

# Analysis and Minimization of $L_2$ -Sensitivity for Linear Systems and Two-Dimensional State-Space Filters Using General Controllability and Observability Gramians

Takao Hinamoto, *Fellow, IEEE*, Shuichi Yokoyama, Takuya Inoue, Wu Zeng, and Wu-Sheng Lu, *Fellow, IEEE*

**Abstract**—A novel expression for the evaluation of  $L_2$ -sensitivity is developed for the cases of linear discrete-time systems, linear continuous-time systems, and two-dimensional (2-D) state-space digital filters. This is accomplished by introducing the concept of general controllability and observability Gramians in each case. Moreover, the  $L_2$ -sensitivity measures obtained here contain the conventional  $L_1/L_2$ -sensitivity measures as a special case. An iterative procedure for constructing the optimal coordinate transformation matrix that minimizes the  $L_2$ -sensitivity measure is then presented in each case. This procedure is advantageous since the initial estimate and the estimate at each iteration can be calculated analytically. Finally, three numerical examples are given to illustrate the utility of the proposed techniques.

**Index Terms**— $L_2$ -sensitivity, linear continuous-time systems, linear discrete-time systems, optimal realization, sensitivity minimization, two-dimensional state-space digital filters.

## I. INTRODUCTION

ONE OF THE primary finite-length register effects in fixed point digital filters is changes in the input–output description of the filter due to approximating real-number parameters with a finite binary representation. Such an effect is loosely called coefficient sensitivity and has been an important research topic since coefficient truncation or rounding may cause an originally stable filter to be an unstable one. It is well-known that the undesirable finite-word-length (FWL) effects can be reduced considerably by the appropriate selection of the filter structure. Several techniques for synthesizing linear discrete-time systems that minimize the coefficient sensitivity have been reported in [1]–[9]. These techniques can be divided into two main classes:  $L_1/L_2$ -sensitivity minimization [1]–[5] and  $L_2$ -sensitivity minimization [6]–[9]. It has been pointed out [6]–[9] that the sensitivity measure based on the  $L_2$  norm only is natural and reasonable. The  $L_1/L_2$ -sensitivity minimization has also been considered in linear continuous-time systems [10]. The problem of minimizing the coefficient sensitivity measure evaluated by using a mixture of  $L_1/L_2$  norms has been studied for two-di-

dimensional (2-D) state-space digital filters [11]–[15]. Recently, the  $L_2$ -sensitivity minimization problem has also been treated in the 2-D case [16]–[18].

This paper addresses the  $L_2$ -sensitivity minimization problem for the cases of linear discrete-time systems, linear continuous-time systems, and 2-D state-space digital filters. To this end, we introduce the concept of general controllability and observability Gramians, and derive a novel expression for evaluating the  $L_2$ -sensitivity for each case. This expression includes the conventional  $L_1/L_2$ -sensitivity measure as a special case. Next, an iterative procedure for the minimization of an  $L_2$ -sensitivity measure with respect to a positive-definite and symmetric matrix is presented in each case. This iteration procedure can be performed analytically. Three numerical examples are presented to demonstrate that the proposed algorithms offer faster convergence as well as improved solutions.

Throughout this paper,  $\mathbf{I}$  denotes the identity matrix of appropriate dimension. The transpose (conjugate transpose) of a matrix  $\mathbf{A}$  is indicated by  $\mathbf{A}^T$  ( $\mathbf{A}^*$ ), and  $\text{tr}[\mathbf{A}]$  and  $\oplus$  are used to denote the trace of a square matrix  $\mathbf{A}$  and the direct sum of matrices, respectively.

## II. $L_2$ -SENSITIVITY MINIMIZATION OF LINEAR SYSTEMS

### A. Linear Discrete-Time Systems

Let  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_n$  be a state-space description of a stable, controllable, and observable linear discrete-time system, i.e.,

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

where  $\mathbf{x}(k)$  is an  $n \times 1$  state-variable vector,  $u(k)$  is a scalar input,  $y(k)$  is a scalar output, and  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ , and  $d$  are real constant matrices of appropriate dimensions. The transfer function of (1) is given by

$$H(z) = \mathbf{c}(z\mathbf{I} - \mathbf{A})^{-1}\mathbf{b} + d. \quad (2)$$

Suppose that (1) is implemented by finite-word-length (FWL) fixed-point arithmetic with a  $B$ -bit fractional representation and is realized by  $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}}, \tilde{\mathbf{c}}, \tilde{d})_n$  with

$$\begin{aligned} \tilde{\mathbf{A}} &= \mathbf{A} + \Delta\mathbf{A} & \tilde{\mathbf{b}} &= \mathbf{b} + \Delta\mathbf{b} \\ \tilde{\mathbf{c}} &= \mathbf{c} + \Delta\mathbf{c} & \tilde{d} &= d + \Delta d \end{aligned} \quad (3)$$

Manuscript received May 21, 2001; revised April 15, 2002. This paper was recommended by Associate Editor R. W. Newcomb.

T. Hinamoto, S. Yokoyama, T. Inoue, and W. Zeng are with the Graduate School of Engineering, Hiroshima University, Higashi-Hiroshima 739-8527, Japan.

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

Publisher Item Identifier 10.1109/TCSI.2002.802362.

where  $\Delta\mathbf{A}$ ,  $\Delta\mathbf{b}$ ,  $\Delta\mathbf{c}$ , and  $\Delta d$  stand for the quantization errors of the coefficient matrices. The transfer function of the FWL realization is then expressed in the form

$$\tilde{H}(z) = \tilde{\mathbf{c}}(z\mathbf{I} - \tilde{\mathbf{A}})^{-1}\tilde{\mathbf{b}} + \tilde{d}. \quad (4)$$

Let  $\{p_i\}$  be the set of the ideal parameters of a realization and let  $\{\tilde{p}_i\}$  be its FWL version where  $\tilde{p}_i = p_i + \Delta p_i$ , and  $\Delta p_i$  indicates the corresponding parameter perturbation. Assuming that this realization has  $N$  parameters, the first-order approximation of Taylor's series expansion yields

$$\begin{aligned} \Delta H(z) &= \tilde{H}(z) - H(z) \\ &= \sum_{i=1}^N \frac{\partial H(z)}{\partial p_i} \Delta p_i. \end{aligned} \quad (5)$$

It is obvious that the smaller  $\partial H(z)/\partial p_i$ ,  $i = 1, 2, \dots, N$  yield the smaller transfer-function error  $\Delta H(z)$ . For a fixed-point implementation of  $B$  bits, the parameter perturbations are considered independent random variables uniformly distributed within the range  $[-2^{-B-1}, 2^{-B-1}]$ . Then a measure of the transfer function error can statistically be defined by

$$\begin{aligned} \sigma_{\Delta H}^2 &= \frac{1}{2\pi j} \oint_{|z|=1} E[|\Delta H(z)|^2] \frac{dz}{z} \\ &= \frac{1}{2\pi} \int_0^{2\pi} E[|\Delta H(e^{j\omega})|^2] d\omega \end{aligned} \quad (6)$$

where  $E(\cdot)$  denotes the ensemble average operation. Since  $\{\Delta p_i\}$  are independent random variables uniformly distributed, it follows that

$$E[|\Delta H(z)|^2] = \sum_{i=1}^N \left| \frac{\partial H(z)}{\partial p_i} \right|^2 \sigma^2 \quad (7)$$

where

$$\sigma^2 = E[(\Delta p_i)^2] = \frac{1}{12} 2^{-2B}.$$

By carrying out a coordinate transformation

$$\bar{\mathbf{x}}(k) = \mathbf{T}^{-1} \mathbf{x}(k) \quad (8)$$

to (1), a new realization  $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_n$  characterized by

$$\bar{\mathbf{A}} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c} \mathbf{T} \quad (9)$$

can be derived. From (2) and (9), it is clear that the transfer function  $H(z)$  is invariant under the coordinate transformation of (8).

*Definition 1:* Let  $\mathbf{X}$  be an  $m \times n$  real matrix and let  $f(\mathbf{X})$  be a scalar complex function of  $\mathbf{X}$ , differentiable with respect to all the entries of  $\mathbf{X}$ . The sensitivity function of  $f$  with respect to  $\mathbf{X}$  is then defined as

$$\mathbf{S}_{\mathbf{X}} = \frac{\partial f}{\partial \mathbf{X}}, \quad (\mathbf{S}_{\mathbf{X}})_{ij} = \frac{\partial f}{\partial x_{ij}} \quad (10)$$

where  $x_{ij}$  denotes the  $(i, j)$ th entry of the matrix  $\mathbf{X}$ .

*Definition 2:* Let  $\mathbf{X}(z)$  be an  $m \times n$  complex matrix-valued function of a complex variable  $z$  and let  $x_{pq}(z)$  be the  $(p, q)$ th

entry of  $\mathbf{X}(z)$ . The  $L_2$ -norm of  $\mathbf{X}(z)$  is then defined as [19, p. 48]

$$\begin{aligned} \|\mathbf{X}(z)\|_2 &= \left[ \frac{1}{2\pi} \int_0^{2\pi} \left( \sum_{p=1}^m \sum_{q=1}^n |x_{pq}(e^{j\omega})|^2 \right) d\omega \right]^{1/2} \\ &= \left( \text{tr} \left[ \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{X}(z) \mathbf{X}^*(z) \frac{dz}{z} \right] \right)^{1/2}. \end{aligned} \quad (11)$$

From the foregoing arguments, the overall  $L_2$ -sensitivity measure is defined by

$$\begin{aligned} S_d &= \frac{\sigma_{\Delta H}^2}{\sigma^2} \\ &= \left\| \frac{\partial H(z)}{\partial \mathbf{A}} \right\|_2^2 + \left\| \frac{\partial H(z)}{\partial \mathbf{b}} \right\|_2^2 + \left\| \frac{\partial H(z)}{\partial \mathbf{c}^T} \right\|_2^2 + \left\| \frac{\partial H(z)}{\partial d} \right\|_2^2 \\ &= \left\| \mathbf{G}^T(z) \mathbf{F}^T(z) \right\|_2^2 + \left\| \mathbf{G}^T(z) \right\|_2^2 + \left\| \mathbf{F}(z) \right\|_2^2 + 1 \end{aligned} \quad (12)$$

where

$$\mathbf{F}(z) = (z\mathbf{I} - \mathbf{A})^{-1} \mathbf{b}, \quad \mathbf{G}(z) = \mathbf{c}(z\mathbf{I} - \mathbf{A})^{-1}.$$

The terminology " $L_2$ -sensitivity" used here reflects the fact that the terms involved in (12) are all  $L_2$ -norms, and the definition differs itself from the mixed  $L_1/L_2$  sensitivity defined by

$$\tilde{S}_d = \left\| \frac{\partial H(z)}{\partial \mathbf{A}} \right\|_1 + \left\| \frac{\partial H(z)}{\partial \mathbf{b}} \right\|_2 + \left\| \frac{\partial H(z)}{\partial \mathbf{c}^T} \right\|_2 + \left\| \frac{\partial H(z)}{\partial d} \right\|_2 \quad (13a)$$

whose upper bound, i.e.,

$$\hat{S}_d = \left\| \mathbf{G}^T(z) \right\|_2^2 + \left\| \mathbf{F}(z) \right\|_2^2 + \left\| \mathbf{G}^T(z) \right\|_2^2 + \left\| \mathbf{F}(z) \right\|_2^2 + 1 \quad (13b)$$

was the sensitivity measure employed by the authors of [1]–[5]. From a technical point of view, the main motivation of investigating the  $L_1/L_2$  sensitivity was to overcome the difficulties introduced if the  $L_1$ -norm term in  $\tilde{S}_d$  would have been replaced by a  $L_2$ -norm term. See [6]–[9] for some detailed accounts of this and other issues concerning the relationship between these two sensitivity measures.

