shown that

$$A_\rho = T_{sc} A_\rho^0 T_{sc}^{-1}, \quad B_\rho = T_{sc} B_\rho^0, \quad C_\rho = C_\rho^0 T_{sc}^{-1}$$

where  $T_{sc}^{-1}$  is a diagonal scaling similarity transformation, and

$$T_{sc} = \text{diag}(d_1, d_2, \dots, d_K), \quad d_k = \prod_{l=1}^k \Delta_l^{-1}, \quad \forall k.$$

Denote  $\bar{\mathcal{X}}_\rho$  and  $\bar{\mathcal{X}}_\rho^0$  as the closed-loop controllability Gramians, corresponding to the controller realizations  $(A_\rho, B_\rho, C_\rho, \beta_0)$  and  $(A_\rho^0, B_\rho^0, C_\rho^0, \beta_0)$ , respectively. Let  $\mathcal{X}_\rho$  be the sub-matrix of  $\bar{\mathcal{X}}_\rho$  with the partition defined in (13), then (14) becomes  $\mathcal{X}_\rho = T_{sc} \mathcal{X}_\rho^0 T_{sc}^{-T}$  with  $\mathcal{X}_\rho^0$  the corresponding sub-matrix of  $\bar{\mathcal{X}}_\rho^0$ . It is easy to see that the  $l_2$ -scaling can be achieved if  $\mathcal{X}_\rho(k, k) = 1, \forall k$ , or equivalently,

$$d_k^2 \mathcal{X}_\rho^0(k, k) = 1, \quad k = 1, 2, \dots, K$$

which leads to

$$\Delta_1 = \sqrt{\mathcal{X}_\rho^0(1, 1)}, \quad \Delta_k = \sqrt{\frac{\mathcal{X}_\rho^0(k, k)}{\mathcal{X}_\rho^0(k-1, k-1)}}, \quad (17)$$

$$k = 2, 3, \dots, K.$$

In the sequel, all the structures under discussion, including the pDFIIt structure, are assumed to have been  $l_2$ -scaled. Here we should note that the  $l_2$ -scaled pDFIIt structure to be analyzed in this paper is different from the structure in (Hao et al., 2006) where  $\{\Delta_k\}$  are free parameters used for maximizing the stability robustness measure.

### 3 PERFORMANCE ANALYSIS AND OPTIMIZED STRUCTURE

In this section, we will analyze the performance of the pDFIIt structure in terms of closed-loop round-off noise gain. The problem of finding the optimized structure will then be formulated and solved.

One notes that for a given digital controller  $C_d(z)$ , there exists a class of  $l_2$ -scaled pDFIIt structures, which are determined by a space, denoted as  $S_\gamma$ , from which the free parameters  $\{\gamma_k\}$  take values. It is easy to see that  $\{\gamma_k\}$  are the parameters to be implemented directly in the structure. Since we are confined to fixed-point implementation for which the FWL effects are more serious, it is desired that  $\gamma_k$  be absolutely not bigger than one and of FWL format. For a fixed-point implementation of  $B_p$  bits, define

$$S_{FWL} \triangleq \{-1, 1\} \cup \left\{ \pm \sum_{l=1}^{B_p} b_l 2^{-l}, b_l = 0, 1, \forall l \right\} \quad (18)$$

which is a discrete space, containing  $2^{B_p+1} + 1$  elements. Therefore, one can choose  $S_\gamma \subset S_{FWL}$ , which means that all  $\gamma_k$  are of exact  $B_\gamma$ -bit format with  $B_\gamma \leq B_p$ .

#### 3.1 Closed-loop Roundoff Noise Gain

In practice, a designed digital controller has to be implemented with finite precision and a rounding operation has to be applied if less-than-double precision fixed-point arithmetic is utilized. Assuming rounding occurs after multiplication (RAM), a variable, say  $x$ , computed with a multiplication, has to be replaced by its quantized version, denoted as  $Q[x]$ , in the ideal computation model. The difference  $Q[x] - x$  is the corresponding roundoff noise, which is usually modelled as a white noise sequence and statistically independent of those produced by other sources.

