

Realization of 3-D Separable-Denominator Digital Filters with Low l_2 -Sensitivity

Takao Hinamoto, *Life Fellow, IEEE*, Akimitsu Doi and Wu-Sheng Lu, *Fellow, IEEE*

Abstract—Three-dimensional (3-D) digital filters find applications in a variety of image and video signal processing problems. This paper presents a coefficient-sensitivity analysis for a wide class of 3-D digital filters with separable denominators in local state space that leads to an analytic formulation for sensitivity minimization, and to present two solution techniques for the sensitivity minimization problem at hand. To this end, a vector-matrix-vector decomposition of a given 3-D transfer function that separates the three variables and leads to a state-space realization in a form convenient for subsequent analysis. An l_2 -sensitivity analysis is then performed, the result is a computationally tractable formula of the overall l_2 -sensitivity for 3-D digital filters. The l_2 -sensitivity is minimized subject to l_2 -scaling constraints by using one of the two solution methods proposed – one relaxes the constraints into a single trace constraint and solves the relaxed problem with an effective matrix iteration scheme; while the other converts the contained optimization problem at hand into an unconstrained problem and solves it using a quasi-Newton algorithm. A case study is presented to illustrate the validity and effectiveness of the proposed techniques.

Index Terms—3-D separable-denominator digital filters, minimal state-space realization, l_2 -sensitivity analysis, low l_2 -sensitivity, l_2 -scaling constraints, no overflow oscillations, Lagrange function, bisection method, quasi-Newton method

I. INTRODUCTION

It is of practical significance in many applications to construct a filter structure so that the coefficient sensitivity of a digital filter is minimum or nearly minimum in a certain sense. Due to finite-word-length (FWL) effects caused by coefficient truncation or rounding, poor sensitivity may lead to degradation of the transfer characteristics in a FWL implementation of the digital filter. For instance, the characteristics of an originally stable filter might be so altered that the filter may become unstable. This motivates the study of the coefficient sensitivity minimization problem for digital filters. Several techniques have been proposed to analyze l_2 -sensitivity and to synthesize the state-space model structures that minimize l_2 -sensitivity [1]–[6]. The minimization of l_2 -sensitivity for two-dimensional (2-D) state-space digital filters has also been investigated [5], [7]–[9]. More recently, the minimization problem of l_2 -sensitivity subject to l_2 -scaling constraints has been treated for 1-D and 2-D state-space digital filters [10]–[13]. It is known that the use of scaling constraints can be beneficial

for suppressing overflow [14], [15]. In addition, considerable research interest has also been observed in the design of multidimensional (M-D) recursive digital filters [16]–[19]. Our study of 3-D separable-denominator digital filters is motivated as it fits naturally into typical time-space digital filtering settings. An example of this scenario is a video processing task such as compression or de-noising of a video clip, in that the signals of interest assume the form of a time series, with a 2-D spatial-domain signal known as image at each sampling time instant. For a signal compression task, since the signal redundancy in the time domain and spatial domain are inherently different, the filters to be used for time-domain processing and spatial-domain processing have to be distinctly designed, this justifies the use of a 3-D filter of the form $H(z_1, z_2, z_3) = H_1(z_1)H_2(z_2, z_3)$, where z_1 and (z_2, z_3) are associated with the time domain and the 2-D spatial domain, respectively. Assume that for processing efficiency one decides to use IIR filters, then $H_2(z_2, z_3)$ is a 2-D IIR filter. Since for most spatial filtering the desired frequency responses are quadrantly symmetrical, it is well known that $H_2(z_2, z_3)$ possesses separable denominators [20]. As a result, it is quite natural to study 3-D IIR separable-denominator digital filters. Relevant recent studies also include state-space realization of general M-D filters and possible applications of state-space realization in uncertainty modeling [21], and 3-D realization and its applications in distributed grid sensor networks [22]. In addition, a state-space model for general M-D spatially distributed dynamic system is proposed in [23], and new results on stability and stability margin for 2-D systems are reported in [24]. On the other hand, the literature offers only handful results on efficient realization of 3-D state-space digital filters with minimum coefficient sensitivity [25]. This is likely due to the fact that the problems encountered in the 3-D filters are considerably more involved and challenging relative to their 1-D and 2-D counterparts because here one deals with 3-D coefficient arrays instead of coefficient vectors for 1-D and coefficient matrices for 2-D filters.

The objectives of this paper are twofold: to present a coefficient-sensitivity analysis for a wide class of 3-D digital filters with separable denominators in local state space that leads to an analytic formulation for sensitivity minimization, and to present two solution techniques for the sensitivity minimization problem at hand. To this end, we present a vector-matrix-vector decomposition of a given 3-D transfer function that separates the three frequency-dependent variables (z_1, z_2, z_3) , and leads in turn to a state-space realization in a form convenient for subsequent analysis. An l_2 -sensitivity analysis is then carried out, a central result of the analysis

Manuscript received January 26, 2012; revised mm dd, 2012.