Note that the  $L_2$ -sensitivity measure, (12), can also be expressed as

$$S_d = \text{tr}[\mathbf{M}] + \text{tr}[\mathbf{W}_0] + \text{tr}[\mathbf{K}_0] + 1 \quad (14a)$$

where

$$\mathbf{M} = \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{F}(z^{-1}) \mathbf{G}(z^{-1}) [\mathbf{F}(z) \mathbf{G}(z)]^T \frac{dz}{z} \quad (14b)$$

$$\mathbf{K}_0 = \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{F}(z) \mathbf{F}^T(z^{-1}) \frac{dz}{z} \quad (14c)$$

$$\mathbf{W}_0 = \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{G}^T(z) \mathbf{G}(z^{-1}) \frac{dz}{z}. \quad (14d)$$

The matrices  $\mathbf{K}_0$  and  $\mathbf{W}_0$  are called the controllability and observability Gramians, respectively, and can be obtained by solving the Lyapunov equations [20]

$$\mathbf{K}_0 = \mathbf{A} \mathbf{K}_0 \mathbf{A}^T + \mathbf{b} \mathbf{b}^T \quad (15a)$$

$$\mathbf{W}_0 = \mathbf{A}^T \mathbf{W}_0 \mathbf{A} + \mathbf{c}^T \mathbf{c}. \quad (15b)$$

*Definition 3:* The general controllability and observability Gramians for a linear discrete-time system (1) are defined by

$$\mathbf{K}_i = \mathbf{A}\mathbf{K}_i\mathbf{A}^T + \frac{1}{2} \left[ \mathbf{A}^i \mathbf{b}\mathbf{b}^T + \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^i \right] \quad (16a)$$

$$\mathbf{W}_i = \mathbf{A}^T \mathbf{W}_i \mathbf{A} + \frac{1}{2} \left[ \mathbf{c}^T \mathbf{c} \mathbf{A}^i + (\mathbf{A}^T)^i \mathbf{c}^T \mathbf{c} \right] \quad (16b)$$

respectively, where  $i = 0, 1, 2, \dots$

Notice that, with  $i = 0$ , (16a) and (16b) become (15a) and (15b), respectively. In other words, the conventional controllability and observability Gramians,  $\mathbf{K}_0$  and  $\mathbf{W}_0$ , can be viewed as the special cases of the general controllability and observability Gramians,  $\mathbf{K}_i$  and  $\mathbf{W}_i$ .

The theorem below relates the evaluation of the matrix  $\mathbf{M}$  to the general controllability and observability Gramians.

*Theorem 1:* The matrix  $\mathbf{M}$  defined in (14b) is symmetric and positive-definite and can be expressed in the form

$$\mathbf{M} = \text{tr}[\mathbf{W}_0]\mathbf{K}_0 + 2 \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i]\mathbf{K}_i. \quad (17)$$

*Proof:* The symmetry of matrix  $\mathbf{M}$  can be verified by evaluating  $\mathbf{M}^T$  using (14b) as

$$\begin{aligned} \mathbf{M}^T &= \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{F}(z)\mathbf{G}(z)[\mathbf{F}(z^{-1})\mathbf{G}(z^{-1})]^T \frac{dz}{z} \\ &= \frac{1}{2\pi} \int_0^{2\pi} \mathbf{F}(e^{j\omega})\mathbf{G}(e^{j\omega})[\mathbf{F}(e^{-j\omega})\mathbf{G}(e^{-j\omega})]^T d\omega \\ &= \frac{1}{2\pi} \int_0^{2\pi} \mathbf{F}(e^{-j\omega})\mathbf{G}(e^{-j\omega})[\mathbf{F}(e^{j\omega})\mathbf{G}(e^{j\omega})]^T d\omega \\ &= \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{F}(z^{-1})\mathbf{G}(z^{-1})[\mathbf{F}(z)\mathbf{G}(z)]^T \frac{dz}{z} \\ &= \mathbf{M}. \end{aligned} \quad (18)$$

To show the positive definiteness of  $\mathbf{M}$ , let  $\mathbf{x}$  be an arbitrary nonzero column vector of dimension  $n$  and use (14b) to compute

$$\begin{aligned} \mathbf{x}^T \mathbf{M} \mathbf{x} &= \frac{1}{2\pi j} \oint_{|z|=1} \mathbf{x}^T \mathbf{F}(z^{-1})\mathbf{G}(z^{-1})[\mathbf{F}(z)\mathbf{G}(z)]^T \mathbf{x} \frac{dz}{z} \\ &= \frac{1}{2\pi j} \oint_{|z|=1} |[\mathbf{F}(z)\mathbf{G}(z)]^T \mathbf{x}|^2 \frac{dz}{z} \\ &= \frac{1}{2\pi} \int_0^{2\pi} |[\mathbf{F}(e^{j\omega})\mathbf{G}(e^{j\omega})]^T \mathbf{x}|^2 d\omega > 0. \end{aligned} \quad (19)$$

Next, noting that

$$\begin{aligned} (z\mathbf{I} - \mathbf{A})^{-1} &= z^{-1}\mathbf{I} + z^{-2}\mathbf{A} + z^{-3}\mathbf{A}^2 + \dots \\ &= \sum_{k=0}^{\infty} z^{-(k+1)}\mathbf{A}^k \end{aligned} \quad (20)$$

and utilizing the *Cauchy integral theorem*

$$\frac{1}{2\pi j} \oint_C z^k \frac{dz}{z} = \begin{cases} 1, & k = 0 \\ 0, & k \neq 0 \end{cases} \quad (21)$$

where  $C$  is a counterclockwise contour that encircles the origin,

it is possible to write the  $\mathbf{M}$  matrix as

$$\begin{aligned} \mathbf{M} &= \frac{1}{2\pi j} \oint_{|z|=1} \sum_{k=0}^{\infty} \sum_{l=0}^{\infty} \sum_{i=0}^{\infty} \sum_{p=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^i \\ &\quad \cdot \mathbf{c}^T \mathbf{b}^T (\mathbf{A}^T)^p z^{k+l-i-p} \frac{dz}{z} \\ &= \sum_{k=0}^{\infty} \sum_{l=0}^{\infty} \sum_{i=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^i \mathbf{c}^T \mathbf{b}^T (\mathbf{A}^T)^{k+l-i} \end{aligned} \quad (22)$$

where  $p = k + l - i \geq 0$ . Since  $\mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^i \mathbf{c}^T$  is a scalar, it follows that

$$\begin{aligned} \mathbf{M} &= \sum_{l=0}^{\infty} \sum_{i=0}^{\infty} \mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^i \mathbf{c}^T \sum_{k=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^{k+l-i} \\ &= \sum_{l=0}^{\infty} \mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^l \mathbf{c}^T \sum_{k=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^k \\ &\quad + \sum_{l=1}^{\infty} \mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^{l-1} \mathbf{c}^T \sum_{k=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^{k+1} \\ &\quad + \sum_{l=2}^{\infty} \mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^{l-2} \mathbf{c}^T \sum_{k=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^{k+2} \\ &\quad + \dots \\ &\quad + \sum_{i=1}^{\infty} \mathbf{c}\mathbf{A}^{i-1} (\mathbf{A}^T)^i \mathbf{c}^T \sum_{k=1}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^{k-1} \\ &\quad + \sum_{i=2}^{\infty} \mathbf{c}\mathbf{A}^{i-2} (\mathbf{A}^T)^i \mathbf{c}^T \sum_{k=2}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^{k-2} \\ &\quad + \dots \end{aligned} \quad (23)$$

Hence

$$\begin{aligned} \mathbf{M} &= \sum_{l=0}^{\infty} \mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^l \mathbf{c}^T \sum_{k=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^k \\ &\quad + \frac{1}{2} \sum_{l=0}^{\infty} \left[ \mathbf{c}\mathbf{A}\mathbf{A}^l (\mathbf{A}^T)^l \mathbf{c}^T + \mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^l \mathbf{A}^T \mathbf{c}^T \right] \\ &\quad \cdot \sum_{k=0}^{\infty} \left[ \mathbf{A}^k \mathbf{b}\mathbf{b}^T \mathbf{A}^T (\mathbf{A}^T)^k + \mathbf{A}^k \mathbf{A} \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^k \right] \\ &\quad + \frac{1}{2} \sum_{l=0}^{\infty} \left[ \mathbf{c}\mathbf{A}^2 \mathbf{A}^l (\mathbf{A}^T)^l \mathbf{c}^T + \mathbf{c}\mathbf{A}^l (\mathbf{A}^T)^l (\mathbf{A}^T)^2 \mathbf{c}^T \right] \\ &\quad \cdot \sum_{k=0}^{\infty} \left[ \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^2 (\mathbf{A}^T)^k + \mathbf{A}^k \mathbf{A}^2 \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^k \right] \\ &\quad + \dots \\ &= \text{tr} \left[ \sum_{l=0}^{\infty} (\mathbf{A}^T)^l \mathbf{c}^T \mathbf{c}\mathbf{A}^l \right] \sum_{k=0}^{\infty} \mathbf{A}^k \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^k \\ &\quad + \frac{1}{2} \text{tr} \left[ \sum_{l=0}^{\infty} (\mathbf{A}^T)^l (\mathbf{c}^T \mathbf{c}\mathbf{A} + \mathbf{A}^T \mathbf{c}^T \mathbf{c}) \mathbf{A}^l \right] \\ &\quad \cdot \sum_{k=0}^{\infty} \mathbf{A}^k \left( \mathbf{b}\mathbf{b}^T \mathbf{A}^T + \mathbf{A} \mathbf{b}\mathbf{b}^T \right) (\mathbf{A}^T)^k \\ &\quad + \frac{1}{2} \text{tr} \left[ \sum_{l=0}^{\infty} (\mathbf{A}^T)^l (\mathbf{c}^T \mathbf{c}\mathbf{A}^2 + (\mathbf{A}^T)^2 \mathbf{c}^T \mathbf{c}) \mathbf{A}^l \right] \\ &\quad \cdot \sum_{k=0}^{\infty} \mathbf{A}^k \left( \mathbf{b}\mathbf{b}^T (\mathbf{A}^T)^2 + \mathbf{A}^2 \mathbf{b}\mathbf{b}^T \right) (\mathbf{A}^T)^k \\ &\quad + \dots \end{aligned} \quad (24)$$

which coincides with (17). This completes the proof of Theorem 1.

