

## Optimization of block-floating-point realizations for digital controllers with finite-word-length considerations\*

WU Jun(吴俊)<sup>†1</sup>, HU Xie-he(胡协和)<sup>1</sup>,  
CHEN Sheng(陈生)<sup>2</sup>, CHU Jian(褚健)<sup>1</sup>

<sup>(1)</sup> *National Key Laboratory of Industrial Control Technology, Institute of Advanced Process Control, Zhejiang University, Hangzhou 310027, China*

<sup>(2)</sup> *Department of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, UK*

<sup>†</sup>E-mail: jwu@ipc.zju.edu.cn

Received Nov.21,2002; revision accepted Mar.3,2003

**Abstract:** The closed-loop stability issue of finite-precision realizations was investigated for digital controllers implemented in block-floating-point format. The controller coefficient perturbation was analyzed resulting from using finite word length (FWL) block-floating-point representation scheme. A block-floating-point FWL closed-loop stability measure was derived which considers both the dynamic range and precision. To facilitate the design of optimal finite-precision controller realizations, a computationally tractable block-floating-point FWL closed-loop stability measure was then introduced and the method of computing the value of this measure for a given controller realization was developed. The optimal controller realization is defined as the solution that maximizes the corresponding measure, and a numerical optimization approach was adopted to solve the resulting optimal realization problem. A numerical example was used to illustrate the design procedure and to compare the optimal controller realization with the initial realization.

**Key words:** Digital controller, Finite word length, Block-floating-point, Closed-loop stability, Optimization  
**Document code:** A **CLC number:** O224

### INTRODUCTION

Due to the finite word length (FWL) effect, a casual controller implementation may degrade the designed closed-loop performance or even destabilize the designed stable closed-loop system, if the controller implementation structure is not carefully chosen. The effects of finite-precision computation have become more critical with the growing popularity of robust controller design methods which focus only on dealing with large plant uncertainty (Keel and Bhattacharyya, 1997).

A control law can be accomplished with different realizations and the parameters of a controller realization are represented by a digital processor of finite bit length in a particular number representation format, such as fixed-point, floating-point or block-floating-point format. In a

given representation format, different controller realizations have different degrees of "robustness" against FWL errors. This property can be utilized to select "optimal" realizations in the given format. In fixed-point format or floating-point format, the optimal controller realization problems were studied in Gevers and Li (1993), Istepanian and Whidborne (2001), Fialho and Georgiou (1994; 1999), Li (1998), Whidborne *et al.* (2000; 2001), Wu *et al.* (2001) and Whidborne and Gu (2001). A comparative study was presented by Istepanian *et al.* (2000) on the stability and performance using block-floating-point and fixed-point implementations for various realizations for a PID digital controller of a steel rolling mill benchmark system. However the optimal controller realization problem in block-floating-point format was not discussed there. To date the true block-floating-point FWL

\* Project supported by the National Natural Science Foundation of China (No.60174026) and the Scientific Research Foundation for Returned Overseas Chinese Scholars of Zhejiang Province (No.J20020546)

closed-loop stability measure has not been seen which can be optimized to obtain the optimal block-floating-point realization. This work is aimed to study the optimal controller realization problem in block-floating-point format.

### BLOCK-FLOATING-POINT

The fixed-point and floating-point formats are the two basic representation schemes for real numbers stored in digital memory and in digital registers. For a group of real numbers stored simultaneously in a digital processor, the so-called block-floating-point format is also available. Suppose that the group of real numbers form a set  $S$ . In the block-floating-point format,  $S$  is divided into some blocks. The block-floating-point scheme may be viewed as aiming to achieve a trade-off between the simplicity of the fixed-point scheme and the accuracy of the floating-point scheme.

For illustrative purpose and without loss of generality, consider the case of dividing  $S$  into two non-empty subsets  $S_1$  and  $S_2$ , which satisfy  $S_1 \cup S_2 = S$  and  $S_1 \cap S_2$  is the empty set. Let  $\eta_1$  be the element in  $S_1$  that has the largest absolute value, and  $\eta_2$  be the element in  $S_2$  that has the largest absolute value. Then, any  $x \in S$  can be expressed uniquely as

$$x = (-1)^s \times u \times 2^h \tag{1}$$

where  $s \in \{0, 1\}$  is the sign of  $x$ ,  $u \in [0, 1)$  is the block mantissa of  $x$ , and the block exponent of  $x$  is

$$h \triangleq \begin{cases} \log_2 \lfloor \eta_1 \rfloor + 1, & \text{for } x \in S_1, \\ \log_2 \lfloor \eta_2 \rfloor + 1, & \text{for } x \in S_2, \end{cases} \tag{2}$$