Let  $\mu$  be a parameter in a controller structure and  $Q[\mu s(n)]$  the quantized version of the product  $\mu s(n)$ . The roundoff noise due to the parameter  $\mu$  can be defined as

$$\psi(\mu) \varepsilon_\mu(n) \triangleq Q[\mu s(n)] - \mu s(n)$$

where  $\psi(\mu) = 1$  if  $\mu$  is nontrivial, otherwise,  $\psi(\mu) = 0$ . In fact, the function  $\psi(\mu)$  is used for indicating the fact that  $\mu$  produces no roundoff noise when it is trivial. Denote  $\Delta u(n)$  as the corresponding output deviation of the closed-loop system to  $\psi(\mu) \varepsilon_\mu(n)$  and  $F(z)$  as the transfer function between  $\psi(\mu) \varepsilon_\mu(n)$  and  $\Delta u(n)$ . It is well known (see, e.g., (Gevers and Li, 1993)) that  $\Delta u(n)$  is a stationary process and the variance  $E[(\Delta u(n))^2] = \psi(\mu) \|F(z)\|_2^2 E[\varepsilon_\mu^2(n)]$ . Then the

roundoff noise gain for  $\mu$  is defined as

$$G_\mu \triangleq \frac{E[(\Delta u(n))^2]}{E[\varepsilon_\mu^2(n)]} = \Psi(\mu) \|F(z)\|_2^2 \quad (19)$$

where  $\|\cdot\|_2$  is the  $L_2$ -norm:

$$\begin{aligned} \|F(z)\|_2 &\triangleq \left\{ \frac{1}{2\pi} \int_0^{2\pi} \sum_{i=1}^l \sum_{k=1}^m |f_{ik}(e^{j\omega})|^2 d\omega \right\}^{1/2} \\ &= \left\{ \text{tr} \left[ \frac{1}{j2\pi} \oint_{|z|=1} F(z) F^{\mathcal{H}}(z) z^{-1} dz \right] \right\}^{1/2} \end{aligned} \quad (20)$$

with  $F(z) = \{f_{ik}(z)\} \in \mathcal{R}^{l \times m}$ , and  $\mathcal{H}$ ,  $\text{tr}[\cdot]$  denoting the conjugate-transpose and trace operators, respectively. Let  $F(z) = D + L(zI - \Phi)^{-1}J$ . It can be shown that

$$\|F(z)\|_2^2 = \text{tr}[DD^T + LW_c L^T] = \text{tr}[D^T D + J^T W_o J] \quad (21)$$

where  $W_c, W_o$  are the controllability and observability Gramians of the realization  $(\Phi, J, L, D)$ , respectively.

Consider a digital controller implemented with a  $\rho$ DFII structure. We note that the parameters in a  $\rho$ DFII structure are  $\{\alpha_k\}$ ,  $\{\beta_k\}$ ,  $\{\Delta_k\}$ , and  $\{\gamma_k\}$ . It follows from (6) that

$$y(n) = \beta_0 u(n) + \sum_{l=1}^K \left[ \frac{p_l(z)}{p_0(z)} \beta_l u(n) - \frac{p_l(z)}{p_0(z)} \alpha_l y(n) \right]. \quad (22)$$

Let us first look at the effect of roundoff noise  $\Psi(\beta_0)\varepsilon_{\beta_0}(n)$  due to  $\beta_0$  on the closed-loop output. Let  $u^*(n)$  and  $y^*(n)$  be the corresponding output of the closed-loop system and the controller, respectively. Clearly, they obey (22) with  $\beta_0 u^*(n)$  replaced by  $\beta_0 u^*(n) + \Psi(\beta_0)\varepsilon_{\beta_0}(n)$ . Denote  $\Delta y(n) \triangleq y^*(n) - y(n)$ . Then one can show that

$$\begin{aligned} \Delta y(n) &= [\beta_0 \Delta u(n) + \Psi(\beta_0)\varepsilon_{\beta_0}(n)] \\ &+ \sum_{l=1}^K \frac{p_l(z)}{p_0(z)} \beta_l \Delta u(n) - \sum_{l=1}^K \frac{p_l(z)}{p_0(z)} \alpha_l \Delta y(n) \end{aligned} \quad (23)$$

where  $\Delta u(n) \triangleq u^*(n) - u(n)$ , satisfying

$$\Delta u(n) = P_d(z) \Delta y(n). \quad (24)$$

Let  $H_{cl}(z)$  be the transfer function of the closed-loop system, which is given by

$$H_{cl}(z) = \frac{P_d(z)}{1 - P_d(z)C_d(z)}$$

where  $P_d(z)$  is the transfer function of plant and  $C_d(z)$  the polynomial parametrized controller transfer function given by (6). It is easy to see that