Making use of (17), we can write (14a) as

$$S_d = (\text{tr}[\mathbf{W}_0] + 1)(\text{tr}[\mathbf{K}_0] + 1) + 2 \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i] \text{tr}[\mathbf{K}_i]. \quad (25)$$

This is a new formula for expressing the  $L_2$ -sensitivity of (12) in terms of the general controllability and observability Gramians. Once  $\mathbf{K}_i$  and  $\mathbf{W}_i$  are obtained, the computation quantity of each term in the right-hand side of (25) is constant independent of  $i$ . Also, the infinite sum can be approximated with the finite sum by truncation provided that the system is stable. It should be noted that the first term in the right-hand side of (25) is identical to the upper bound of the  $L_1/L_2$ -mixed sensitivity measure [1]–[5]. Therefore, the new formula for the  $L_2$ -sensitivity reveals its connection to the  $L_1/L_2$ -mixed sensitivity measure in a natural way.

Next, we consider the problem of synthesizing the optimal state-space model with minimum  $L_2$ -sensitivity.

From (16a), (16b), and (9), the general controllability and observability Gramians  $\bar{\mathbf{K}}_i$  and  $\bar{\mathbf{W}}_i$  in the new realization  $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_n$  are related to the original ones by

$$\bar{\mathbf{K}}_i = \mathbf{T}^{-1} \mathbf{K}_i \mathbf{T}^{-T}, \quad \bar{\mathbf{W}}_i = \mathbf{T}^T \mathbf{W}_i \mathbf{T} \quad (26)$$

where  $i = 0, 1, 2, \dots$ . For the new realization  $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_n$ , (25) is changed to

$$S_d(\mathbf{T}\mathbf{T}^T) = (\text{tr}[\bar{\mathbf{W}}_0] + 1)(\text{tr}[\bar{\mathbf{K}}_0] + 1) + 2 \sum_{i=1}^{\infty} \text{tr}[\bar{\mathbf{W}}_i] \text{tr}[\bar{\mathbf{K}}_i] \quad (27)$$

or equivalently

$$S_d(\mathbf{P}) = (\text{tr}[\mathbf{W}_0\mathbf{P}] + 1)(\text{tr}[\mathbf{K}_0\mathbf{P}^{-1}] + 1) + 2 \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i\mathbf{P}] \text{tr}[\mathbf{K}_i\mathbf{P}^{-1}] \quad (28)$$

where  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$  is an  $n \times n$  positive-definite symmetric matrix.

Note that in (27) and (28) the dependence of the  $L_2$ -sensitivity measure on the coordinate transformation matrix  $\mathbf{T}$  (equivalently  $\mathbf{P}$ ) has been explicitly indicated. In this regard, it can be readily verified that the  $L_2$ -sensitivity measure of the original realization as defined in (12) [and evaluated in (14a) and (25)] can be referred to as  $S_d(\mathbf{I})$  where  $\mathbf{I}$  is the identity matrix of dimension  $n$ .

Using the formula for evaluating matrix gradient [21, p. 275]

$$\begin{aligned} \frac{\partial[\text{tr}(\mathbf{M}\mathbf{X})]}{\partial \mathbf{X}} &= \mathbf{M}^T \\ \frac{\partial[\text{tr}(\mathbf{M}\mathbf{X}^{-1})]}{\partial \mathbf{X}} &= -(\mathbf{X}^{-1}\mathbf{M}\mathbf{X}^{-1})^T \end{aligned} \quad (29)$$

it follows from (28) that

$$\begin{aligned} \frac{\partial S_d(\mathbf{P})}{\partial \mathbf{P}} &= (\text{tr}[\mathbf{K}_0\mathbf{P}^{-1}] + 1)\mathbf{W}_0 \\ &\quad - (\text{tr}[\mathbf{W}_0\mathbf{P}] + 1)\mathbf{P}^{-1}\mathbf{K}_0\mathbf{P}^{-1} \\ &\quad + 2 \sum_{i=1}^{\infty} (\text{tr}[\mathbf{K}_i\mathbf{P}^{-1}]\mathbf{W}_i \\ &\quad \quad - \text{tr}[\mathbf{W}_i\mathbf{P}]\mathbf{P}^{-1}\mathbf{K}_i\mathbf{P}^{-1}). \end{aligned} \quad (30)$$

Letting  $\partial S_d(\mathbf{P})/\partial \mathbf{P} = \mathbf{0}$  provides

$$\mathbf{P}\mathbf{G}(\mathbf{P})\mathbf{P} = \mathbf{H}(\mathbf{P}) \quad (31a)$$

where

$$\mathbf{H}(\mathbf{P}) = (\text{tr}[\mathbf{W}_0\mathbf{P}] + 1)\mathbf{K}_0 + 2 \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i\mathbf{P}]\mathbf{K}_i \quad (31b)$$

$$\begin{aligned} \mathbf{G}(\mathbf{P}) &= (\text{tr}[\mathbf{K}_0\mathbf{P}^{-1}] + 1)\mathbf{W}_0 \\ &\quad + 2 \sum_{i=1}^{\infty} \text{tr}[\mathbf{K}_i\mathbf{P}^{-1}]\mathbf{W}_i. \end{aligned} \quad (31c)$$

The matrix  $\mathbf{G}(\mathbf{P})$  in (31c) can be expressed as

$$\begin{aligned} \mathbf{G}(\mathbf{P}) &= \text{tr}[\mathbf{K}_0\mathbf{P}^{-1}]\mathbf{W}_0 + 2 \sum_{i=1}^{\infty} \text{tr}[\mathbf{K}_i\mathbf{P}^{-1}]\mathbf{W}_i + \mathbf{W}_0 \\ &= \text{tr}[\bar{\mathbf{K}}_0]\mathbf{W}_0 + 2 \sum_{i=1}^{\infty} \text{tr}[\bar{\mathbf{K}}_i]\mathbf{W}_i + \mathbf{W}_0. \end{aligned} \quad (32)$$

Using an argument similar to that of the proof for Theorem 1, it can be shown that the sum of the first two terms in the above expression gives a symmetric and positive-definite matrix. Since the third term, matrix  $\mathbf{W}_0$  [see (14d)], is also symmetric and positive-definite, matrix  $\mathbf{G}(\mathbf{P})$  is positive-definite.

It is noted that the sensitivity measure  $S_d(\mathbf{P})$  here is essentially equivalent to a function  $R(\mathbf{P})$  defined in [8, p. 104] or  $f(\mathbf{P})$  defined in [17, p. 811], which has been shown to be unimodal whose stationary point is unique, symmetric, and positive-definite and is the global minimum of the function. Consequently, in what follows, our attention will be focused on solving (31a).

Note that  $\mathbf{P}\mathbf{W}\mathbf{P} = \mathbf{M}$  has the unique solution [5] as

$$\mathbf{P} = \mathbf{W}^{-(1/2)}[\mathbf{W}^{1/2}\mathbf{M}\mathbf{W}^{1/2}]^{1/2}\mathbf{W}^{-(1/2)} \quad (33)$$

in the case when  $\mathbf{W} > 0$  and  $\mathbf{M} \geq 0$  are symmetric. A natural iteration scheme suggested by (31a) is

$$\mathbf{P}_{i+1}\mathbf{G}(\mathbf{P}_i)\mathbf{P}_{i+1} = \mathbf{H}(\mathbf{P}_i). \quad (34)$$

By (33), the matrix  $\mathbf{P}_{i+1}$  satisfying (34) can be obtained as

$$\begin{aligned} \mathbf{P}_{i+1} &= \mathbf{G}^{-(1/2)}(\mathbf{P}_i) \left[ \mathbf{G}^{1/2}(\mathbf{P}_i) \mathbf{H}(\mathbf{P}_i) \right. \\ &\quad \left. \cdot \mathbf{G}^{1/2}(\mathbf{P}_i) \right]^{1/2} \mathbf{G}^{-(1/2)}(\mathbf{P}_i). \end{aligned} \quad (35)$$

This iteration process continues until

$$|S_d(\mathbf{P}_{i+1}) - S_d(\mathbf{P}_i)| < \varepsilon \quad (36)$$

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

To obtain a reasonable initial estimate  $\mathbf{P}_0$  analytically, we denote the principal term of  $S_d(\mathbf{P})$  in (28) as

$$S_d^{(0)}(\mathbf{P}) = (\text{tr}[\mathbf{W}_0\mathbf{P}] + 1)(\text{tr}[\mathbf{K}_0\mathbf{P}^{-1}] + 1). \quad (37)$$

Note that  $S_d^{(0)}(\mathbf{P})$  in (37) corresponds to the upper bound of  $L_1/L_2$ -mixed sensitivity measure [1]–[5]. The positive-definite symmetric matrix  $\mathbf{P}$ , say  $\mathbf{P}_0$ , minimizing (37) is then derived uniquely as

$$\mathbf{P}_0 = \mathbf{W}_0^{-(1/2)} \left[ \mathbf{W}_0^{1/2} \mathbf{K}_0 \mathbf{W}_0^{1/2} \right]^{1/2} \mathbf{W}_0^{-(1/2)} \quad (38)$$

which serves as the initial estimate in the iteration process (35).

Once the optimal positive-definite symmetric matrix  $\mathbf{P}$  that minimizes (28) is obtained, the optimal coordinate transformation matrix  $\mathbf{T}$  is constructed as

$$\mathbf{T} = \mathbf{P}^{1/2}\mathbf{U} \quad (39)$$

where  $U$  is any  $n \times n$  orthogonal matrix. The optimal state-space model with minimum  $L_2$ -sensitivity can then be realized by substituting (39) into (9).

### B. Linear Continuous-Time Systems

In a linear continuous-time system, the transfer function is given by

$$H(s) = \mathbf{c}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{b} + d \quad (40)$$

instead of (2). Equation (14a) is then written in the form

$$S_c = \text{tr}[\mathbf{M}] + \text{tr}[\mathbf{W}_0] + \text{tr}[\mathbf{K}_0] + 1 \quad (41)$$

with

$$\begin{aligned} \mathbf{M} &= \frac{1}{2\pi j} \int_{-j\infty}^{j\infty} \mathbf{F}(-s)\mathbf{G}(-s)[\mathbf{F}(s)\mathbf{G}(s)]^T ds \\ \mathbf{K}_0 &= \int_0^\infty e^{\mathbf{A}t}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T t} dt \\ \mathbf{W}_0 &= \int_0^\infty e^{\mathbf{A}^T t}\mathbf{c}^T \mathbf{c} e^{\mathbf{A}t} dt \\ \mathbf{F}(s) &= (s\mathbf{I} - \mathbf{A})^{-1}\mathbf{b} \\ \mathbf{G}(s) &= \mathbf{c}(s\mathbf{I} - \mathbf{A})^{-1} \end{aligned}$$

where the controllability and observability Gramians  $\mathbf{K}_0$  and  $\mathbf{W}_0$  can be obtained by solving the Lyapunov equations [20]

$$\mathbf{A}\mathbf{K}_0 + \mathbf{K}_0\mathbf{A}^T = -\mathbf{b}\mathbf{b}^T \quad (42a)$$

$$\mathbf{A}^T\mathbf{W}_0 + \mathbf{W}_0\mathbf{A} = -\mathbf{c}^T\mathbf{c} \quad (42b)$$

in place of (15a) and (15b).