•  $\lfloor \cdot \rfloor$  denotes the floor function, i.e.,  $x \rfloor$  is the closest integer less than or equal to  $x$ . Obviously, all the elements in the same block have the same exponent value of  $h$ . When all the elements in  $S$  are stored in a digital processor of the bit length

$$\beta = 1 + \beta_u + \beta_h \tag{3}$$

in a block-floating-point scheme, the bits are assigned as follows: 1 bit for the sign,  $\beta_u$  bits for  $u$  which is represented in fixed-point with the two's complement system, and  $\beta_h$  bits for  $h$ . Thus the set of all the block-floating-point num-

bers that can be represented by the bit length  $\beta$  is given by

$$\mathcal{F} \triangleq \left\{ \left( \sum_{j=1}^{\beta_u} b_j 2^{-j} - s \right) \times 2^h : s \in \{0, 1\}, b_j \in \{0, 1\}, h \in Z, \underline{h} \leq h \leq \bar{h} \right\} \tag{4}$$

where  $Z$  denotes the set of integers,  $\underline{h}$  and  $\bar{h}$  represent the lower and upper limits of the block exponent, respectively, and  $\bar{h} - \underline{h} = 2^{\beta_h} - 1$ . Obviously, when  $h > \bar{h}$  or  $h < \underline{h}$ , overflow or underflow will occur in the block-floating-point representation.

When no underflow or overflow occurs, that is,  $h$  is within  $Z_{[\underline{h}, \bar{h}]}$ , the block-floating-point quantization operator  $\mathcal{Q}: S \rightarrow \mathcal{F}$  is defined as

$$\mathcal{Q}(x) \triangleq (-1)^s 2^{(h-\beta_u)} \lfloor 2^{(\beta_u-h)} |x| + 0.5 \rfloor. \tag{5}$$

The quantization error of the block-floating-point representation is defined as

$$\epsilon \triangleq |x - \mathcal{Q}(x)|. \tag{6}$$

Denote

$$r(x) \triangleq \begin{cases} 2^{\log_2 \lfloor \eta_1 \rfloor + 1}, & \text{for } x \in S_1, \\ 2^{\log_2 \lfloor \eta_2 \rfloor + 1}, & \text{for } x \in S_2. \end{cases} \tag{7}$$

It can be shown easily that the quantization error is bounded by

$$\epsilon < r(x) 2^{-(\beta_u+1)} \tag{8}$$

Thus, when  $x \in S$  is implemented in the block-floating-point format of  $\beta_u$  block mantissa bits, assuming no underflow or overflow, it is perturbed to

$$Q(x) = x + r(x)\delta, |\delta| < 2^{-(\beta_u+1)}. \tag{9}$$

Hence the perturbation resulting from FWL block-floating-point representation is neither multiplicative nor additive. The perturbation depends on the set  $S$  and how  $S$  is divided into blocks. It can also be seen that the dynamic range of block-floating-point representation is determined by  $\beta_h$  bits, and the precision of block-floating-point representation is determined by  $\beta_u$  bits.

### PROBLEM STATEMENT

Consider the discrete-time closed-loop con-

trol system consisting of a linear time-invariant plant  $P$  and a digital controller  $C$ . The plant model  $P$  is assumed to be strictly proper with a state-space description

$$\begin{cases} \mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{e}(k) \\ \mathbf{y}(k) = \mathbf{C}\mathbf{x}(k) \end{cases} \quad (10)$$

which is completely state controllable and observable with  $\mathbf{A} \in R^{n \times n}$ ,  $\mathbf{B} \in R^{n \times p}$  and  $\mathbf{C} \in R^{q \times n}$ . The digital controller  $C$  is described by

$$\begin{cases} \mathbf{v}(k+1) = \mathbf{F}\mathbf{v}(k) + \mathbf{G}\mathbf{y}(k) \\ \mathbf{e}(k) = \mathbf{J}\mathbf{v}(k) + \mathbf{M}\mathbf{y}(k) \end{cases} \quad (11)$$

with  $\mathbf{F} \in R^{m \times m}$ ,  $\mathbf{G} \in R^{m \times q}$ ,  $\mathbf{J} \in R^{p \times m}$  and  $\mathbf{M} \in R^{p \times q}$ .

Assume that a realization  $(\mathbf{F}_0, \mathbf{G}_0, \mathbf{J}_0, \mathbf{M}_0)$  of  $C$  has been designed. It is well-known that the realizations of  $C$  are not unique. All the realizations of  $C$  form the realization set

$$\begin{aligned} S_C &\triangleq \{(\mathbf{F}, \mathbf{G}, \mathbf{J}, \mathbf{M}) : \mathbf{F} = \mathbf{T}^{-1}\mathbf{F}_0\mathbf{T}, \\ &\mathbf{G} = \mathbf{T}^{-1}\mathbf{G}_0, \mathbf{J} = \mathbf{J}_0\mathbf{T}, \mathbf{M} = \mathbf{M}_0\} \end{aligned} \quad (12)$$

where  $\mathbf{T} \in R^{m \times m}$  is any real-valued nonsingular matrix, called a similarity transformation. Denote

$$\mathbf{X} = [x_{j,k}] \triangleq \begin{bmatrix} \mathbf{M} & \mathbf{J} \\ \mathbf{G} & \mathbf{F} \end{bmatrix}. \quad (13)$$