T. Hinamoto and A. Doi are with Hiroshima Institute of Technology Hiroshima 731-5193, Japan (e-mails: hinamoto@ieee.org, doi@cc.it-hiroshima.ac.jp, Phone:+81-82-921-4338, Fax:+81-82-921-8978)

W.-S. Lu is with the Department of Electrical and Computer Engineering, University of Victoria, Victoria, B.C, Canada, V8W 3P6. (e-mail: wslu@ece.uvic.ca, Phone:+1-250-721-8692, Fax:+1-250-721-6052)

is a computationally tractable formula of the overall l_2 -sensitivity of a 3-D separable-denominator state-space digital filter. We remark that the formula is of considerable difference from that of [25] as it takes into consideration the particular structure of the realization, especially with its fixed 0 and 1 components subject to dynamic-range l_2 -scaling constraints. As a result, this paper deals with a non-convex constrained optimization problem to produce a state-space realization with reduced number of nontrivial parameters, guaranteed freedom of internal overflow and reasonably low coefficient sensitivity. We present two solution methods – one relaxes the constraints into a single trace constraint and solves the relaxed problem with an effective matrix iteration scheme where the Lagrange multiplier is determined via a bisection technique; while the other converts the contained optimization problem at hand into an unconstrained problem and solves it using a quasi-Newton algorithm. Closed-form formula for the gradient is derived for efficient evaluation. A numerical example is presented to illustrate the validity and effectiveness of the proposed techniques in Section IV.

Throughout, \mathbf{I}_n denotes the identity matrix of dimension $n \times n$. The transpose (conjugate transpose) of a matrix \mathbf{A} is indicated by \mathbf{A}^T (\mathbf{A}^*). $\text{tr}[\mathbf{A}]$ is used to denote the trace of a square matrix \mathbf{A} . Moreover, bold uppercase, bold lowercase and plain lowercase are used to make the distinction between matrices, vectors and scalar values, respectively.

II. REALIZATION AND SENSITIVITY ANALYSIS

A. Minimal Realization

Consider a stable 3-D separable-denominator digital filter

$$H(z_1, z_2, z_3) = \frac{N(z_1, z_2, z_3)}{D_1(z_1)D_2(z_2)D_3(z_3)} \quad (1a)$$

where the denominator and the numerator are assumed to be coprime and

$$N(z_1, z_2, z_3) = \sum_{i=0}^{N_1} \sum_{j=0}^{N_2} \sum_{k=0}^{N_3} a_{ijk} z_1^{-i} z_2^{-j} z_3^{-k} \quad (1b)$$

$$D_l(z_l) = 1 + \sum_{\xi=1}^{N_l} b_{l\xi} z_l^{-\xi} \quad \text{for } l = 1, 2, 3.$$

Although $H(z_1, z_2, z_3)$ has separable denominator, it is not a separable transfer function because the numerator $N(z_1, z_2, z_3)$ is a non-separable polynomial. A key step towards a minimal state-space realization of $H(z_1, z_2, z_3)$ and subsequent minimization of its l_2 -sensitivity is to separate the terms in three variables in $H(z_1, z_2, z_3)$. To this end, algebraic manipulations are performed on $N(z_1, z_2, z_3)$, which lead (1a) to a vector-matrix-vector expression where $H(z_1, z_2, z_3)$ are completely separated in its three variables. Namely,

$$H(z_1, z_2, z_3) = \mathbf{f}_1(z_1) \mathbf{H}_2(z_2) \mathbf{g}_3(z_3) \quad (2a)$$

where

$$\mathbf{f}_1(z_1) = \frac{[1, z_1^{-1}, \dots, z_1^{-N_1}]}{D_1(z_1)}, \quad \mathbf{g}_3(z_3) = \frac{[1, z_3^{-1}, \dots, z_3^{-N_3}]^T}{D_3(z_3)}$$

$$\mathbf{H}_2(z_2) = \frac{\Delta_0 + \Delta_1 z_2^{-1} + \dots + \Delta_{N_2} z_2^{-N_2}}{D_2(z_2)}$$

$$\Delta_l = \begin{bmatrix} a_{0l0} & a_{0l1} & \cdots & a_{0lN_3} \\ a_{1l0} & a_{1l1} & \cdots & a_{1lN_3} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N_1l0} & a_{N_1l1} & \cdots & a_{N_1lN_3} \end{bmatrix}, \quad l=0, 1, \dots, N_2. \quad (2b)$$

We note that the decomposition in (2a) is somewhat similar to that of a finite 3-D array into the product of two finite 1-D vector arrays and another 1-D finite matrix array in the middle, as done in [26]. Furthermore, if the numerator in (1a) is separable, that is, $N(z_1, z_2, z_3) = N_1(z_1)N_2(z_2)N_3(z_3)$ with $N_l(z_l) = \sum_{\xi=0}^{N_l} a_{l\xi} z_l^{-\xi}$ for $l = 1, 2, 3$, then $a_{ijk} = a_{1i}a_{2j}a_{3k}$ holds for $i = 0, 1, \dots, N_1$, $j = 0, 1, \dots, N_2$ and $k = 0, 1, \dots, N_3$. In such a case, matrices Δ_l for $l = 0, 1, \dots, N_2$ will possess considerably higher sparsity.