*Definition 4:* The general controllability and observability Gramians for a linear continuous-time system are defined by

$$\mathbf{A}\mathbf{K}_i + \mathbf{K}_i\mathbf{A}^T = -\frac{1}{2} \left( \mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T i\Delta} + e^{\mathbf{A}i\Delta} \mathbf{b}\mathbf{b}^T \right) \quad (43a)$$

$$\mathbf{A}^T\mathbf{W}_i + \mathbf{W}_i\mathbf{A} = -\frac{1}{2} \left( \mathbf{c}^T \mathbf{c} e^{\mathbf{A}i\Delta} + e^{\mathbf{A}^T i\Delta} \mathbf{c}^T \mathbf{c} \right) \quad (43b)$$

where  $i = 0, 1, 2, \dots$ , and  $\Delta$  stands for a sampling interval.

*Theorem 2:* The  $\mathbf{M}$  matrix defined in (41) can be approximated by

$$\mathbf{M} \simeq \text{tr}[\mathbf{W}_0]\mathbf{K}_0 + \sqrt{2}\Delta \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i]\mathbf{K}_i. \quad (44)$$

*Proof:* Applying the inverse Laplace transform yields

$$\mathcal{L}^{-1}[\mathbf{F}(s)\mathbf{G}(s)] = \int_0^\infty e^{\mathbf{A}(t-\tau)}\mathbf{b}\mathbf{c}e^{\mathbf{A}^T\tau} d\tau \quad (45a)$$

$$\mathcal{L}^{-1}[\mathbf{F}(s)\mathbf{G}(s)]^T = \int_0^\infty e^{\mathbf{A}^T\tau}\mathbf{c}^T\mathbf{b}^T e^{\mathbf{A}(t-\tau)} d\tau \quad (45b)$$

where [20, p. 67]

$$\mathbf{e}^{\mathbf{A}t} = \begin{cases} \sum_{k=0}^{\infty} \frac{1}{k!} \mathbf{A}^k t^k, & t \geq 0 \\ \mathbf{0}, & t < 0. \end{cases}$$

Using Parseval's theorem, we can write the matrix  $\mathbf{M}$  in (41) as

$$\begin{aligned} \mathbf{M} &= \int_0^\infty \int_0^\infty \int_0^\infty e^{\mathbf{A}(t-\lambda)}\mathbf{b}\mathbf{c}e^{\mathbf{A}^T\lambda} \\ &\quad \cdot e^{\mathbf{A}^T\tau}\mathbf{c}^T\mathbf{b}^T e^{\mathbf{A}(t-\tau)} d\lambda d\tau dt \\ &= \int_0^\infty \int_0^\infty \mathbf{c}e^{\mathbf{A}^T\tau}\mathbf{c}^T \\ &\quad \cdot \left[ \int_0^\infty e^{\mathbf{A}(t-\lambda)}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T(t-\tau)} dt \right] d\lambda d\tau \quad (46) \end{aligned}$$

which can be approximated by

$$\begin{aligned} \mathbf{M} &\simeq \int_0^\infty \mathbf{c}e^{\mathbf{A}^T\lambda}\mathbf{c}^T \\ &\quad \cdot \left[ \int_0^\infty e^{\mathbf{A}(t-\lambda)}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T(t-\lambda)} dt \right] d\lambda \\ &\quad + \frac{\Delta}{\sqrt{2}} \sum_{i=1}^{\infty} \int_0^\infty \mathbf{c}e^{\mathbf{A}^T(\lambda+i\Delta)}\mathbf{c}^T \\ &\quad \cdot \left[ \int_0^\infty e^{\mathbf{A}(t-\lambda)}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T(t-\lambda-i\Delta)} dt \right] d\lambda \\ &\quad + \frac{\Delta}{\sqrt{2}} \sum_{i=1}^{\infty} \int_0^\infty \mathbf{c}e^{\mathbf{A}^T\tau}\mathbf{c}^T \\ &\quad \cdot \left[ \int_0^\infty e^{\mathbf{A}(t-\tau-i\Delta)}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T(t-\tau)} dt \right] d\tau \\ &= \int_0^\infty \mathbf{c}e^{\mathbf{A}^T\lambda}\mathbf{c}^T \\ &\quad \cdot \left[ \int_0^\infty e^{\mathbf{A}(t-\lambda)}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T(t-\lambda)} dt \right] d\lambda \\ &\quad + \frac{\Delta}{\sqrt{2}} \sum_{i=1}^{\infty} \int_0^\infty \mathbf{c}e^{\mathbf{A}^T(\lambda+i\Delta)}\mathbf{c}^T \\ &\quad \cdot \left[ \int_0^\infty e^{\mathbf{A}(t+i\Delta)}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T t} dt \right] d\lambda \\ &\quad + \frac{\Delta}{\sqrt{2}} \sum_{i=1}^{\infty} \int_0^\infty \mathbf{c}e^{\mathbf{A}^T\tau}\mathbf{c}^T \\ &\quad \cdot \left[ \int_0^\infty e^{\mathbf{A}t}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T(t+i\Delta)} dt \right] d\tau \quad (47) \end{aligned}$$

where  $\Delta$  denotes a sampling interval. From (47) it follows that

$$\begin{aligned} \mathbf{M} &\simeq \int_0^\infty \mathbf{c}e^{\mathbf{A}^T\lambda}\mathbf{c}^T d\lambda \int_0^\infty e^{\mathbf{A}t}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T t} dt \\ &\quad + \frac{\Delta}{\sqrt{2}} \sum_{i=1}^{\infty} \int_0^\infty \mathbf{c}e^{\mathbf{A}^T\lambda}\mathbf{c}^T e^{\mathbf{A}i\Delta} d\lambda \\ &\quad \cdot \int_0^\infty e^{\mathbf{A}t} e^{\mathbf{A}i\Delta} \mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T t} dt \\ &\quad + \frac{\Delta}{\sqrt{2}} \sum_{i=1}^{\infty} \int_0^\infty \mathbf{c}e^{\mathbf{A}i\Delta} \mathbf{c}^T e^{\mathbf{A}^T\tau} d\tau \\ &\quad \cdot \int_0^\infty e^{\mathbf{A}t} \mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T i\Delta} e^{\mathbf{A}^T t} dt \\ &= \text{tr} \left[ \int_0^\infty e^{\mathbf{A}^T\lambda}\mathbf{c}^T \mathbf{c} e^{\mathbf{A}^T\lambda} d\lambda \right] \int_0^\infty e^{\mathbf{A}t}\mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T t} dt + \frac{\Delta}{2\sqrt{2}} \\ &\quad \cdot \sum_{i=1}^{\infty} \text{tr} \left[ \int_0^\infty e^{\mathbf{A}^T\lambda} \left( e^{\mathbf{A}^T i\Delta} \mathbf{c}^T \mathbf{c} + \mathbf{c}^T \mathbf{c} e^{\mathbf{A}i\Delta} \right) e^{\mathbf{A}^T\lambda} d\lambda \right] \\ &\quad \cdot \int_0^\infty e^{\mathbf{A}t} \left( e^{\mathbf{A}i\Delta} \mathbf{b}\mathbf{b}^T + \mathbf{b}\mathbf{b}^T e^{\mathbf{A}^T i\Delta} \right) e^{\mathbf{A}^T t} dt \quad (48) \end{aligned}$$

which is identical to (44). This completes the proof of Theorem 2.

Making use of (44), one can write (41) as

$$\begin{aligned} S_c &\simeq (\text{tr}[\mathbf{W}_0] + 1)(\text{tr}[\mathbf{K}_0] + 1) \\ &\quad + \sqrt{2}\Delta \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i]\text{tr}[\mathbf{K}_i]. \quad (49) \end{aligned}$$

This can be viewed as the continuous-time version of (25), but in the present case (49) is an approximate expression for evaluating the  $L_2$ -sensitivity of a linear continuous-time system in the state-space description. The infinite sum in (49) can be approximated with a finite sum by truncation provided that the system is stable. Note that the first term in the right-hand side of (49) is the upper bound of the  $L_1/L_2$ -mixed sensitivity measure [10], thus the  $L_2$ -sensitivity expression in (49) offers an analytical relation to the  $L_1/L_2$ -mixed sensitivity measure.

In the realization of the optimal linear continuous-time system with minimum  $L_2$ -sensitivity, (49) is changed to

$$S_c(\mathbf{T}\mathbf{T}^T) \simeq (\text{tr}[\mathbf{T}^T \mathbf{W}_0 \mathbf{T}] + 1)(\text{tr}[\mathbf{T}^{-1} \mathbf{K}_0 \mathbf{T}^{-T}] + 1) + \sqrt{2\Delta} \sum_{i=1}^{\infty} \text{tr}[\mathbf{T}^T \mathbf{W}_i \mathbf{T}] \text{tr}[\mathbf{T}^{-1} \mathbf{K}_i \mathbf{T}^{-T}] \quad (50)$$

or equivalently

$$S_c(\mathbf{P}) \simeq (\text{tr}[\mathbf{W}_0 \mathbf{P}] + 1)(\text{tr}[\mathbf{K}_0 \mathbf{P}^{-1}] + 1) + \sqrt{2\Delta} \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i \mathbf{P}] \text{tr}[\mathbf{K}_i \mathbf{P}^{-1}] \quad (51)$$

where  $\mathbf{P} = \mathbf{T}\mathbf{T}^T$ . Similar to the discrete-time case, the  $L_2$ -sensitivity measure evaluated in (41) and (49) may be referred to as  $S_c(\mathbf{I})$ . It follows that the formulas in (31a), (34), and (35) also hold for the present case provided that  $\mathbf{H}(\mathbf{P})$  and  $\mathbf{G}(\mathbf{P})$  are modified to

$$\mathbf{H}(\mathbf{P}) = (\text{tr}[\mathbf{W}_0 \mathbf{P}] + 1) \mathbf{K}_0 + \sqrt{2\Delta} \sum_{i=1}^{\infty} \text{tr}[\mathbf{W}_i \mathbf{P}] \mathbf{K}_i \quad (52a)$$

$$\mathbf{G}(\mathbf{P}) = (\text{tr}[\mathbf{K}_0 \mathbf{P}^{-1}] + 1) \mathbf{W}_0 + \sqrt{2\Delta} \sum_{i=1}^{\infty} \text{tr}[\mathbf{K}_i \mathbf{P}^{-1}] \mathbf{W}_i. \quad (52b)$$

Consequently, the iteration process described by (34) and (35) and the initial estimate given by (38) can also apply to the linear continuous-time systems.