We also refer to  $\mathbf{X}$  as a realization of  $C$ . The stability of the closed-loop system depends on the eigenvalues of the matrix

$$\begin{aligned} \bar{\mathbf{A}}(\mathbf{X}) &= \begin{bmatrix} \mathbf{A} + \mathbf{B}\mathbf{M}\mathbf{C} & \mathbf{B}\mathbf{J} \\ \mathbf{G}\mathbf{C} & \mathbf{F} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \mathbf{X} \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \\ &\triangleq \mathbf{M}_0 + \mathbf{M}_1 \mathbf{X} \mathbf{M}_2 \end{aligned} \quad (14)$$

where  $\mathbf{0}$  and  $\mathbf{I}$  denote the zero and identity matrices of appropriate dimensions, respectively. All the different realizations  $\mathbf{X}$  have the same set of closed-loop poles if they are implemented with infinite precision. Since the closed-loop system is designed to be stable, the eigenvalues

$$|\lambda_i(\bar{\mathbf{A}}(\mathbf{X}))| < 1, \quad \forall i \in \{1, \dots, m+n\}. \quad (15)$$

For a matrix  $\mathbf{W} = [w_{j,k}]$ , let  $\mathbf{U}(\mathbf{W})$  be the matrix of the same dimension whose elements are all 1s,

$$\|\mathbf{W}\|_{\max} \triangleq \max_{j,k} |w_{j,k}|, \quad (16)$$

$$\pi(\mathbf{W}) \triangleq \min_{j,k} \{|w_{j,k}| : w_{j,k} \neq 0\}, \quad (17)$$

For two matrices  $\mathbf{W} = [w_{j,k}]$  and  $\mathbf{Z} = [z_{j,k}]$  of the same dimension, define the Hadamard product of  $\mathbf{W}$  and  $\mathbf{Z}$

$$\mathbf{W} \circ \mathbf{Z} \triangleq [w_{j,k} z_{j,k}]. \quad (18)$$

We have known that the controller realization  $\mathbf{X}$  is implemented in block-floating-point format of  $\beta_h$  block exponent bits,  $\beta_u$  block mantissa bits and one sign bit. In the remainder of this paper, it is assumed that  $\mathbf{X}$  stored in the block-floating-point format is divided into "natural" blocks of  $\mathbf{F}$ ,  $\mathbf{G}$ ,  $\mathbf{J}$  and  $\mathbf{M}$ . Let  $\xi_1$  be the element in  $\mathbf{F}$  which has the largest absolute value. Similarly  $\xi_2$ ,  $\xi_3$  and  $\xi_4$  is defined in  $\mathbf{G}$ ,  $\mathbf{J}$  and  $\mathbf{M}$  respectively. Denote

$$\mathbf{q}(\mathbf{X}) \triangleq [\xi_1 \quad \xi_2 \quad \xi_3 \quad \xi_4]^T \quad (19)$$

with  $^T$  being the transpose operator.

Firstly, the dynamic range of  $\beta_h$  bits must be large enough for  $\mathbf{X}$ . We define a dynamic range measure for realization  $\mathbf{X}$  in block-floating-point format as

$$\gamma(\mathbf{X}) \triangleq \log_2 \frac{4 \|\mathbf{q}(\mathbf{X})\|_{\max}}{\pi(\mathbf{q}(\mathbf{X}))}. \quad (20)$$

The rationale of this dynamic range measure becomes clear in the following (obvious) proposition.

**Proposition 1** The realization  $\mathbf{X}$  can be represented in the block-floating-point format of  $\beta_h$  block exponent bits without underflow or overflow, if  $2^{\beta_h} \geq \log_2 \left( \frac{\|\mathbf{q}(\mathbf{X})\|_{\max}}{\pi(\mathbf{q}(\mathbf{X}))} \right) + 2$ .

Let  $\beta_h^{\min}$  be the smallest block exponent bit length that, when used to implement  $\mathbf{X}$ , does not cause overflow or underflow. The minimum required block exponent bit length can easily be computed by

$$\beta_h^{\min}(\mathbf{X}) = \lceil \log_2 (\lfloor \log_2 \|\mathbf{q}(\mathbf{X})\|_{\max} \rfloor - \lfloor \log_2 \pi(\mathbf{q}(\mathbf{X})) \rfloor + 1) \rceil, \quad (21)$$

where  $\lceil \cdot \rceil$  denotes the ceil function, i.e.,  $\lceil x \rceil$  is the closest integer greater than or equal to  $x$ . Note that the measure  $\gamma(\mathbf{X})$  defined in Eq. (20) provides an estimate of  $\beta_h^{\min}$  as

$$\hat{\beta}_h^{\min}(\mathbf{X}) \triangleq \lceil \log_2 \gamma(\mathbf{X}) \rceil. \quad (22)$$

It can easily be seen that  $\hat{\beta}_h^{\min} \geq \beta_h^{\min}$ .