The 1-D transfer function $\mathbf{H}_2(z_2)$ in (2b) has $(N_3 + 1)$ inputs and $(N_1 + 1)$ outputs. It can be realized with a minimal state-space model $(\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \Delta_0)_p$ as

$$\begin{aligned} \mathbf{x}(k+1) &= \mathbf{A}_2 \mathbf{x}(k) + \mathbf{B}_2 \mathbf{u}(k) \\ \mathbf{y}(k) &= \mathbf{C}_2 \mathbf{x}(k) + \Delta_0 \mathbf{u}(k) \end{aligned} \quad (3)$$

where $\mathbf{x}(k)$ is a $p \times 1$ state-variable vector, $\mathbf{u}(k)$ is an $(N_3 + 1) \times 1$ input vector, $\mathbf{y}(k)$ is an $(N_1 + 1) \times 1$ output vector, and $\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2$ and Δ_0 are real constant matrices of appropriate dimensions. Here, p is the least dimension such that a state-space realization of $\mathbf{H}_2(z_2)$ is controllable and observable. Such a realization is called a minimal realization [27]. The reader is referred to Appendix I that explains how this realization can actually be constructed. The transfer function of the linear system in (3) can be expressed as

$$\mathbf{H}_2(z_2) = \mathbf{C}_2(z_2 \mathbf{I}_p - \mathbf{A}_2)^{-1} \mathbf{B}_2 + \Delta_0. \quad (4)$$

The 1-D transfer function $\mathbf{f}_1(z_1)$ in (2b) has $(N_1 + 1)$ inputs and a single output, while the 1-D transfer function $\mathbf{g}_3(z_3)$ in (2b) has a single input and $(N_3 + 1)$ outputs. Consequently they can be realized with minimal state-space models $(\mathbf{A}_1, \mathbf{B}_1, \bar{\mathbf{e}}_{N_1}^T, \bar{\mathbf{e}}_1^T)_{N_1}$ and $(\mathbf{A}_3, \hat{\mathbf{e}}_{N_3}, \mathbf{C}_3, \hat{\mathbf{e}}_1)_{N_3}$ as

$$\begin{aligned} \mathbf{f}_1(z_1) &= \frac{(\bar{\mathbf{e}}_2^T - b_{11} \bar{\mathbf{e}}_1^T) z_1^{-1} + \dots + (\bar{\mathbf{e}}_{N_1}^T - b_{1N_1} \bar{\mathbf{e}}_1^T) z_1^{-N_1}}{1 + b_{11} z_1^{-1} + \dots + b_{1N_1} z_1^{-N_1}} + \bar{\mathbf{e}}_1^T \\ &= \bar{\mathbf{e}}_{N_1}^T (z_1 \mathbf{I}_{N_1} - \mathbf{A}_1)^{-1} \mathbf{B}_1 + \bar{\mathbf{e}}_1^T \\ \mathbf{g}_3(z_3) &= \frac{(\hat{\mathbf{e}}_2 - b_{31} \hat{\mathbf{e}}_1) z_3^{-1} + \dots + (\hat{\mathbf{e}}_{N_3} - b_{3N_3} \hat{\mathbf{e}}_1) z_3^{-N_3}}{1 + b_{31} z_3^{-1} + \dots + b_{3N_3} z_3^{-N_3}} + \hat{\mathbf{e}}_1 \\ &= \mathbf{C}_3(z_3 \mathbf{I}_{N_3} - \mathbf{A}_3)^{-1} \hat{\mathbf{e}}_{N_3} + \hat{\mathbf{e}}_1, \end{aligned} \quad (5a)$$

respectively, where [27]

$$\mathbf{A}_1 = \begin{bmatrix} \mathbf{0} & \vdots \\ \cdots & \mathbf{b}_1 \\ \mathbf{I}_{N_1-1} & \vdots \end{bmatrix}, \quad \mathbf{B}_1 = \begin{bmatrix} 0 & \cdots & 1 \\ \mathbf{b}_1 & \vdots & \ddots & \vdots \\ 1 & \cdots & 0 \end{bmatrix}$$

$$\mathbf{A}_3 = \begin{bmatrix} \mathbf{0} & \vdots & \mathbf{I}_{N_3-1} \\ \cdots & \cdots & \cdots \\ \mathbf{c}_3 & & \end{bmatrix}, \quad \mathbf{C}_3 = \begin{bmatrix} \mathbf{c}_3 \\ \cdots & \cdots & \cdots \\ 0 & \cdots & 0 & 1 \\ \vdots & \ddots & \vdots \\ 1 & 0 & \cdots & 0 \end{bmatrix}$$

$$\begin{aligned} [\bar{e}_1, \bar{e}_2, \dots, \bar{e}_{N_1}] &= \mathbf{I}_{N_1}, & [\hat{e}_1, \hat{e}_2, \dots, \hat{e}_{N_3}] &= \mathbf{I}_{N_3} \\ \mathbf{b}_1 &= -[b_{1N_1}, \dots, b_{12}, b_{11}]^T, & \mathbf{c}_3 &= -[b_{3N_3}, \dots, b_{32}, b_{31}]^T. \end{aligned} \quad (5b)$$