$$H_{cl}(z) = D_{cl} + C_{cl}(zI - A_{cl})^{-1}B_{cl} \quad (25)$$

with  $(A_{cl}, B_{cl}, C_{cl}, D_{cl})$  the realization of closed-loop system. It then follows from (23) and (24) that

$$\Delta u(n) = S_0(z) \Psi(\beta_0) \varepsilon_{\beta_0}(n)$$

where  $S_0(z)$  is the transfer function between  $\Psi(\beta_0)\varepsilon_{\beta_0}(n)$  and  $\Delta u(n)$ , which is given by

$$S_0(z) = H_{cl}(z) V_0(z)$$

with

$$V_0(z) \triangleq \frac{p_0(z)}{p_0(z) + \sum_{l=1}^K \alpha_l p_l(z)}.$$

Comparing  $V_0(z)$  with (6), it follows from (8) that

$$\begin{aligned} V_0(z) &= [\beta_0 + C_\rho(zI - A_\rho)^{-1}B_\rho]_{\beta_0=1, \tilde{V}_\beta=0} \\ &= 1 - C_\rho(zI - A_\rho)^{-1}\tilde{V}_\alpha. \end{aligned}$$

One observes that  $S_0(z)$  is of the form  $S_0(z) = [D_2 + C_2(zI_2 - A_2)^{-1}B_2][D_1 + C_1(zI_1 - A_1)^{-1}B_1]$ , where  $A_1 = A_\rho, B_1 = -\tilde{V}_\alpha, C_1 = C_\rho, D_1 = 1, A_2 = A_{cl}, B_2 = B_{cl}, C_2 = C_{cl}, D_2 = D_{cl}$ , and  $I_k, k = 1, 2$  denotes the identity matrix of a proper dimension. It is easy to verify that

$$S_0(z) \triangleq \tilde{D} + \tilde{C}(z\tilde{I} - \tilde{A})^{-1}\tilde{B}$$

where

$$\begin{aligned} \tilde{D} &= D_2 D_1, \quad \tilde{C} = [D_2 C_1 \quad C_2] \\ \tilde{I} &= \begin{bmatrix} I_1 & \mathbf{0} \\ \mathbf{0} & I_2 \end{bmatrix} \\ \tilde{A} &= \begin{bmatrix} A_1 & \mathbf{0} \\ B_2 C_1 & A_2 \end{bmatrix}, \quad \tilde{B} = \begin{bmatrix} B_1 \\ B_2 D_1 \end{bmatrix}. \end{aligned}$$

According to (19) and (21), the roundoff noise gain due to parameter  $\beta_0$  is given by

$$\begin{aligned} G_{\beta_0} &= \Psi(\beta_0) \|S_0(z)\|_2^2 = \Psi(\beta_0) \text{tr}(\tilde{D}^T \tilde{D} + \tilde{B}^T \tilde{W} \tilde{B}) \\ &\triangleq \Psi(\beta_0) G_0 \end{aligned}$$

where  $\tilde{W}$  is the observability Gramian of the realization  $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D})$ .

Using the same procedure, one can analyze the roundoff noise gain due to the parameter  $\beta_k$ . Let  $\Psi(\beta_k)\varepsilon_{\beta_k}(n)$  be the corresponding roundoff noise. It can be shown that the transfer function from  $\Psi(\beta_k)\varepsilon_{\beta_k}(n)$  to  $\Delta u(n)$ , denoted as  $S_k(z)$ , is

$$S_k(z) = H_{cl}(z) V_k(z)$$

with  $H_{cl}(z)$  given by (25) and

$$V_k(z) = \frac{p_k(z)}{p_0(z) + \sum_{l=1}^K \alpha_l p_l(z)} = C_\rho(zI - A_\rho)^{-1}e_k$$

for  $k = 1, 2, \dots, K$ , where  $e_k$  is the  $k$ th elementary vector whose elements are all zero except the  $k$ th one which is 1. Therefore,

$$G_{\beta_k} = \Psi(\beta_k) \|S_k(z)\|_2^2 \triangleq \Psi(\beta_k) G_k, \quad \forall k$$

with

$$G_k = \text{tr}(\tilde{D}_k^T \tilde{D}_k + \tilde{B}_k^T \tilde{W}_k \tilde{B}_k)$$

where  $(\tilde{A}_k, \tilde{B}_k, \tilde{C}_k, \tilde{D}_k)$  is the realization of  $S_k(z)$  and  $\tilde{W}_k$  is the corresponding observability Gramian.