When the dynamic range is sufficient, according to Eq. (9),  $\mathbf{X}$  is perturbed to  $\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta}$  due to the effect of finite  $\beta_u$  where

$$\mathbf{E}(\mathbf{X}) \triangleq \begin{bmatrix} 2^{-\log_2 |\xi_4| + 1} \mathbf{U}(\mathbf{M}) & 2^{-\log_2 |\xi_3| + 1} \mathbf{U}(\mathbf{J}) \\ 2^{-\log_2 |\xi_2| + 1} \mathbf{U}(\mathbf{G}) & 2^{-\log_2 |\xi_1| + 1} \mathbf{U}(\mathbf{F}) \end{bmatrix} \quad (23)$$

Each element  $\delta_{j,k}$  of  $\mathbf{\Delta}$  is bounded by  $\pm 2^{-(\beta_u + 1)}$ , that is,

$$\|\mathbf{\Delta}\|_{\max} < 2^{-(\beta_u + 1)}. \quad (24)$$

With the perturbation  $\mathbf{\Delta}$ ,  $\lambda_i(\bar{\mathbf{A}}(\mathbf{X}))$  is moved to  $\lambda_i(\mathbf{A}(\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta}))$ . If an eigenvalue of  $\mathbf{A}(\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta})$  is outside the open unit disk, the closed-loop system, designed to be stable, becomes unstable with the finite-precision implemented  $\mathbf{X}$ .

It is therefore critical to know when the FWL error will cause closed-loop instability. This means that we would like to know the largest open "cube" in the perturbation space within which the closed-loop system remains stable. Based on this consideration, a precision measure for realization  $\mathbf{X}$  in block-floating-point format can be defined as

$$\mu_0(\mathbf{X}) \triangleq \inf \{ \|\mathbf{\Delta}\|_{\max} : \bar{\mathbf{A}}(\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta}) \text{ is unstable} \}. \quad (25)$$

From the above definition, the following proposition is obvious.

**Proposition 2**  $\bar{\mathbf{A}}(\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta})$  is stable if  $\|\mathbf{\Delta}\|_{\max} < \mu_0(\mathbf{X})$ .

Thus under the condition that the dynamic range is sufficient, that is,  $\beta_h \geq \beta_h^{\min}$ , the perturbation  $\|\mathbf{\Delta}\|_{\max}$  and therefore the block mantissa bit length  $\beta_u$  determines whether the closed-loop remains stable. Let  $\beta_u^{\min}$  be the block mantissa bit length such that  $\forall \beta_u \geq \beta_u^{\min}$ , the closed-loop system is stable with  $\mathbf{X}$  implemented by  $\beta_u$  block mantissa bits and the closed-loop system is unstable with  $\mathbf{X}$  implemented by  $\beta_u^{\min} - 1$  block mantissa bits. The precision measure  $\mu_0(\mathbf{X})$  provides an estimate of  $\beta_u^{\min}$  as

$$\hat{\beta}_{u0}^{\min}(\mathbf{X}) \triangleq -\lfloor \log_2 \mu_0(\mathbf{X}) \rfloor - 1. \quad (26)$$

It can be seen that  $\hat{\beta}_{u0}^{\min} \geq \beta_u^{\min}$ .

Define the minimum total bit length required in the implementation of  $\mathbf{X}$  as

$$\beta^{\min} \triangleq \beta_h^{\min} + \beta_u^{\min} + 1. \quad (27)$$

Clearly,  $\mathbf{X}$  implemented with a bit length  $\beta \geq \beta^{\min}$  can guarantee a sufficient dynamic range and closed-loop stability. Combining the measures  $\gamma(\mathbf{X})$  and  $\mu_0(\mathbf{X})$  results in the following true FWL closed-loop stability measure for the given realization  $\mathbf{X}$  in block-floating-point format

$$\rho_0(\mathbf{X}) \triangleq \mu_0(\mathbf{X}) / \gamma(\mathbf{X}). \quad (28)$$

An estimate of  $\beta^{\min}$  is given by  $\rho_0(\mathbf{X})$  as

$$\hat{\beta}_0^{\min}(\mathbf{X}) \triangleq -\lfloor \log_2 \rho_0(\mathbf{X}) \rfloor + 1. \quad (29)$$

It is clear that  $\hat{\beta}_0^{\min} \geq \beta^{\min}$ . The following proposition summarizes the usefulness of  $\rho_0(\mathbf{X})$  as a measure for the FWL characteristics of  $\mathbf{X}$  in block-floating-point format.