The state-space model in (3) contains $(p + N_1 + 1)(p + N_3 + 1)$ nontrivial and independent parameters, while the numbers of independent parameters in (5a) are N_1 and N_3 , respectively. We stress that although the parameter vectors in each of $(\mathbf{A}_1, \mathbf{B}_1)$ and $(\mathbf{A}_3, \mathbf{C}_3)$ appear twice, as independent parameters we count each of them only once because, when it varies, the two copies vary in exactly the same way.

As a result, a local state-space model for the 3-D digital filter with separable denominator in (1a) can be realized by

$$\begin{aligned} \mathbf{x}'(i, j, k) &= \mathbf{A}\mathbf{x}(i, j, k) + \mathbf{b}u(i, j, k) \\ y(i, j, k) &= \mathbf{c}\mathbf{x}(i, j, k) + du(i, j, k) \end{aligned} \quad (6a)$$

where $i, j, k \geq 0$ and

$$\begin{aligned} \mathbf{x}'(i, j, k) &= \begin{bmatrix} \mathbf{x}^h(i+1, j, k) \\ \mathbf{x}^v(i, j+1, k) \\ \mathbf{x}^a(i, j, k+1) \end{bmatrix}, & \mathbf{x}(i, j, k) &= \begin{bmatrix} \mathbf{x}^h(i, j, k) \\ \mathbf{x}^v(i, j, k) \\ \mathbf{x}^a(i, j, k) \end{bmatrix} \\ \mathbf{A} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_4 & \mathbf{A}_6 \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{A}_5 \\ \mathbf{0} & \mathbf{0} & \mathbf{A}_3 \end{bmatrix}, & \mathbf{b} &= \begin{bmatrix} \mathbf{B}_1 \Delta_0 \hat{\mathbf{e}}_1 \\ \mathbf{B}_2 \hat{\mathbf{e}}_1 \\ \hat{\mathbf{e}}_{N_3} \end{bmatrix} \\ \mathbf{c} &= [\bar{\mathbf{e}}_{N_1}^T \mathbf{C}_2 \quad \bar{\mathbf{e}}_1^T \Delta_0 \mathbf{C}_3], & d &= a_{000} \\ \mathbf{A}_4 &= \mathbf{B}_1 \mathbf{C}_2, & \mathbf{A}_5 &= \mathbf{B}_2 \mathbf{C}_3, & \mathbf{A}_6 &= \mathbf{B}_1 \Delta_0 \mathbf{C}_3. \end{aligned} \quad (6b)$$

Here, $\mathbf{x}^h(i, j, k)$ is an $N_1 \times 1$ horizontal state vector, $\mathbf{x}^v(i, j, k)$ is a $p \times 1$ vertical state vector, $\mathbf{x}^a(i, j, k)$ is an $N_3 \times 1$ additional state vector, $u(i, j, k)$ is a scalar input, and $y(i, j, k)$ is a scalar output.

In summary, the 3-D digital filter under consideration admits an implementation scheme as illustrated in Fig. 1 showing vividly a system structure that allows one to focus on optimizing its dominating 1-D MIMO subsystem $(\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \Delta_0)_p$ in order to reduce its coefficient sensitivity.

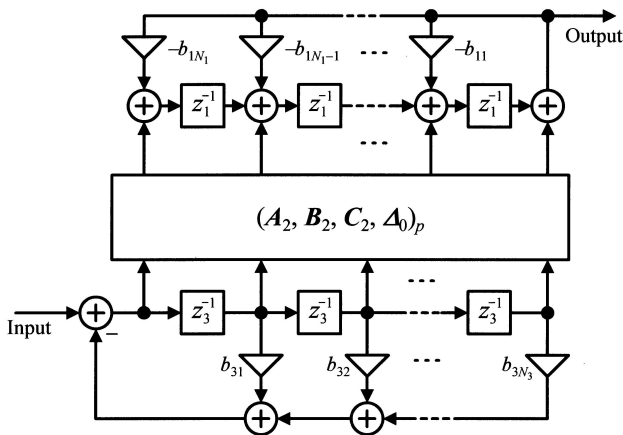


Fig. 1. Block diagram of a 3-D separable-denominator digital filter.

B. l_2 -Sensitivity Analysis

The l_2 -sensitivities of the 3-D transfer function in (2a) with respect to coefficient matrices $\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \Delta_0, \mathbf{b}_1$, and \mathbf{c}_3^T are computed as follows.