Comparing the positions of  $\alpha_k, \gamma_k$  and  $\Delta_{k+1}$  with that of  $\beta_k$  in Fig. 3, one can see easily that

$$G_{\alpha_k} = \Psi(\alpha_k)G_k, G_{\gamma_k} = \Psi(\gamma_k)G_k, G_{\Delta_k} = \Psi(\Delta_k)G_{k-1}$$

for  $k = 1, 2, \dots, K$ .

Therefore, the total closed-loop roundoff noise gain of the pDFIIt structure is

$$G_p \triangleq \sum_{k=1}^K [G_{\alpha_k} + G_{\gamma_k} + G_{\Delta_k}] + \sum_{k=0}^K G_{\beta_k} \triangleq \sum_{k=0}^K \nu_k G_k \quad (26)$$

where the coefficients  $\nu_k$  can be specified easily with the expressions, obtained above, of roundoff noise gain for all the parameters.

### 3.2 Structure Optimization

For a given digital controller  $C_d(z)$  and any given free parameters  $\{\gamma_k\}$ , one can obtain the  $l_2$ -scaled pDFIIt structure with the procedure presented in Section II. The roundoff noise gain  $G_p$  can then be evaluated with (26). Since different sets of  $\{\gamma_k\}$  yield different pDFIIt structures and hence lead to different roundoff noise gain  $G_p$ , an interesting problem is to minimize  $G_p$  with respect to these free parameters, which leads to the following optimal pDFIIt structure problem:

$$\min_{\gamma_k \in S_\gamma} G_p. \quad (27)$$

It seems impossible to obtain analytical solutions to the problem (27) due to the high nonlinearity of  $G_p$  in  $\{\gamma_k\}$ . However, noting that  $S_\gamma$  is of finite number of elements, the problem can be well solved using the exhaustive searching method.

## 4 A DESIGN EXAMPLE

In this section, we illustrate our design procedure and the performance of the proposed structure with a numerical example, in which  $S_\gamma = \{\pm 1, \pm(2^{-1} + 2^{-2}), \pm 2^{-1}, \pm 2^{-2}, 0\}$ . The elements in the set  $S_\gamma$  are of exact 3-bit fixed-point format (including one bit for the sign). Using more bits or floating-point formats will lead to a further improved performance, which can also confirm the effectiveness of our design procedure.

Consider a discrete-time control system, where the digital plant  $P_d(z) = 10^{-1} \times \frac{0.0181z^4 + 0.0033z^3 - 0.1628z^2 + 0.0111z + 0.0163}{z^5 - 3.7174z^4 + 5.7458z^3 - 4.6673z^2 + 2.0336z - 0.3953}$

Table 1: Comparison of Different Structures.

	zDFIIt	$\delta$ DFIIt	$R_f$	$\rho$ DFIIt
$G$	$1.5191 \times 10^4$	7.1763	4.9919	1.0085
$N_p$	19	19	49	24

and controller  $C_d(z) = 0.0577 + \frac{0.2258z^5 - 0.6588z^4 + 0.8195z^3 - 0.5320z^2 + 0.1814z - 0.0234}{z^6 - 3.6172z^5 + 5.9513z^4 - 5.6335z^3 + 3.2509z^2 - 1.0895z + 0.1690}$ .

The corresponding poles of the closed-loop system are  $\{0.4523 \pm j0.5315, 0.4837 \pm j0.4556, 0.6055 \pm j0.4108, 0.7814 \pm j0.3099, 0.8886 \pm j0.3326, 0.9113\}$ .

Applying exhaustive searching to (27), one gets the optimal pDFIIt structure, denoted as pDFIIt, for which  $\gamma_1 = 1, \gamma_5 = 0.5, \gamma_k = 0.75, k \in \{2, 3, 4, 6\}$ . For comparison, an optimal fully parametrized state-space realization, denoted by  $R_f$ , is obtained using the procedure in (Li et al., 2002). zDFIIt and  $\delta$ DFIIt are the traditional DFIIt structures in the shift- and  $\delta$ -operators, corresponding to  $\gamma_k = 0, \forall k$  and  $\gamma_k = 1, \forall k$ , respectively.

The comparative results of different structures are presented in Table I, where  $G$  is the roundoff noise gain and  $N_p$  is the number of nontrivial parameters in each structure.