**Proposition 3** The controller realization  $\mathbf{X}$  implemented in block-floating-point format with a bit length  $\beta$  can guarantee a sufficient dynamic range and closed-loop stability, if

$$2^{\beta-1} \geq \frac{1}{\rho_0(\mathbf{X})}. \quad (30)$$

The closed-loop stability measure  $\rho_0(\mathbf{X})$  depends on the controller realization  $\mathbf{X}$  only. Consequently, an optimal realization can in theory be found by maximizing  $\rho_0(\mathbf{X})$  over  $S_C$ , leading to the following optimal controller realization problem

$$\nu_{\text{true}} \triangleq \max_{\mathbf{X} \in S_C} \rho_0(\mathbf{X}). \quad (31)$$

However, the difficulty with this approach is that computing the value of  $\mu_0(\mathbf{X})$  is an unsolved open problem. Thus, the true FWL closed-loop stability measure  $\rho_0(\mathbf{X})$  and the optimal realization problem (31) have limited practical significance. In the next section, an alternative measure is developed which not only can quantify the FWL characteristics of  $\mathbf{X}$  in block-floating-point format but also is computationally tractable.

## A TRACTABLE FWL CLOSED-LOOP STABILITY MEASURE

When the FWL error  $\mathbf{\Delta}$  is small, from a

first-order approximation,  $\forall i \in \{1, \dots, m+n\}$   
 $|\lambda_i(\mathbf{A}(\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta}))| - |\lambda_i(\bar{\mathbf{A}}(\mathbf{X}))| \approx$   
 $\sum_{j,k} \frac{\partial |\lambda_i|}{\partial \delta_{j,k}} \Big|_{\mathbf{\Delta}=\mathbf{0}} \delta_{j,k}$ . (32)

For the derivative  $\frac{\partial |\lambda_i|}{\partial \mathbf{\Delta}} = \left[ \frac{\partial |\lambda_i|}{\partial \delta_{j,k}} \right]$ , define

$$\left\| \frac{\partial |\lambda_i|}{\partial \mathbf{\Delta}} \right\|_{\text{sum}} \triangleq \sum_{j,k} \left| \frac{\partial |\lambda_i|}{\partial \delta_{j,k}} \right|. \quad (33)$$

Then

$$|\lambda_i(\bar{\mathbf{A}}(\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta}))| - |\lambda_i(\bar{\mathbf{A}}(\mathbf{X}))| \leq$$

$$\|\mathbf{\Delta}\|_{\max} \left\| \frac{\partial |\lambda_i|}{\partial \mathbf{\Delta}} \Big|_{\mathbf{\Delta}=\mathbf{0}} \right\|_{\text{sum}}. \quad (34)$$

This leads to the following precision measure for realization  $\mathbf{X}$  in block-floating-point format

$$\mu_1(\mathbf{X}) \triangleq \min_{i \in \{1, \dots, m+n\}} \frac{1 - |\lambda_i(\bar{\mathbf{A}}(\mathbf{X}))|}{\left\| \frac{\partial |\lambda_i|}{\partial \mathbf{\Delta}} \Big|_{\mathbf{\Delta}=\mathbf{0}} \right\|_{\text{sum}}}. \quad (35)$$

Obviously, if  $\|\mathbf{\Delta}\|_{\max} < \mu_1(\mathbf{X})$ , then  $|\lambda_i(\bar{\mathbf{A}}(\mathbf{X} + \mathbf{E}(\mathbf{X}) \circ \mathbf{\Delta}))| < 1$  which means that the closed-loop remains stable under the FWL error  $\mathbf{\Delta}$ . In other words, for a given  $\mathbf{X}$  implemented in block-floating-point format with a sufficient dynamic range, the closed-loop can tolerate those FWL perturbations  $\mathbf{\Delta}$  whose norms  $\|\mathbf{\Delta}\|_{\max}$  are less than  $\mu_1(\mathbf{X})$ . The larger  $\mu_1(\mathbf{X})$  is, the larger the FWL errors the closed-loop system can tolerate. Similar to Eq. (26), from the precision measure  $\mu_1(\mathbf{X})$ , an estimate of  $\beta_u^{\min}$  is given as

$$\hat{\beta}_{u1}^{\min}(\mathbf{X}) \triangleq -\lfloor \log_2 \mu_1(\mathbf{X}) \rfloor - 1. \quad (36)$$

The assumption of small  $\mathbf{\Delta}$  is usually valid in practical implementation of digital controllers. Generally speaking, there is no rigorous relationship between  $\mu_0(\mathbf{X})$  and  $\mu_1(\mathbf{X})$ , but  $\mu_1(\mathbf{X})$  is connected with a lower bound of  $\mu_0(\mathbf{X})$  in some manner; there are “stable perturbation cubes” larger than  $\{\mathbf{\Delta} : \|\mathbf{\Delta}\|_{\max} < \mu_1(\mathbf{X})\}$  while there is no “stable perturbation cube” larger than  $\{\mathbf{\Delta} : \|\mathbf{\Delta}\|_{\max} < \mu_0(\mathbf{X})\}$  (Wu *et al.*, 2001). Hence, in most cases, it is reasonable to take that  $\mu_1(\mathbf{X}) \leq \mu_0(\mathbf{X})$  and  $\hat{\beta}_{u1}^{\min} \geq \hat{\beta}_{u0}^{\min}$ . More importantly, unlike the measure  $\mu_0(\mathbf{X})$ , the value of  $\mu_1(\mathbf{X})$  can be computed explicitly. It is easy to see that

$$\frac{\partial |\lambda_i|}{\partial \mathbf{\Delta}} \Big|_{\mathbf{\Delta}=\mathbf{0}} = \mathbf{E}(\mathbf{X}) \circ \frac{\partial |\lambda_i|}{\partial \mathbf{X}} \quad (37)$$