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

$$\mathbf{S}_{\mathbf{X}} = \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}}, \quad (\mathbf{S}_{\mathbf{X}})_{ij} = \frac{\partial f(\mathbf{X})}{\partial x_{ij}} \quad (7)$$

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

By means of (2a), (4), (5a), Definition 1, and the formula

$$\frac{\partial \mathbf{A}^{-1}}{\partial t} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial t} \mathbf{A}^{-1}, \quad (8)$$

the sensitivities of $H(z_1, z_2, z_3)$ with respect to matrices $\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \Delta_0, \mathbf{b}_1$, and \mathbf{c}_3^T are evaluated as

$$\begin{aligned} \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{A}_2} &= [\mathbf{f}(z_2, z_3) \mathbf{g}(z_1, z_2)]^T \\ \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{B}_2} &= [\mathbf{g}_3(z_3) \mathbf{g}(z_1, z_2)]^T \\ \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{C}_2} &= [\mathbf{f}(z_2, z_3) \mathbf{f}_1(z_1)]^T \\ \frac{\partial H(z_1, z_2, z_3)}{\partial \Delta_0} &= [\mathbf{g}_3(z_3) \mathbf{f}_1(z_1)]^T \\ \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{b}_1} &= [\mathbf{f}_1(z_1) \mathbf{H}_2(z_2) \mathbf{g}_3(z_3) \mathbf{g}_1(z_1)]^T \\ \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{c}_3^T} &= \mathbf{f}_3(z_3) \mathbf{f}_1(z_1) \mathbf{H}_2(z_2) \mathbf{g}_3(z_3) \end{aligned} \quad (9a)$$

where

$$\begin{aligned} \mathbf{f}(z_2, z_3) &= (z_2 \mathbf{I}_p - \mathbf{A}_2)^{-1} \mathbf{B}_2 \mathbf{g}_3(z_3) \\ \mathbf{g}(z_1, z_2) &= \mathbf{f}_1(z_1) \mathbf{C}_2 (z_2 \mathbf{I}_p - \mathbf{A}_2)^{-1} \\ \mathbf{g}_1(z_1) &= \bar{\mathbf{e}}_{N_1}^T (z_1 \mathbf{I}_{N_1} - \mathbf{A}_1)^{-1} \\ \mathbf{f}_3(z_3) &= (z_3 \mathbf{I}_{N_3} - \mathbf{A}_3)^{-1} \hat{\mathbf{e}}_{N_3}. \end{aligned} \quad (9b)$$

Definition 2: Let $\mathbf{X}(z_1, z_2, z_3)$ be an $m \times n$ complex-valued matrix function of complex variables z_1, z_2 and z_3 , and $x_{pq}(z_1, z_2, z_3)$ be the (p, q) th entry of $\mathbf{X}(z_1, z_2, z_3)$. The l_2 -norm of $\mathbf{X}(z_1, z_2, z_3)$ is defined as

$$\begin{aligned} \|\mathbf{X}(z_1, z_2, z_3)\|_2 &= \left[\frac{1}{(2\pi)^3} \int_0^{2\pi} \int_0^{2\pi} \int_0^{2\pi} \sum_{p=1}^m \sum_{q=1}^n |x_{pq}(e^{j\omega_1}, e^{j\omega_2}, e^{j\omega_3})|^2 d\omega_1 d\omega_2 d\omega_3 \right]^{\frac{1}{2}} \\ &= \left(\text{tr} \left[\frac{1}{(2\pi j)^3} \oint_{|z_1|=1} \oint_{|z_2|=1} \oint_{|z_3|=1} \mathbf{X}(z_1, z_2, z_3) \mathbf{X}^*(z_1, z_2, z_3) \frac{dz_1}{z_1} \frac{dz_2}{z_2} \frac{dz_3}{z_3} \right] \right)^{\frac{1}{2}}. \end{aligned} \quad (10)$$

With (10), the overall l_2 -sensitivity measure for the 3-D transfer function in (2a) is defined using (4) and (5a) by

$$S = \left\| \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{A}_2} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{B}_2} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{C}_2} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{\Delta}_0} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{b}_1} \right\|_2^2 + \left\| \frac{\partial H(z_1, z_2, z_3)}{\partial \mathbf{c}_3^T} \right\|_2^2. \quad (11)$$

We remark that measure S in (11) is a natural 3-D extension of the classical l_2 -sensitivity measure for 1-D state-space digital filters [1]. We also remark that coefficient sensitivity is closely related to FWL effects, the reader is referred to Appendix II for further details.