From this example, one can see that zDFIIt yields a very large roundoff noise gain, though it has only 19 parameters to implement, while  $\delta$ DFIIt has a much better performance. The fully parametrized optimal realization  $R_f$  yields a further better performance, however, all the 49 parameters in  $R_f$  are nontrivial. It is interesting to see that pDFIIt beats  $R_f$  in terms of the roundoff noise performance. Moreover, pDFIIt is very sparse and has only 24 nontrivial parameters, which is less than half of those in  $R_f$ .

## 5 CONCLUSIONS

In this paper, we have addressed the optimal controller structure problem in a discrete-time control system with roundoff noise consideration. Our major contribution is twofold. Firstly, a sparse controller structure, which is a  $l_2$ -scaled generalized DFIIt structure, has been derived. Secondly, the performance of the proposed structure has been analyzed by deriving the corresponding expression of closed-loop roundoff noise gain and the problem of finding optimized sparse structures has been solved. Finally, a numerical example has been given, which shows that the proposed structure can achieve much better performance than some well-known structures and particularly, outperforms the traditional optimal fully

parametrized realization greatly in terms of reducing roundoff noise and implementation complexity. This optimal controller design strategy with high precision arithmetic can be utilized to develop suitable control systems for robotic platforms performing complex movements, where efficiency, accuracy and fast speed are essential.

## REFERENCES

- Fialho, I. J. and Georgiou, T. T. (1994). On stability and performance of sampled-data systems subject to wordlength constraint. *IEEE Trans. Automat. Contr.*, 39:2476–2481.
- Gevers, M. and Li, G. (1993). *Parametrizations in Control, Estimation and Filtering Problems: Accuracy Aspects*. Springer-Verlag, London, U.K.
- Hao, J., Li, G., and Wan, C. (2006). Two classes of efficient digital controller structures with stability consideration. *IEEE Trans. Automat. Contr.*, 51:164–170.
- Hwang, S. Y. (1977). Minimum uncorrelated unit noise in state-space digital filtering. *IEEE Trans. Acoust., Speech, Signal Processing*, 25:273–281.
- Li, G. (1998). On the structure of digital controllers with finite word length consideration. *IEEE Trans. Automat. Contr.*, 43:689–693.
- Li, G. and Gevers, M. (1990). Optimal finite-precision implementation of a state-estimate feedback controller. *IEEE Trans. Circuits Syst.*, 38:1487–1499.
- Li, G., Wu, J., Chen, S., and Zhao, K. Y. (2002). Optimum structures of digital controllers in sampled-data systems: a roundoff noise analysis. *IEE Proc. Control Theory Appl.*, 149:247–255.
- Liu, K., Skelton, R., and Grigoriadis, K. (1992). Optimal controllers for finite wordlength implementation. *IEEE Trans. Automat. Contr.*, 37:1294–1304.
- Mullis, C. T. and Roberts, R. A. (1976). Synthesis of minimum roundoff noise fixed-point digital filters. *IEEE Trans. Circuits Syst.*, 23:551–562.
- Roberts, R. A. and Mullis, C. T. (1987). *Digital Signal Processing*. Addison Wesley.
- Williamson, D. and Kadiman, K. (1989). Optimal finite wordlength linear quadratic regulation. *IEEE Trans. Automat. Contr.*, 34:1218–1228.
- Wong, N. and Ng, T. S. (2000). Roundoff noise minimization in a modified direct form delta operator iir structure. *IEEE Trans. Circuits Syst. II*, 47:1533–1536.
- Wong, N. and Ng, T. S. (2001). A generalized direct-form delta operator-based iir filter with minimum noise gain and sensitivity. *IEEE Trans. Circuit Syst. II*, 48:425–431.
- Wu, J., Chen, S., Li, G., and Chu, J. (2003). Constructing sparse realisations of finite-precision digital controllers based on a closed-loop stability related measure. *IEE Proc. Control Theory Appl.*, 150:61–68.
- Wu, J., Chen, S., Li, G., Istepanian, R. H., and Chu, J. (2001). An improved closed-loop stability related measure for finite-precision digital controller realizations. *IEEE Trans. Automat. Contr.*, 46:1162–1166.
- Yu, W. S. and Ko, H. J. (2003). Improved eigenvalue sensitivity for finite-precision digital controller realizations via orthogonal hermitian transform. *IEE Proc. Control Theory Appl.*, 50:365–375.