### III. $L_2$ -SENSITIVITY MINIMIZATION OF 2-D DIGITAL FILTERS

A 2-D digital filter can be represented by the following local state-space (LSS) model  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{m,n}$  which is stable, separately locally controllable, and separately locally observable [22]:

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^h(i+1, j) \\ \mathbf{x}^v(i, j+1) \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix} \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} + \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} u(i, j) \\ &\triangleq \mathbf{A} \mathbf{x}(i, j) + \mathbf{b} u(i, j) \\ y(i, j) &= [\mathbf{c}_1 \quad \mathbf{c}_2] \begin{bmatrix} \mathbf{x}^h(i, j) \\ \mathbf{x}^v(i, j) \end{bmatrix} + d u(i, j) \\ &\triangleq \mathbf{c} \mathbf{x}(i, j) + d u(i, j) \end{aligned} \quad (53)$$

where  $\mathbf{x}^h(i, j)$  is an  $m \times 1$  horizontal state vector,  $\mathbf{x}^v(i, j)$  is an  $n \times 1$  vertical state vector,  $u(i, j)$  is a scalar input,  $y(i, j)$  is a scalar output, and  $\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3, \mathbf{A}_4, \mathbf{b}_1, \mathbf{b}_2, \mathbf{c}_1, \mathbf{c}_2$ , and  $d$  are real constant matrices of appropriate dimensions. The transfer function of (53) is given by

$$H(z_1, z_2) = \mathbf{c}(z_1 \mathbf{I}_m \oplus z_2 \mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b} + d. \quad (54)$$

Let the coordinate transformation be specified by

$$\bar{\mathbf{x}}(i, j) = \mathbf{T}^{-1} \mathbf{x}(i, j) \quad (55)$$

where  $\mathbf{T} = \mathbf{T}_1 \oplus \mathbf{T}_4$  with  $\mathbf{T}_1$  and  $\mathbf{T}_4$  nonsingular. Then the LSS model (53) is equivalent to a new realization  $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{m,n}$  characterized by

$$\bar{\mathbf{A}} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \quad \bar{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b}, \quad \bar{\mathbf{c}} = \mathbf{c} \mathbf{T} \quad (56)$$

in the sense that the transfer function  $H(z_1, z_2)$  remains invariant under such a transformation.

*Definition 5:* Let  $\mathbf{X}(z_1, z_2)$  be an  $m \times n$  complex matrix valued function of the complex variables  $z_1$  and  $z_2$ . The  $L_p$  norm of  $\mathbf{X}(z_1, z_2)$  is then defined as

$$\|\mathbf{X}\|_p = \left[ \frac{1}{(2\pi j)^2} \oint \oint_{\Gamma^2} \|\mathbf{X}(z_1, z_2)\|_F^p \frac{dz_1 dz_2}{z_1 z_2} \right]^{1/p} \quad (57)$$

where  $\|\mathbf{X}(z_1, z_2)\|_F$  is the Frobenius norm of the matrix  $\mathbf{X}(z_1, z_2)$  defined by

$$\|\mathbf{X}(z_1, z_2)\|_F = \left[ \sum_{p=1}^m \sum_{q=1}^n |x_{pq}(z_1, z_2)|^2 \right]^{1/2}$$

and  $\Gamma^2 = \{(z_1, z_2) : |z_1| = 1, |z_2| = 1\}$ . An  $L_2$ -sensitivity is then evaluated by

$$\begin{aligned} S_{2d} &= \left\| \frac{\partial H(z_1, z_2)}{\partial \mathbf{A}} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2)}{\partial \mathbf{b}} \right\|_2^2 \\ &\quad + \left\| \frac{\partial H(z_1, z_2)}{\partial \mathbf{c}^T} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2)}{\partial d} \right\|_2^2 \\ &= \|\mathbf{G}^T(z_1, z_2) \mathbf{F}^T(z_1, z_2)\|_2^2 + \|\mathbf{G}^T(z_1, z_2)\|_2^2 \\ &\quad + \|\mathbf{F}(z_1, z_2)\|_2^2 + 1 \end{aligned} \quad (58)$$

where

$$\begin{aligned} \mathbf{F}(z_1, z_2) &= (z_1 \mathbf{I}_m \oplus z_2 \mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b} \\ \mathbf{G}(z_1, z_2) &= \mathbf{c}(z_1 \mathbf{I}_m \oplus z_2 \mathbf{I}_n - \mathbf{A})^{-1}. \end{aligned}$$

By Definition 5, we can write (58) as

$$S_{2d} = \text{tr}[\mathbf{M}] + \text{tr}[\mathbf{W}_{00}] + \text{tr}[\mathbf{K}_{00}] + 1 \quad (59)$$

where

$$\begin{aligned} \mathbf{M} &= \frac{1}{(2\pi j)^2} \oint \oint_{\Gamma^2} \mathbf{F}(z_1^{-1}, z_2^{-1}) \mathbf{G}(z_1^{-1}, z_2^{-1}) \\ &\quad \cdot [\mathbf{F}(z_1, z_2) \mathbf{G}(z_1, z_2)]^T \frac{dz_1 dz_2}{z_1 z_2} \\ \mathbf{K}_{00} &= \frac{1}{(2\pi j)^2} \oint \oint_{\Gamma^2} \mathbf{F}(z_1, z_2) \mathbf{F}^T(z_1^{-1}, z_2^{-1}) \frac{dz_1 dz_2}{z_1 z_2} \\ \mathbf{W}_{00} &= \frac{1}{(2\pi j)^2} \oint \oint_{\Gamma^2} \mathbf{G}^T(z_1, z_2) \mathbf{G}(z_1^{-1}, z_2^{-1}) \frac{dz_1 dz_2}{z_1 z_2}. \end{aligned}$$

The matrices  $\mathbf{K}_{00}$  and  $\mathbf{W}_{00}$  here are called 2-D local controllability and observability Gramians, respectively, and can be derived from

$$\mathbf{F}(z_1, z_2) = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{f}(i, j) z_1^{-i} z_2^{-j} \quad (60a)$$

$$\mathbf{G}(z_1, z_2) = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{g}^T(i, j) z_1^{-i} z_2^{-j} \quad (60b)$$

where

$$\begin{aligned} \mathbf{f}(i, j) &= \mathbf{A}^{(i-1, j)} \begin{bmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{0} \end{bmatrix} \\ &\quad + \mathbf{A}^{(i, j-1)} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ \mathbf{b}_2 \end{bmatrix} \\ \mathbf{g}^T(i, j) &= \mathbf{c}\mathbf{A}^{(i-1, j)} \begin{bmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \mathbf{c}\mathbf{A}^{(i, j-1)} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \end{aligned}$$

and the state transition matrix  $\mathbf{A}^{(i, j)}$  is derived from

$$\begin{aligned} \mathbf{A}^{(1, 0)} &= \begin{bmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{A} \quad \mathbf{A}^{(0, 1)} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \mathbf{A} \\ \mathbf{A}^{(0, 0)} &= \mathbf{I}_{m+n}, \quad \mathbf{A}^{(-i, j)} = \mathbf{0} \quad (i \geq 1) \\ \mathbf{A}^{(i, -j)} &= \mathbf{0} \quad (j \geq 1) \\ \mathbf{A}^{(i, j)} &= \mathbf{A}^{(1, 0)} \mathbf{A}^{(i-1, j)} + \mathbf{A}^{(0, 1)} \mathbf{A}^{(i, j-1)} \\ &= \mathbf{A}^{(i-1, j)} \mathbf{A}^{(1, 0)} + \mathbf{A}^{(i, j-1)} \mathbf{A}^{(0, 1)}, \\ &\quad (i, j) > (0, 0). \end{aligned}$$

By utilizing the following 2-D Cauchy integral theorem:

$$\frac{1}{(2\pi j)^2} \oint \oint_{C^2} z_1^m z_2^n \frac{dz_1 dz_2}{z_1 z_2} = \begin{cases} 1, & (m, n) = (0, 0) \\ 0, & (m, n) \neq (0, 0) \end{cases} \quad (61)$$

where  $C$  is a counterclockwise contour that encircles the origin, it is possible to write the  $\mathbf{K}_{00}$  and  $\mathbf{W}_{00}$  in the form

$$\mathbf{K}_{00} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{f}(i, j) \mathbf{f}^T(i, j) \quad (62a)$$

$$\mathbf{W}_{00} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{g}(i, j) \mathbf{g}^T(i, j). \quad (62b)$$

*Definition 6:* The general 2-D local controllability and observability Gramians for a 2-D state-space digital filter, (53), are defined by

$$\mathbf{K}_{\xi\zeta} = \frac{1}{2} \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \left[ \mathbf{f}(i, j) \mathbf{f}^T(i + \xi, j + \zeta) + \mathbf{f}(i + \xi, j + \zeta) \mathbf{f}^T(i, j) \right] \quad (63a)$$

$$\mathbf{W}_{\xi\zeta} = \frac{1}{2} \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \left[ \mathbf{g}(i, j) \mathbf{g}^T(i + \xi, j + \zeta) + \mathbf{g}(i + \xi, j + \zeta) \mathbf{g}^T(i, j) \right] \quad (63b)$$

where  $(0, 0) \leq (\xi, \zeta) < (\infty, \infty)$ .

*Theorem 3:* The matrix  $\mathbf{M}$  defined in (59) can be evaluated as

$$\mathbf{M} = \text{tr}[\mathbf{W}_{00}] \mathbf{K}_{00} + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\mathbf{W}_{\xi\zeta}] \mathbf{K}_{\xi\zeta}. \quad (64)$$

*Proof:* This can be proved by substituting (60a) and (60b) into the  $\mathbf{M}$  matrix defined in (59) and utilizing the 2-D Cauchy integral theorem stated in (61).

Substituting (64) into (59) yields

$$\begin{aligned} S_{2d} &= (\text{tr}[\mathbf{W}_{00}] + 1)(\text{tr}[\mathbf{K}_{00}] + 1) \\ &\quad + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\mathbf{W}_{\xi\zeta}] \text{tr}[\mathbf{K}_{\xi\zeta}]. \end{aligned} \quad (65)$$

For the new realization  $(\bar{\mathbf{A}}, \bar{\mathbf{b}}, \bar{\mathbf{c}}, d)_{m,n}$  characterized by (56), it is easy to show that (65) is changed to

$$\begin{aligned} S_{2d}(\mathbf{P}) &= (\text{tr}[\mathbf{W}_{00}\mathbf{P}] + 1)(\text{tr}[\mathbf{K}_{00}\mathbf{P}^{-1}] + 1) \\ &\quad + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\mathbf{W}_{\xi\zeta}\mathbf{P}] \text{tr}[\mathbf{K}_{\xi\zeta}\mathbf{P}^{-1}] \end{aligned} \quad (66)$$

where

$$\mathbf{P} = \mathbf{T}\mathbf{T}^T, \quad \mathbf{P} = \begin{bmatrix} \mathbf{P}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_4 \end{bmatrix}.$$

According to the partition of  $\mathbf{P}$ , the Gramians  $\mathbf{W}_{\xi\zeta}$  and  $\mathbf{K}_{\xi\zeta}$  become

$$\mathbf{W}_{\xi\zeta} = \begin{bmatrix} \mathbf{W}_{\xi\zeta}^{(1)} & \mathbf{W}_{\xi\zeta}^{(2)} \\ \mathbf{W}_{\xi\zeta}^{(3)} & \mathbf{W}_{\xi\zeta}^{(4)} \end{bmatrix}, \quad \mathbf{K}_{\xi\zeta} = \begin{bmatrix} \mathbf{K}_{\xi\zeta}^{(1)} & \mathbf{K}_{\xi\zeta}^{(2)} \\ \mathbf{K}_{\xi\zeta}^{(3)} & \mathbf{K}_{\xi\zeta}^{(4)} \end{bmatrix}. \quad (67)$$

By introducing

$$\hat{\mathbf{W}}_{\xi\zeta} = \begin{bmatrix} \mathbf{W}_{\xi\zeta}^{(1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{W}_{\xi\zeta}^{(4)} \end{bmatrix}, \quad \hat{\mathbf{K}}_{\xi\zeta} = \begin{bmatrix} \mathbf{K}_{\xi\zeta}^{(1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_{\xi\zeta}^{(4)} \end{bmatrix} \quad (68)$$

we can express  $S_{2d}(\mathbf{P})$  in (66) as

$$\begin{aligned} S_{2d}(\mathbf{P}) &= (\text{tr}[\hat{\mathbf{W}}_{00}\mathbf{P}] + 1)(\text{tr}[\hat{\mathbf{K}}_{00}\mathbf{P}^{-1}] + 1) \\ &\quad + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\hat{\mathbf{W}}_{\xi\zeta}\mathbf{P}] \text{tr}[\hat{\mathbf{K}}_{\xi\zeta}\mathbf{P}^{-1}]. \end{aligned} \quad (69)$$

Hence, it suffices to deal with the matrix  $\hat{\mathbf{K}}$  instead of  $\mathbf{K}$ . To make the subsequent derivation simple, in the following we omit the hat and write  $\mathbf{K}$  for  $\hat{\mathbf{K}}$ .