To derive a computationally tractable formula for sensitivity S , we need to evaluate the terms in (9a) explicitly. To this end, we write the impulse responses of the following three 2-D transfer functions as

$$\begin{aligned} g_3(z_3)f_1(z_1) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{R}_{ij} z_1^{-i} z_3^{-j} \\ \mathbf{H}_2(z_2)g_3(z_3) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{r}_{ij} z_2^{-i} z_3^{-j} \\ f_1(z_1)\mathbf{H}_2(z_2) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \hat{\mathbf{r}}_{ij} z_1^{-i} z_2^{-j} \end{aligned} \quad (12a)$$

where for $i \geq 1$ and $j \geq 1$,

$$\begin{aligned} \mathbf{R}_{00} &= \hat{\mathbf{e}}_1 \bar{\mathbf{e}}_1^T, \quad \mathbf{R}_{ij} = \mathbf{C}_3 \mathbf{A}_3^{j-1} \hat{\mathbf{e}}_{N_3} \bar{\mathbf{e}}_{N_1}^T \mathbf{A}_1^{i-1} \mathbf{B}_1 \\ \mathbf{R}_{i0} &= \hat{\mathbf{e}}_1 \bar{\mathbf{e}}_{N_1}^T \mathbf{A}_1^{i-1} \mathbf{B}_1, \quad \mathbf{R}_{0j} = \mathbf{C}_3 \mathbf{A}_3^{j-1} \hat{\mathbf{e}}_{N_3} \bar{\mathbf{e}}_1^T \\ \mathbf{r}_{00} &= \mathbf{\Delta}_0 \hat{\mathbf{e}}_1, \quad \mathbf{r}_{ij} = \mathbf{C}_2 \mathbf{A}_2^{i-1} \mathbf{B}_2 \mathbf{C}_3 \mathbf{A}_3^{j-1} \hat{\mathbf{e}}_{N_3} \\ \mathbf{r}_{i0} &= \mathbf{C}_2 \mathbf{A}_2^{i-1} \mathbf{B}_2 \hat{\mathbf{e}}_1, \quad \mathbf{r}_{0j} = \mathbf{\Delta}_0 \mathbf{C}_3 \mathbf{A}_3^{j-1} \hat{\mathbf{e}}_{N_3} \\ \hat{\mathbf{r}}_{00} &= \bar{\mathbf{e}}_1^T \mathbf{\Delta}_0, \quad \hat{\mathbf{r}}_{ij} = \bar{\mathbf{e}}_{N_1}^T \mathbf{A}_1^{i-1} \mathbf{B}_1 \mathbf{C}_2 \mathbf{A}_2^{j-1} \mathbf{B}_2 \\ \hat{\mathbf{r}}_{i0} &= \bar{\mathbf{e}}_{N_1}^T \mathbf{A}_1^{i-1} \mathbf{B}_1 \mathbf{\Delta}_0, \quad \hat{\mathbf{r}}_{0j} = \bar{\mathbf{e}}_1^T \mathbf{C}_2 \mathbf{A}_2^{j-1} \mathbf{B}_2. \end{aligned} \quad (12b)$$

Since the 3-D filter in (1a) is assumed to be stable, each series in (12a) is convergent, hence the infinite sum can be approximated by a finite sum with $(0,0) \leq (i,j) \leq (I,J)$ provided that positive integers I and J are sufficiently large. From (12b), it follows that the adequate numerical values of such I and J depend on the spectral radii of matrices \mathbf{A}_1 , \mathbf{A}_2 , and \mathbf{A}_3 . A practical approach to identify the right values of I and J is by trial-and-error in that one computes a truncated series with certain (I,J) , then repeating the computation with both I and J increased by a certain amount and compare the two results. The process continues until the difference in norm between the two resulting matrices is less than a prescribed tolerance.

From (9a), (9b) and (12a) it follows that

$$\begin{aligned} f(z_2, z_3)g(z_1, z_2) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} \mathbf{I}_p & \mathbf{0} \end{bmatrix} \\ &\cdot \left(z_2 \mathbf{I}_{2p} - \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_p \end{bmatrix} z_1^{-i} z_3^{-j} \\ g_3(z_3)g(z_1, z_2) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{R}_{ij} \mathbf{C}_2 (z_2 \mathbf{I}_p - \mathbf{A}_2)^{-1} z_1^{-i} z_3^{-j} \end{aligned}$$

$$\begin{aligned} f(z_2, z_3)f_1(z_1) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} (z_2 \mathbf{I}_p - \mathbf{A}_2)^{-1} \mathbf{B}_2 \mathbf{R}_{ij} z_1^{-i} z_3^{-j} \\ f_1(z_1)\mathbf{H}_2(z_2)g_3(z_3)g_1(z_1) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} [\bar{\mathbf{e}}_{N_1}^T \quad \bar{\mathbf{e}}_1^T \mathbf{r}_{ij} \bar{\mathbf{e}}_{N_1}^T] \\ &\cdot \left(z_1 \mathbf{I}_{2N_1} - \begin{bmatrix} \mathbf{A}_1 & \mathbf{B}_1 \mathbf{r}_{ij} \bar{\mathbf{e}}_{N_1}^T \\ \mathbf{0} & \mathbf{A}_1 \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_{N_1} \end{bmatrix} z_2^{-i} z_3^{-j} \\ f_3(z_3)f_1(z_1)\mathbf{H}_2(z_2)g_3(z_3) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} [\mathbf{I}_{N_3} \quad \mathbf{0}] \\ &\cdot \left(z_3 \mathbf{I}_{2N_3} - \begin{bmatrix} \mathbf{A}_3 & \hat{\mathbf{e}}_{N_3} \hat{\mathbf{r}}_{ij} \mathbf{C}_3 \\ \mathbf{0} & \mathbf{A}_3 \end{bmatrix} \right)^{-1} \begin{bmatrix} \hat{\mathbf{e}}_{N_3} \hat{\mathbf{r}}_{ij} \hat{\mathbf{e}}_1 \\ \hat{\mathbf{e}}_{N_3} \end{bmatrix} z_1^{-i} z_2^{-j}. \end{aligned} \quad (13)$$