Let  $\mathbf{p}_i$  be a right eigenvector of  $\mathbf{A}(\mathbf{X})$  corresponding to the eigenvalue  $\lambda_i$ . Define

$$\mathbf{M}_p \triangleq [\mathbf{p}_1 \quad \mathbf{p}_2 \quad \dots \quad \mathbf{p}_{m+n}] \quad (38)$$

and

$$\mathbf{M}_y \triangleq [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \dots \quad \mathbf{y}_{m+n}] = \mathbf{M}_p^{-H} \quad (39)$$

where the superscript  $H$  denotes the conjugate transpose operator and  $\mathbf{y}_i$  is called the reciprocal left eigenvector related to  $\mathbf{p}_i$ . The following lemma is due to Li (1998).

**Lemma 1** Let  $\bar{\mathbf{A}}(\mathbf{X}) = \mathbf{M}_0 + \mathbf{M}_1 \mathbf{X} \mathbf{M}_2$  given in Eq. (14) be diagonalizable. Then

$$\frac{\partial \lambda_i}{\partial \mathbf{X}} = \mathbf{M}_1^T \mathbf{y}_i^* \mathbf{p}_i^T \mathbf{M}_2^T \quad (40)$$

where the superscript  $*$  denotes the conjugate operation.

The following proposition shows that, given a  $\mathbf{X}$ , the value of  $\mu_1(\mathbf{X})$  can easily be calculated.

**Proposition 4** Let  $\mathbf{A}(\mathbf{X})$  be diagonalizable. Then

$$\mu_1(\mathbf{X}) = \min_{i \in \{1, \dots, m+n\}} \frac{1 - |\lambda_i|}{\left\| (\mathbf{M}_1^T \text{Re}[\lambda_i^* \mathbf{y}_i^* \mathbf{p}_i^T] \mathbf{M}_2^T) \circ \mathbf{E}(\mathbf{X}) \right\|_{\text{sum}}} \quad (41)$$

**Proof** Noting  $|\lambda_i| = \sqrt{\lambda_i^* \lambda_i}$  leads to

$$\frac{\partial |\lambda_i|}{\partial \mathbf{X}} = \frac{1}{2\sqrt{\lambda_i^* \lambda_i}} \left( \frac{\partial \lambda_i^*}{\partial \mathbf{X}} \lambda_i + \lambda_i^* \frac{\partial \lambda_i}{\partial \mathbf{X}} \right) =$$

$$\frac{1}{2|\lambda_i|} \left( \left( \frac{\partial \lambda_i}{\partial \mathbf{X}} \right)^* \lambda_i + \lambda_i^* \frac{\partial \lambda_i}{\partial \mathbf{X}} \right) =$$

$$\frac{1}{|\lambda_i|} \text{Re} \left[ \lambda_i^* \frac{\partial \lambda_i}{\partial \mathbf{X}} \right]. \quad (42)$$

Combining Eqs. (35), (37), (42) and Lemma 1 results in this proposition.

Replacing  $\mu_0(\mathbf{X})$  with  $\mu_1(\mathbf{X})$  in Eq. (28) leads to a computationally tractable FWL closed-loop stability measure

$$\rho_1(\mathbf{X}) \triangleq \mu_1(\mathbf{X}) / \gamma(\mathbf{X}). \quad (43)$$

From the measure  $\rho_1(\mathbf{X})$ , an estimate of  $\beta^{\min}$  is given as

$$\hat{\beta}_1^{\min}(\mathbf{X}) \triangleq -\lfloor \log_2 \rho_1(\mathbf{X}) \rfloor + 1. \quad (44)$$

OPTIMIZATION PROCEDURE

As different realizations  $\mathbf{X}$  have different values of the FWL closed-loop stability measure  $\rho_1(\mathbf{X})$ , it is of practical importance to find an “optimal” realization  $\mathbf{X}_{opt}$  that maximizes  $\rho_1(\mathbf{X})$ . The controller implemented with this optimal realization  $\mathbf{X}_{opt}$  needs a minimum bit length and has a maximum tolerance to the FWL error in block-floating-point format. This optimal controller realization problem is formally defined as

$$v \triangleq \max_{\mathbf{X} \in S_c} \rho_1(\mathbf{X}). \tag{45}$$

Assume that a controller has been designed using some standard controller design method. This controller, denoted as

$$\mathbf{X}_0 \triangleq \begin{bmatrix} \mathbf{M}_0 & \mathbf{J}_0 \\ \mathbf{G}_0 & \mathbf{F}_0 \end{bmatrix}, \tag{46}$$