By letting  $\partial S_{2d}(\mathbf{P})/\partial \mathbf{P}_1 = \mathbf{0}$  and  $\partial S_{2d}(\mathbf{P})/\partial \mathbf{P}_4 = \mathbf{0}$ , we obtain

$$\mathbf{P}_1 \mathbf{F}_1(\mathbf{P}) \mathbf{P}_1 = \mathbf{F}_2(\mathbf{P}) \quad (70a)$$

$$\mathbf{P}_4 \mathbf{F}_3(\mathbf{P}) \mathbf{P}_4 = \mathbf{F}_4(\mathbf{P}) \quad (70b)$$

respectively, where

$$\begin{aligned} \mathbf{F}_1(\mathbf{P}) &= \text{tr}[\mathbf{K}_{00}\mathbf{P}^{-1}] \mathbf{W}_{00}^{(1)} \\ &\quad + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\mathbf{K}_{\xi\zeta}\mathbf{P}^{-1}] \mathbf{W}_{\xi\zeta}^{(1)} \end{aligned}$$

$$\begin{aligned} \mathbf{F}_2(\mathbf{P}) &= \text{tr}[\mathbf{W}_{00}\mathbf{P}] \mathbf{K}_{00}^{(1)} \\ &\quad + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\mathbf{W}_{\xi\zeta}\mathbf{P}] \mathbf{K}_{\xi\zeta}^{(1)} \end{aligned}$$

$$\begin{aligned} \mathbf{F}_3(\mathbf{P}) &= \text{tr}[\mathbf{K}_{00}\mathbf{P}^{-1}] \mathbf{W}_{00}^{(4)} \\ &\quad + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\mathbf{K}_{\xi\zeta}\mathbf{P}^{-1}] \mathbf{W}_{\xi\zeta}^{(4)} \end{aligned}$$

$$\begin{aligned} \mathbf{F}_4(\mathbf{P}) &= \text{tr}[\mathbf{W}_{00}\mathbf{P}] \mathbf{K}_{00}^{(4)} \\ &\quad + 2 \sum_{(0,0) < (\xi,\zeta) < (\infty,\infty)} \text{tr}[\mathbf{W}_{\xi\zeta}\mathbf{P}] \mathbf{K}_{\xi\zeta}^{(4)}. \end{aligned}$$

The iterative procedure reported in [18] can be used to solve (70a) and (70b) for  $\mathbf{P}_1$  and  $\mathbf{P}_4$ , that is,

$$\begin{aligned} \mathbf{P}_1^{(i+1)} &= \mathbf{F}_1^{-1/2}(\mathbf{P}^{(i)}) [\mathbf{F}_1^{1/2}(\mathbf{P}^{(i)}) \mathbf{F}_2(\mathbf{P}^{(i)}) \\ &\quad \cdot \mathbf{F}_1^{1/2}(\mathbf{P}^{(i)})]^{1/2} \mathbf{F}_1^{-1/2}(\mathbf{P}^{(i)}) \end{aligned} \quad (71a)$$

$$\begin{aligned} \mathbf{P}_4^{(i+1)} &= \mathbf{F}_3^{-1/2}(\mathbf{P}^{(i)}) [\mathbf{F}_3^{1/2}(\mathbf{P}^{(i)}) \mathbf{F}_4(\mathbf{P}^{(i)}) \\ &\quad \cdot \mathbf{F}_3^{1/2}(\mathbf{P}^{(i)})]^{1/2} \mathbf{F}_3^{-1/2}(\mathbf{P}^{(i)}) \end{aligned} \quad (71b)$$

where  $\mathbf{P}^{(i)}$  is derived from the previous iteration. The iteration continues until

$$|S_{2d}(\mathbf{P}^{(i+1)}) - S_{2d}(\mathbf{P}^{(i)})| < \varepsilon \quad (72)$$

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

The initial estimate  $\mathbf{P}^{(0)} = \mathbf{P}_1^{(0)} \oplus \mathbf{P}_4^{(0)}$  can be obtained by minimizing the upper bound of  $L_1/L_2$ -norms mixed sensitivity measure [11], [14]

$$S_{2d}^{(0)}(\mathbf{P}) = (\text{tr}[\mathbf{W}_{00}\mathbf{P}] + 1)(\text{tr}[\mathbf{K}_{00}\mathbf{P}^{-1}] + 1) \quad (73)$$

which yields a closed-form solution as

$$\mathbf{P}^{(0)} = \mathbf{W}_{00}^{-(1/2)} \left[ \mathbf{W}_{00}^{1/2} \mathbf{K}_{00} \mathbf{W}_{00}^{1/2} \right]^{1/2} \mathbf{W}_{00}^{-(1/2)}. \quad (74)$$

Once the positive-definite  $L_2$ -optimal symmetric matrix  $\mathbf{P} = \mathbf{P}_1 \oplus \mathbf{P}_4$  is obtained, the  $L_2$ -optimal coordinate transformation matrix can be constructed as

$$\mathbf{T} = \left[ \mathbf{P}_1^{1/2} \oplus \mathbf{P}_4^{1/2} \right] \left[ \mathbf{U}_1 \oplus \mathbf{U}_4 \right] \quad (75)$$

where  $\mathbf{U}_1$  and  $\mathbf{U}_4$  are arbitrary  $m \times m$  and  $n \times n$  orthogonal matrices, respectively. The  $L_2$ -optimal filter structures that minimize (69) can be synthesized by substituting (75) into (56).

#### IV. NUMERICAL EXAMPLES

*Example 1:* Consider a linear discrete-time system, (1), specified by

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0.45377 & -1.55616 & 1.97486 \end{bmatrix} \\ \mathbf{b} &= [0 \ 0 \ 1]^T \\ \mathbf{c} &= [0.02317 \ 0.02302 \ 0.07930] \\ d &= 0.01594 \end{aligned}$$

whose poles are at  $z = 0.6578817$  and  $z = 0.6584892 \pm j0.5060989$ .

Using (15a) and (15b), its controllability and observability Gramians are calculated as

$$\begin{aligned} \mathbf{K}_0 &= \begin{bmatrix} 17.06183537 & 14.88646408 & 9.60276753 \\ 14.88646408 & 17.06183537 & 14.88646408 \\ 9.60276753 & 14.88646408 & 17.06183537 \end{bmatrix} \\ \mathbf{W}_0 &= \begin{bmatrix} 0.048103939 & -0.119291362 & 0.095427125 \\ -0.119291362 & 0.311061239 & -0.249967573 \\ 0.095427125 & -0.249967573 & 0.231012260 \end{bmatrix}. \end{aligned}$$

By using (16a), (16b), (25), and (37), the  $L_2$ -sensitivity measure and the first term in it are computed as

$$\begin{aligned} S_d &= 159.89091094 \\ S_d^{(0)}(\mathbf{I}) &= 82.98421445 \end{aligned}$$

where the infinite sum in (25) was truncated with  $i = 100$  (the same truncation will be applied hereafter in this example).

The initial estimate  $\mathbf{P}_0$  is derived from (38) as

$$\mathbf{P}_0 = \begin{bmatrix} 68.33559995 & 39.19845769 & 13.72288921 \\ 39.19845769 & 33.21750517 & 21.48361875 \\ 13.72288921 & 21.48361875 & 22.73726342 \end{bmatrix}$$

TABLE I  
 $L_2$ -SENSITIVITY PROFILE OF THE FIRST TEN ITERATIONS

$i$	$L_2$ -Sensitivity	$i$	$L_2$ -Sensitivity
	159.8909109417472	5	8.832683342813343
0	8.943594607334033	6	8.832683342812404
1	8.833335506929652	7	8.832683342812398
2	8.832687360817378	8	8.832683342812402
3	8.832683367564366	9	8.832683342812402
4	8.832683342964870	10	8.832683342812393

or equivalently

$$\mathbf{T}_0 = \begin{bmatrix} 1.538415708 & 4.703727215 & 6.621467159 \\ 3.039406990 & 4.443971631 & 2.056848672 \\ 4.246846977 & 2.126308653 & -0.424694816 \end{bmatrix}.$$

Using (28) and (37), we evaluate

$$\begin{aligned} S_d(\mathbf{P}_0) &= 8.94359461 \\ S_d^{(0)}(\mathbf{P}_0) &= 5.75547618. \end{aligned}$$

Applying (35) to minimize (28), after ten iterations we obtain

$$\mathbf{P}_{10} = \begin{bmatrix} 81.46913871 & 48.38423145 & 17.98009171 \\ 48.38423145 & 39.03244413 & 24.14333809 \\ 17.98009171 & 24.14333809 & 23.95155207 \end{bmatrix}$$

or equivalently

$$\mathbf{T}_{10} = \begin{bmatrix} 1.283114259 & 5.356474502 & 7.150590005 \\ 3.056235553 & 4.759869047 & 2.652454530 \\ 4.338392561 & 2.264582260 & 0.039613686 \end{bmatrix}.$$

Using (28) and (37), we evaluate

$$\begin{aligned} S_d(\mathbf{P}_{10}) &= 8.83268334 \\ S_d^{(0)}(\mathbf{P}_{10}) &= 5.81485521. \end{aligned}$$

The  $L_2$ -sensitivity profile of the first ten iterations is given in Table I, from which we see that with a tolerance  $\varepsilon = 10^{-7}$  the algorithm converges after four iterations. By substituting  $\mathbf{T} = \mathbf{T}_{10}$  into (9), the optimal state-space model with minimum  $L_2$ -sensitivity is then obtained as

$$\begin{aligned} \bar{\mathbf{A}} &= \begin{bmatrix} 0.796804001 & -0.414267323 & -0.025895398 \\ 0.414267323 & 0.565420531 & -0.316443915 \\ -0.025895398 & 0.316443915 & 0.612635467 \end{bmatrix} \\ \bar{\mathbf{b}} &= \begin{bmatrix} 0.444118830 \\ -0.413263073 \\ 0.229880039 \end{bmatrix} \\ \bar{\mathbf{c}} &= [0.444118830 \ 0.413263073 \ 0.229880039] \\ d &= 0.015940000. \end{aligned}$$

For comparison purposes, the gradient-flow-based algorithm [8] in which  $\mathbf{P}_{i+1}$  is generated by

$$\mathbf{P}_{i+1} = \mathbf{P}_i - \mu \frac{\partial S_d(\mathbf{P}_i)}{\partial \mathbf{P}_i} \quad (76)$$



was applied to the above example with the same initial estimate. In the simulation, the step-size  $\mu$  in (76) is set to 0.28 (which seems the best), and the algorithm after 1000 iterations yields

$$S_d(\mathbf{P}_{1000}) = 8.832\,746\,57.$$

This is slightly larger than

$$S_d(\mathbf{P}_{10}) = 8.832\,683\,34$$

obtained by our method.