Referring to (9a), (10) and (13), the l_2 -sensitivity measure in (11) can be expressed as

$$S = \text{tr}[\mathbf{M}_A(\mathbf{I}_p)] + \text{tr}[\mathbf{W}_B] + \text{tr}[\mathbf{K}_C] + \text{tr}[\mathbf{N}_{\Delta_0}] + \text{tr}[\mathbf{W}_1] + \text{tr}[\mathbf{K}_3] \quad (14a)$$

where Gramians $\mathbf{M}_A(P)$, \mathbf{W}_B , \mathbf{K}_C , \mathbf{N}_{Δ_0} , \mathbf{W}_1 , and \mathbf{K}_3 can be computed using

$$\begin{aligned} \mathbf{M}_A(P) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} \mathbf{0} & \mathbf{I}_p \end{bmatrix} \mathbf{M}_{ij}^A \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_p \end{bmatrix} \\ \mathbf{M}_{ij}^A &= \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}^T \mathbf{M}_{ij}^A \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{P}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{W}_B &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{W}_{ij}^B, \quad \mathbf{K}_C = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{K}_{ij}^C \\ \mathbf{W}_{ij}^B &= \mathbf{A}_2^T \mathbf{W}_{ij}^B \mathbf{A}_2 + \mathbf{C}_2^T \mathbf{R}_{ij}^T \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{K}_{ij}^C &= \mathbf{A}_2 \mathbf{K}_{ij}^C \mathbf{A}_2^T + \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{R}_{ij}^T \mathbf{B}_2^T \\ \mathbf{N}_{\Delta_0} &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \mathbf{R}_{ij}^T \mathbf{R}_{ij} \\ \mathbf{W}_1 &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} \mathbf{0} & \mathbf{I}_{N_1} \end{bmatrix} \mathbf{W}_{ij}^1 \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_{N_1} \end{bmatrix} \\ \mathbf{W}_{ij}^1 &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{B}_1 \mathbf{r}_{ij} \bar{\mathbf{e}}_{N_1}^T \\ \mathbf{0} & \mathbf{A}_1 \end{bmatrix}^T \mathbf{W}_{ij}^1 \begin{bmatrix} \mathbf{A}_1 & \mathbf{B}_1 \mathbf{r}_{ij} \bar{\mathbf{e}}_{N_1}^T \\ \mathbf{0} & \mathbf{A}_1 \end{bmatrix} + \begin{bmatrix} \bar{\mathbf{e}}_{N_1}^T & \bar{\mathbf{e}}_1^T \mathbf{r}_{ij} \bar{\mathbf{e}}_{N_1}^T \end{bmatrix}^T \begin{bmatrix} \bar{\mathbf{e}}_{N_1}^T & \bar{\mathbf{e}}_1^T \mathbf{r}_{ij} \bar{\mathbf{e}}_{N_1}^T \end{bmatrix} \\ \mathbf{K}_3 &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} \mathbf{I}_{N_3} & \mathbf{0} \end{bmatrix} \mathbf{K}_{ij}^3 \begin{bmatrix} \mathbf{I}_{N_3} \\ \mathbf{0} \end{bmatrix} \\ \mathbf{K}_{ij}^3 &= \begin{bmatrix} \mathbf{A}_3 & \hat{\mathbf{e}}_{N_3} \hat{\mathbf{r}}_{ij} \mathbf{C}_3 \\ \mathbf{0} & \mathbf{A}_3 \end{bmatrix} \mathbf{K}_{ij}^3 \begin{bmatrix} \mathbf{A}_3 & \hat{\mathbf{e}}_{N_3} \hat{\mathbf{r}}_{ij} \mathbf{C}_3 \\ \mathbf{0} & \mathbf{A}_3 \end{bmatrix}^T + \begin{bmatrix} \hat{\mathbf{e}}_{N_3} \hat{\mathbf{r}}_{ij} \hat{\mathbf{e}}_1 \\ \hat{\mathbf{e}}_{N_3} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{e}}_{N_3} \hat{\mathbf{r}}_{ij} \hat{\mathbf{e}}_1 \\ \hat{\mathbf{e}}_{N_3} \end{bmatrix}^T. \end{aligned} \quad (14b)$$