is used as the initial controller realization in the above optimal controller realization problem. Let  $\mathbf{p}_{0i}$  be a right eigenvector of  $\bar{\mathbf{A}}(\mathbf{X}_0)$  corresponding to the eigenvalue  $\lambda_i$ , and  $\mathbf{y}_{0i}$  be the reciprocal left eigenvector related to  $\mathbf{p}_{0i}$ . The definition of  $S_c$  in Eq.(12) means that

$$\mathbf{X} \triangleq \mathbf{X}(\mathbf{T}) = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{-1} \end{bmatrix} \mathbf{X}_0 \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T} \end{bmatrix} \tag{47}$$

where  $\det(\mathbf{T}) \neq 0$ . It can then be shown that

$$\bar{\mathbf{A}}(\mathbf{X}) = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{-1} \end{bmatrix} \bar{\mathbf{A}}(\mathbf{X}_0) \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T} \end{bmatrix} \tag{48}$$

which implies that

$$\mathbf{p}_i = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{-1} \end{bmatrix} \mathbf{p}_{0i}, \quad \mathbf{y}_i = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^\top \end{bmatrix} \mathbf{y}_{0i}. \tag{49}$$

Hence

$$\begin{aligned} & \mathbf{M}_1^\top \text{Re}[\lambda_i^* \mathbf{y}_i^* \mathbf{p}_i^\top] \mathbf{M}_2^\top = \\ & \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^\top \end{bmatrix} \mathbf{M}_1^\top \text{Re}[\lambda_i^* \mathbf{y}_{0i}^* \mathbf{p}_{0i}^\top] \mathbf{M}_2^\top \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{-\top} \end{bmatrix} \triangleq \\ & \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^\top \end{bmatrix} \Phi_i \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{-\top} \end{bmatrix} \end{aligned} \tag{50}$$

with  $\Phi_i = \mathbf{M}_1^\top \text{Re}[\lambda_i^* \mathbf{y}_{0i}^* \mathbf{p}_{0i}^\top] \mathbf{M}_2^\top$ . Define the following cost function:

$$f(\mathbf{T}) \triangleq \min_{i \in \{1, \dots, m+n\}} \bullet$$

$$\left( \frac{\left\| \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^\top \end{bmatrix} \Phi_i \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}^{-\top} \end{bmatrix} \circ \mathbf{E}(\mathbf{X}(\mathbf{T})) \right\|_{\text{sum}}}{|\lambda_i| (1 - |\lambda_i|)} \log_2 \frac{4 \left\| \mathbf{q}(\mathbf{X}(\mathbf{T})) \right\|_{\text{max}}}{\pi(\mathbf{q}(\mathbf{X}(\mathbf{T})))} \right)^{-1}. \tag{51}$$

Then the optimal controller realization problem Eq.(45) can be posed as the following optimization problem:

$$v = \max_{\substack{\mathbf{T} \in R^{(m+n)} \\ \det \mathbf{T} \neq 0}} f(\mathbf{T}). \tag{52}$$

As the optimization problem Eq. (52) is highly nonlinear, global optimization algorithms, such as the genetic algorithm and adaptive simulated annealing, can be adopted to provide a (sub)optimal similarity transformation  $\mathbf{T}_{opt}(\alpha)$ . Global optimization methods are however computationally demanding. Local optimization algorithms, such as Rosenbrock and Simplex algorithms, are computationally simpler but run more risks of only attaining a local solution. Our experience with the optimization problem Eq. (29) suggests that local optimization methods are usually efficient in controllers of low order while global optimization methods have to be adopted in controllers of high order. It also helps to choose a “good” initial controller realization, such as Li’s closed-loop sub-optimal realization (Li, 1998), as the initial guess for the optimization routine. With the solution  $\mathbf{T}_{opt}$  of optimization problem Eq.(52), the optimal realization  $\mathbf{X}_{opt}$  can readily be computed.

DESIGN EXAMPLE

An example is used to illustrate the design procedure based on the FWL closed-loop stability measure. In this example, the discrete-time plant is given by

$$\mathbf{A} = \begin{bmatrix} 3.7156e+0 & -5.4143e+0 & 3.6525e+0 & -9.6420e-1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix},$$

$$\mathbf{B} = [1 \ 0 \ 0 \ 0]^\top,$$

$$\mathbf{C} = [1.1160e-6 \ 4.3000e-8 \ 1.0880e-6 \ 1.4000e-8].$$

The initial realization of the digital controller is given by

$$\begin{aligned}
 F_0 &= \\
 &\begin{bmatrix} 2.6963e+2 & -4.2709e+1 & 2.2873e+1 & 2.6184e+2 \\ 2.5561e+2 & -4.0497e+1 & 2.1052e+1 & 2.4806e+2 \\ 5.6096e+1 & -8.5715e+0 & 5.2162e+0 & 5.4920e+1 \\ -2.3907e+2 & 3.7998e+1 & -2.0338e+1 & -2.3203e+2 \end{bmatrix}, \\
 G_0 &= \\
 &[-4.6765e+1 \quad -4.5625e+1 \quad -9.5195e+0 \quad 4.1609e+1]^T, \\
 J_0 &= \\
 &[-2.5548e+2 \quad -2.7185e+2 \quad -2.7188e+2 \quad 2.7188e+2], \\
 M_0 &=[0].
 \end{aligned}$$