*Example 2:* Now we consider a linear continuous-time system, (40), characterized by

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -0.162\,484\,210 & -0.526\,273\,961 & -0.628\,831\,301 \end{bmatrix}$$

$$\mathbf{b} = [0 \ 0 \ 1]^T$$

$$\mathbf{c} = [0.250\,120\,406 \ 0 \ 0.191\,751\,308]$$

$$d = -0.539\,352\,1$$

whose poles are given by  $s = -0.376\,731\,3$  and  $s = -0.126\,050 \pm j0.644\,524\,2$ .

The controllability and observability Gramians  $\mathbf{K}_0$  and  $\mathbf{W}_0$  are computed from (42a) and (42b) as

$$\mathbf{W}_0 = \begin{bmatrix} 0.148\,531\,757 & 0.152\,935\,083 & 0.192\,511\,683 \\ 0.152\,935\,083 & 0.248\,817\,754 & 0.290\,599\,751 \\ 0.192\,511\,683 & 0.290\,599\,752 & 0.491\,362\,363 \end{bmatrix}$$

$$\mathbf{K}_0 = \begin{bmatrix} 11.487\,179\,114 & 0.000\,000\,000 & -2.968\,181\,165 \\ 0.000\,000\,000 & 2.968\,181\,165 & 0.000\,000\,000 \\ -2.968\,181\,165 & 0.000\,000\,000 & 1.562\,076\,459 \end{bmatrix}.$$

By letting  $\Delta = 0.01$  and using (43a) and (43b), the formula (49), where the infinite sum was truncated to its first 8000 terms, gives

$$S_c \simeq 67.383\,280$$

$$S_c^{(0)}(\mathbf{I}) = 32.141\,035$$

where

$$S_c^{(0)}(\mathbf{P}) = (\text{tr}[\mathbf{W}_0\mathbf{P}] + 1)(\text{tr}[\mathbf{K}_0\mathbf{P}^{-1}] + 1).$$

The same truncation will be applied hereafter in Example 2. The initial estimate  $\mathbf{P}_0$  was obtained using (38) as

$$\mathbf{P}_0 = \begin{bmatrix} 13.645\,073\,461 & -2.672\,229\,834 & -3.167\,232\,201 \\ -2.672\,229\,834 & 6.099\,162\,847 & -1.239\,205\,730 \\ -3.167\,232\,201 & -1.239\,205\,730 & 2.880\,278\,340 \end{bmatrix}$$

or equivalently

$$\mathbf{T}_0 = \begin{bmatrix} 2.058\,386\,621 & 2.912\,443\,826 & 0.962\,179\,266 \\ 0.993\,607\,653 & -0.940\,504\,298 & -2.056\,054\,071 \\ -0.223\,906\,291 & -1.287\,114\,675 & 1.083\,272\,877 \end{bmatrix}.$$

In this case, (51) yields

$$S_c(\mathbf{P}_0) \simeq 20.704\,544$$

$$S_c^{(0)}(\mathbf{P}_0) = 10.256\,183.$$

In order to minimize  $S_c(\mathbf{P})$  in (51), we apply (35) where  $\mathbf{H}(\mathbf{P})$

TABLE II  
 $L_2$  SENSITIVITY PROFILE OF THE FIRST TEN ITERATIONS

$i$	$L_2$ -Sensitivity	$i$	$L_2$ -Sensitivity
	67.38327994750175	5	20.57670199135822
0	20.70454380351688	6	20.57670199135815
1	20.57708100856075	7	20.57670199135823
2	20.57670318054942	8	20.57670199135817
3	20.57670199508960	9	20.57670199135819
4	20.57670199136992	10	20.57670199135813

and  $\mathbf{G}(\mathbf{P})$  are given by (52a) and (52b). Ten iterations yield

$$\mathbf{P}_{10} = \begin{bmatrix} 14.932\,872\,358 & -3.417\,647\,079 & -3.136\,150\,798 \\ -3.417\,647\,079 & 6.452\,249\,132 & -1.114\,438\,061 \\ -3.136\,150\,798 & -1.114\,438\,061 & 2.621\,026\,248 \end{bmatrix}$$

or equivalently

$$\mathbf{T}_{10} = \begin{bmatrix} 2.001\,302\,974 & 3.020\,964\,828 & 1.342\,173\,712 \\ 1.185\,944\,380 & -1.030\,092\,946 & -1.996\,169\,728 \\ -0.360\,253\,310 & -1.235\,733\,847 & 0.981\,939\,744 \end{bmatrix}.$$

Now (51) gives

$$S_c(\mathbf{P}_{10}) \simeq 20.576\,702$$

$$S_c^{(0)}(\mathbf{P}_{10}) = 10.308\,948.$$

The  $L_2$ -sensitivity profile of the first ten iterations is given in Table II. It is observed that, with a tolerance  $\varepsilon = 10^{-8}$ , the algorithm converges after four iterations.

By substituting  $\mathbf{T} = \mathbf{T}_{10}$  into (9), the optimal state-space model with minimum  $L_2$ -sensitivity is then obtained as

$$\bar{\mathbf{A}} = \begin{bmatrix} -0.083\,598\,330 & -0.505\,801\,407 & -0.130\,205\,421 \\ 0.505\,801\,407 & -0.191\,285\,835 & -0.417\,260\,962 \\ -0.130\,205\,421 & 0.417\,260\,962 & -0.353\,947\,136 \end{bmatrix}$$

$$\bar{\mathbf{b}} = [0.431\,487\,669 \ -0.518\,651\,368 \ 0.523\,993\,264]^T$$

$$\bar{\mathbf{c}} = [0.431\,487\,669 \ 0.518\,651\,368 \ 0.523\,993\,264]$$

$$d = -0.539\,352\,1.$$

*Example 3:* Consider a 2-D stable state-space digital filter of order (2,2) modeled by (53) where

$$\mathbf{A} = \begin{bmatrix} 1.888\,99 & -0.912\,19 & -1.0 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 \\ 0.027\,71 & -0.025\,80 & 1.888\,99 & 1.0 \\ -0.025\,80 & 0.024\,31 & -0.912\,19 & 0.0 \end{bmatrix}$$

$$\mathbf{b} = [0.219\,089 \ 0.0 \ -0.028\,889 \ 0.091\,219]^T$$

$$\mathbf{c} = [0.288\,89 \ -0.091\,219 \ -0.219\,089 \ 0.0]$$

$$d = 0.870\,21.$$

First, (60a), (60b), (62a), and (62b) were used with the truncation  $(0,0) \leq (i,j) \leq (100,100)$  to evaluate  $\mathbf{K}_{00}$  and  $\mathbf{W}_{00}$  as shown at the bottom of the next page. Next,  $\mathbf{K}_{\xi\zeta}$  and  $\mathbf{W}_{\xi\zeta}$  were calculated using (63a) and (63b) with truncations  $(0,0) \leq (\xi,\zeta) \leq (100,100)$  and  $(0,0) \leq (i,j) \leq (100,100)$ . The  $L_2$ -sensitivity measure was then computed using (65) as

$$S_{2d} = 1.049\,619\,37 \times 10^7$$

$$S_{2d}^{(0)}(\mathbf{I}) = 2.304\,669\,43 \times 10^5$$

where  $S_{2d}^{(0)}(\mathbf{I})$  is defined in (73).

The initial estimate  $\mathbf{P}^{(0)} = \mathbf{P}_1^{(0)} \oplus \mathbf{P}_4^{(0)}$  was computed using (74) as

$$\mathbf{P}_1^{(0)} = \begin{bmatrix} 18.23031272 & 18.79385826 \\ 18.79385826 & 20.02037679 \end{bmatrix}$$

$$\mathbf{P}_4^{(0)} = \begin{bmatrix} 0.18668352 & -0.17520854 \\ -0.17520854 & 0.17117480 \end{bmatrix}.$$

Then, using (69) and (73), we obtain

$$S_{2d}(\mathbf{P}^{(0)}) = 7.90009787 \times 10^3$$

$$S_{2d}^{(0)}(\mathbf{P}^{(0)}) = 501.39575676.$$

Finally, by applying (71a) and (71b), with two iterations we obtain  $\mathbf{P}^{(2)} = \mathbf{P}_1^{(2)} \oplus \mathbf{P}_4^{(2)}$  where

$$\mathbf{P}_1^{(2)} = \begin{bmatrix} 39.18732693 & 40.66977497 \\ 40.66977497 & 43.44455278 \end{bmatrix}$$

$$\mathbf{P}_4^{(2)} = \begin{bmatrix} 0.08808260 & -0.08208280 \\ -0.08208280 & 0.07946989 \end{bmatrix}$$

or equivalently

$$\mathbf{T}_1^{(2)} = \begin{bmatrix} 4.69642605 & 4.13895027 \\ 4.13895027 & 5.12968259 \end{bmatrix}$$

$$\mathbf{T}_4^{(2)} = \begin{bmatrix} 0.23320700 & -0.18356769 \\ -0.18356769 & 0.21394578 \end{bmatrix}.$$

The  $L_2$ -sensitivity profile of the first ten iterations is given in Table III. Then, (69) and (73) were calculated as

$$S_{2d}(\mathbf{P}^{(2)}) = 4.77995163 \times 10^3$$

$$S_{2d}^{(0)}(\mathbf{P}^{(2)}) = 807.46953086.$$

TABLE III  
 $L_2$  SENSITIVITY PROFILE OF THE FIRST TEN ITERATIONS

i	$L_2$ -sensitivity	i	$L_2$ -sensitivity
	$1.04961937 \times 10^7$	5	$4.77998761 \times 10^3$
0	$7.90009787 \times 10^3$	6	$4.77998760 \times 10^3$
1	$4.78462354 \times 10^3$	7	$4.77998760 \times 10^3$
2	$4.77995163 \times 10^3$	8	$4.77998760 \times 10^3$
3	$4.77998982 \times 10^3$	9	$4.77998760 \times 10^3$
4	$4.77998749 \times 10^3$	10	$4.77998760 \times 10^3$

By substituting  $\mathbf{T}^{(2)} = \mathbf{T}_1^{(2)} \oplus \mathbf{T}_4^{(2)}$  into (56), the optimal filter structure with minimum  $L_2$ -sensitivity was obtained as shown at the bottom of the page.