Based on the proposed FWL closed-loop stability measure, the optimization problem Eq. (52) is formed. Using the MATLAB routine *fminsearch*. *m* for the controller of 4 order, this optimization problem is solved to obtain the optimal similarity transformation

$$\begin{aligned}
 T_{\text{opt}} &= \\
 &\begin{bmatrix} -1.0345e-001 & 1.2904e-001 & 3.8329e-003 & 1.0911e-002 \\ -1.1078e-001 & 1.1742e-001 & 2.9461e-003 & 8.1639e-003 \\ -2.3775e-002 & 2.3815e-002 & 4.9498e-004 & 1.8293e-003 \\ 9.2138e-002 & -1.1474e-001 & -3.4007e-003 & -9.6780e-003 \end{bmatrix}.
 \end{aligned}$$

It is obvious that the true minimum block exponent bit length  $\beta_h^{\min}(X)$  for a realization  $X$  can directly be obtained by examining the elements of  $X$ . The true minimum block mantissa bit length  $\beta_u^{\min}(X)$  however can only be obtained through simulation. That is, starting from a very large  $\beta_u$ , reduce  $\beta_u$  by one bit and check the closed-loop stability. The process is repeated until there appears closed-loop instability at  $\beta_u = \beta_{uu}$ . Then  $\beta_u^{\min} = \beta_{uu} + 1$ . Table 1 summarizes

**Table 1** Various measures and bit lengths for  $X$  and  $X_{\text{opt}}$

	$X_0$	$X_{\text{opt}}$
$\rho_1(X)$	1.5154e-11	4.7787e-6
$\hat{\beta}_1^{\min}(X)$	37	19
$\mu_1(X)$	6.8793e-11	3.6388e-5
$\hat{\beta}_{u1}^{\min}(X)$	33	14
$\gamma(X)$	4.5395e+0	7.6146e+0
$\hat{\beta}_h^{\min}(X)$	3	3
$\beta^{\min}(X)$	33	16
$\beta_u^{\min}(X)$	30	12
$\beta_h^{\min}(X)$	2	3

the various measures, the corresponding estimated minimum bit lengths and the true minimum bit lengths for the controller realizations  $X_0$  and  $X_{\text{opt}}$ . It can be seen that  $X_{\text{opt}}$  improves the measure  $\rho_1$  by a factor of 300000 over  $X_0$  and that the block-floating-point implemented  $X_0$  needs at least 33 bits while the implementation of  $X_{\text{opt}}$  needs at least 16 bits. More than half of the bit length is saved.

## References

- Fialho, I. J. and Georgiou, T. T., 1994. On stability and performance of sampled-data systems subject to word-length constraint. *IEEE Trans. Automatic Control*, **39** (12): 2476 – 2481.
- Fialho, I. J. and Georgiou, T. T., 1999. Optimal Finite Wordlength Digital Controller Realization. Proc. American Control Conf., San Diego, p.4326 – 4327.
- Gevers, M. and Li, G., 1993. Parameterizations in Control, Estimation and Filtering Problems: Accuracy Aspects. Springer Verlag, London.
- Istefanian, R. S. H., Whidborne, J. F. and Bauer, P., 2000. Stability Analysis of Block Floating Point Digital Controllers. Proc. UKACC Int. Conf. Control, Cambridge, CD-ROM, 6 pages.
- Istefanian, R. S. H. and Whidborne, J. F., 2001. Digital Controller Implementation and Fragility: A Modern Perspective. Springer Verlag, London.
- Keel, L. H. and Bhattacharyya, S. P., 1997. Robust, fragile, or optimal? *IEEE Trans. Automatic Control*, **42**(8): 1098 – 1105.
- Li, G., 1998. On the structure of digital controllers with finite word length consideration. *IEEE Trans. Automatic Control*, **43**(5): 689 – 693.
- Whidborne, J. F., Wu, J. and Istefanian, R. S. H., 2000. Finite word length stability issues in an  $l_1$  framework. *Int. J. Control*, **73**(2): 166 – 176.
- Whidborne, J. F., Istefanian, R. S. H. and Wu, J., 2001. Reduction of controller fragility by pole sensitivity minimization. *IEEE Trans. Automatic Control*, **46** (2): 320 – 325.
- Whidborne, J. F. and Gu, D., 2001. Optimal finite-precision controller and filter realizations using floating-point arithmetic. *King's College London Mechanical Engineering Department Report EM/2001/07*, London.
- Wu, J., Chen, S., Li, G., Istefanian, R. S. H. and Chu, J., 2001. An improved closed-loop stability related measure for finite-precision digital controller realizations. *IEEE Trans. Automatic Control*, **46** (7): 1162 – 1166.