It is noted that the  $L_2$ -sensitivity measure  $S_{2d}(\mathbf{P}^{(2)}) = 4.77995163 \times 10^3$  in the optimal filter structure obtained here is considerably smaller than  $S_{2d} = 6.22257856 \times 10^3$  in the optimal filter structure obtained by Li [17].

### V. CONCLUSION

Novel expressions for the  $L_2$ -sensitivity of linear discrete-time systems, linear continuous-time systems, and 2-D state-space digital filters have been developed. Each expression is derived in terms of the general controllability and observability Gramians. It has been shown that, for each class of systems, the new formula offers a natural connection between the  $L_2$ -sensitivity and the  $L_1/L_2$ -mixed sensitivity. An iterative procedure has also been proposed to facilitate the minimization of the new expressions of the  $L_2$ -sensitivity measure. Our computer simulation results have demonstrated the effectiveness of the proposed techniques compared with several existing methods.

$$\mathbf{K}_{00} = \begin{bmatrix} 87.09615905 & 85.18196473 & 1.63649020 & -1.53601192 \\ 85.18196473 & 87.09270706 & 1.31792548 & -1.23006710 \\ 1.63649020 & 1.31792548 & 1.13375598 & -1.03484650 \\ -1.53601192 & -1.23006710 & -1.03484650 & 0.97280328 \end{bmatrix}$$

$$\mathbf{W}_{00} = \begin{bmatrix} 11.33767913 & -10.35207667 & 15.85678762 & 12.35025475 \\ -10.35207667 & 9.66197314 & -15.00759446 & -11.56613634 \\ 15.85678762 & -15.00759446 & 638.95928802 & 622.90731771 \\ 12.35025475 & -11.56613634 & 622.90731771 & 638.94455368 \end{bmatrix}$$

$$\bar{\mathbf{A}} = \begin{bmatrix} 0.96297355 & -0.14767517 & -0.17187245 & 0.13528851 \\ 0.13855174 & 0.92601645 & 0.13867749 & -0.10915927 \\ 0.07557078 & -0.03014312 & 0.98322700 & 0.14346512 \\ -0.03121145 & 0.05788553 & -0.15069421 & 0.90576300 \end{bmatrix}$$

$$\bar{\mathbf{b}} = [0.16146755 \quad -0.13028217 \quad 0.65224855 \quad 0.98600103]^T$$

$$\bar{\mathbf{c}} = [0.97919962 \quad 0.72777683 \quad -0.05109309 \quad 0.04021766]$$

## REFERENCES

- [1] L. Thiele, "Design of sensitivity and round-off noise optimal state-space discrete systems," *Int. J. Circuit Theory Appl.*, vol. 12, pp. 39–46, Jan. 1984.
- [2] —, "On the sensitivity of linear state-space systems," *IEEE Trans. Circuits Syst.*, vol. CAS-33, pp. 502–510, May 1986.
- [3] M. Iwatsuki, M. Kawamata, and T. Higuchi, "Statistical sensitivity and minimum sensitivity structures with fewer coefficients in discrete time linear systems," *IEEE Trans. Circuits Syst.*, vol. 37, pp. 72–80, Jan. 1989.
- [4] G. Li and M. Gevers, "Optimal finite precision implementation of a state-estimate feedback controller," *IEEE Trans. Circuits Syst.*, vol. 37, pp. 1487–1498, Dec. 1990.
- [5] G. Li, B. D. O. Anderson, M. Gevers, and J. E. Perkins, "Optimal FWL design of state-space digital systems with weighted sensitivity minimization and sparseness consideration," *IEEE Trans. Circuits Syst. I*, vol. 39, pp. 365–377, May 1992.
- [6] W.-Y. Yan and J. B. Moore, "On  $L^2$ -sensitivity minimization of linear state-space systems," *IEEE Trans. Circuits Syst. I*, vol. 39, pp. 641–648, Aug. 1992.
- [7] G. Li and M. Gevers, "Optimal synthetic FWL design of state-space digital filters," in *Proc. 1992 IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, vol. 4, pp. 429–432.
- [8] M. Gevers and G. Li, *Parametrizations in Control, Estimation and Filtering Problems: Accuracy Aspects*. Berlin, Germany: Springer-Verlag, 1993.
- [9] U. Helmke and J. B. Moore, *Optimization and Dynamical Systems*. London, U.K.: Springer-Verlag, 1994.
- [10] W. J. Lutz and S. L. Hakimi, "Design of multi-input multi-output systems with minimum sensitivity," *IEEE Trans. Circuits Syst.*, vol. 35, pp. 1114–1122, Sept. 1988.
- [11] M. Kawamata, T. Lin, and T. Higuchi, "Minimization of sensitivity of 2-D state-space digital filters and its relation to 2-D balanced realizations," in *Proc. 1987 IEEE Int. Symp. Circuits and Systems*, pp. 710–713.
- [12] T. Hinamoto, T. Hamanaka, and S. Maekawa, "Synthesis of 2-D state-space digital filters with low sensitivity based on the Fornasini-Marchesini model," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 38, pp. 1587–1594, Sept. 1990.
- [13] T. Hinamoto, T. Takao, and M. Muneyasu, "Synthesis of 2-D separable-denominator digital filters with low sensitivity," *J. Franklin Inst.*, vol. 329, pp. 1063–1080, 1992.
- [14] T. Hinamoto and T. Takao, "Synthesis of 2-D state-space filter structures with low frequency-weighted sensitivity," *IEEE Trans. Circuits Syst. II*, vol. 39, pp. 646–651, Sept. 1992.
- [15] —, "Minimization of frequency-weighting sensitivity in 2-D systems based on the Fornasini–Marchesini second model," in *1992 IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, pp. 401–404.
- [16] G. Li, "On frequency weighted minimal  $L_2$  sensitivity of 2-D systems using Fornasini-Marchesini LSS model," *IEEE Trans. Circuits Syst. I*, vol. 44, pp. 642–646, July 1997.
- [17] —, "Two-dimensional system optimal realizations with  $L_2$ -sensitivity minimization," *IEEE Trans. Signal Processing*, vol. 46, pp. 809–813, Mar. 1998.
- [18] T. Hinamoto, Y. Zempo, Y. Nishino, and W.-S. Lu, "An analytical approach for the synthesis of two-dimensional state-space filter structures with minimum weighted sensitivity," *IEEE Trans. Circuits Syst. I*, vol. 46, pp. 1172–1183, Oct. 1999.
- [19] K. Zhou and J. C. Doyle, *Essentials of Robust Control*. Upper Saddle River, NJ: Prentice-Hall, 1998.
- [20] T. Kailath, *Linear System*. Englewood Cliffs, N.J.: Prentice-Hall, 1980.
- [21] L. L. Scharf, *Statistical Signal Processing*. Reading, MA: Addison-Wesley, 1991.
- [22] R. P. Roesser, "A discrete state-space model for linear image processing," *IEEE Trans. Automat. Contr.*, vol. AC-20, pp. 1–10, Feb. 1975.



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

From 1972 to 1988, he was with the Faculty of Engineering, Kobe University. From 1979 to 1981, he was on leave from Kobe University as a visiting member of staff in the Department of Electrical Engineering, Queen's University, Kingston, ON, Canada. During 1988–1991, he was a Professor of electronic

circuits in the Faculty of Engineering, Tottori University, Japan. Since January 1992, he has been a Professor of electronic control in the Department of Elec-

trical Engineering, Hiroshima University, Higashi-Hiroshima, Japan. His research interests include digital signal processing, system theory, and control engineering. He has published more than 250 papers in these areas and is the co-editor and coauthor of *Two-Dimensional Signal and Image Processing* (Tokyo, Japan: SICE, 1996).

Dr. Hinamoto served as an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II from 1993 to 1995 and presently serves as an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS I. He also served as Chair of the 12th Digital Signal Processing (DSP) Symposium held in Hiroshima in November 1997, sponsored by the DSP Technical Committee of IEICE. He was the Guest Editor of the special section of DSP in the August 1998 issue of the *IEICE Transactions on Fundamentals*. Since 1995, he has been a member of the steering committee of the IEEE Midwest Symposium on Circuits and Systems and, since 1998, a member of the Digital Signal Processing Technical Committee in the IEEE Circuits and Systems Society. He served as a member of the Technical Program Committee for ISCAS'99. From 1993 to 2000, he served as a senator or member of the Board of Directors in the Society of Instrument and Control Engineers (SICE), and from 1999 to 2001, he was Chair of the Chugoku Chapter of SICE. He played a leading role in establishing the Hiroshima Section of IEEE, and served as the Interim Chair of the section. He is a recipient of the IEEE Third Millennium Medal.



**Shuichi Yokoyama** received the B.E. and M.E. degrees in electrical engineering from Hiroshima University, Higashi-Hiroshima, Japan, in 1996 and 1998, respectively.

He was engaged in research on digital signal processing during his graduate studies. Since 1998, he has been with NTT DoCoMo Chugoku Inc., Hiroshima, Japan.



**Takuya Inoue** received the B.E. and M.E. degrees in electrical engineering from Hiroshima University, Higashi-Hiroshima, Japan, in 1999 and 2001, respectively.

He was engaged in research on digital signal processing during his graduate studies. Since 2001, he has been with Fujitsu System Solutions Limited, Tokyo, Japan.



**Wu Zeng** received the B.E. degree in automatic control engineering from Beijing Institute of Mechanical Industry, Beijing, China, and the M.E. degree in electrical engineering from Hiroshima University, Higashi-Hiroshima, Japan, in 1996 and 2000, respectively. He was engaged in research on digital signal processing during his graduate studies. Since 2000, he has been with Precision Digital Corporation, Osaka, Japan.



**Wu-Sheng Lu** (S'81–M'85–SM'90–F'99) received his undergraduate education in mathematics from Fudan University, Shanghai, China, from 1959 to 1964, and the M.S. degree in electrical engineering and the Ph.D. degree in control science from the University of Minnesota, Minneapolis, in 1983, and 1984, respectively.

He was a Post-Doctoral Fellow at University of Victoria, Victoria, BC, Canada, in 1985 and a Visiting Assistant Professor at the University of Minnesota in 1986. Since 1987, he has been with University of Victoria where he is currently a Professor. His teaching and research interests are in the areas of digital signal processing and application of optimization methods. He is the coauthor, with A. Antoniou, of *Two-Dimensional Digital Filters* (New York: Marcel Dekker, 1992).

Dr. Lu was an Associate Editor of the *Canadian Journal of Electrical and Computer Engineering* in 1989, and the Editor of the same journal from 1990 to 1992. He served as an Associate Editor for IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II from 1993 to 1995 and for IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS I from 1999 to 2001. He is presently an Associate Editor for the *International Journal of Multidimensional Systems and Signal Processing*. He is a Fellow of the Engineering Institute of Canada.