III. REALIZATION WITH LOW SENSITIVITY

Transfer function $H(z_1, z_2, z_3)$ in (2a) consists of three 1-D factors. As shown in (5a) and (5b), realizations of the first and last factors, i.e. $f_1(z_1)$ and $g_3(z_3)$, possess rather

simple structures as they involve only N_1 and N_3 independent parameters, respectively, in addition to the fixed ones and zeros. We note, however, that if coordinate transformations would be applied to each of these two sub-transfer functions so as to reduce coefficient sensitivity, the resulting state-space realizations would in general involve as many as $2[N_1(N_1 + 1) + N_3(N_3 + 1)]$ parameters, which implies a drastic increase in realization complexity. On the contrary, the sub-transfer function $\mathbf{H}_2(z_2)$ contains most of the filter's nontrivial coefficients and its state-space realization $(\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \mathbf{\Delta}_0)_p$ hardly contains many fixed trivial components like ones and zeros (see (A.4) in Appendix I), therefore, application of a coordinate transformation to $(\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \mathbf{\Delta}_0)_p$ will not increase realization complexity in a significant manner. For these reasons, we in the rest of the paper shall seek to minimize the sensitivity associated with matrix transfer function $\mathbf{H}_2(z_2)$.

Applying a coordinate transformation $\bar{\mathbf{x}}(k) = \mathbf{T}^{-1}\mathbf{x}(k)$ to the linear system $(\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \mathbf{\Delta}_0)_p$ in (3), we obtain a new realization $(\bar{\mathbf{A}}_2, \bar{\mathbf{B}}_2, \bar{\mathbf{C}}_2, \mathbf{\Delta}_0)_p$ characterized by

$$\bar{\mathbf{A}}_2 = \mathbf{T}^{-1}\mathbf{A}_2\mathbf{T}, \quad \bar{\mathbf{B}}_2 = \mathbf{T}^{-1}\mathbf{B}_2, \quad \bar{\mathbf{C}}_2 = \mathbf{C}_2\mathbf{T} \quad (15)$$

where \mathbf{T} is a $p \times p$ nonsingular matrix. For the new realization, the l_2 -sensitivity measure in (14a) is written as

$$S(\mathbf{P}) = J(\mathbf{P}) + \text{tr}[\mathbf{N}_{\Delta_0}] + \text{tr}[\mathbf{W}_1] + \text{tr}[\mathbf{K}_3] \quad (16a)$$

where

$$J(\mathbf{P}) = \text{tr}[\mathbf{M}_A(\mathbf{P})\mathbf{P}] + \text{tr}[\mathbf{W}_B\mathbf{P}] + \text{tr}[\mathbf{K}_C\mathbf{P}^{-1}] \quad (16b)$$

with $\mathbf{P} = \mathbf{T}\mathbf{T}^T$. For the reasons stated above, our attention shall now be focused on the minimization of sensitivity measure $J(\mathbf{P})$.

Since $\mathbf{f}(z_2, z_3)$ is the transfer function from the filter input to the state-variable vector $\mathbf{x}(k)$, a controllability Gramian \mathbf{K} is defined by

$$\mathbf{K} = \frac{1}{(2\pi j)^2} \oint_{|z_2|=1} \oint_{|z_3|=1} \mathbf{f}(z_2, z_3) \mathbf{f}^*(z_2, z_3) \frac{dz_2}{z_2} \frac{dz_3}{z_3} \quad (17)$$

that can be obtained by solving the following Lyapunov equations:

$$\begin{aligned} \mathbf{K}_0 &= \mathbf{A}_3\mathbf{K}_0\mathbf{A}_3^T + \hat{\mathbf{e}}_{N_3}\hat{\mathbf{e}}_{N_3}^T \\ \mathbf{K} &= \mathbf{A}_2\mathbf{K}\mathbf{A}_2^T + \mathbf{B}_2(\mathbf{C}_3\mathbf{K}_0\mathbf{C}_3^T + \hat{\mathbf{e}}_1\hat{\mathbf{e}}_1^T)\mathbf{B}_2^T. \end{aligned} \quad (18)$$

For the new realization $(\bar{\mathbf{A}}_2, \bar{\mathbf{B}}_2, \bar{\mathbf{C}}_2, \mathbf{\Delta}_0)_p$, l_2 -scaling constraints on the state-variable vector $\bar{\mathbf{x}}(k)$ are given by

$$(\bar{\mathbf{K}})_{ii} = (\mathbf{T}^{-1}\mathbf{K}\mathbf{T}^{-T})_{ii} = 1 \quad \text{for } i = 1, 2, \dots, p. \quad (19)$$

The reader is referred to Appendix III that explains (19) in details. Summarizing, our sensitivity minimization problem is to find a coordinate transformation matrix \mathbf{T} that minimizes $J(\mathbf{P})$ in (16b) subject to the l_2 -scaling constraints in (19). Two solution methods for this problem are presented below.

We stress that apart from the difference in technical details in implementing these methods, they are analogous to their 1-D and 2-D counterparts investigated in [10]–[11].

A. A Constrained Optimization Method

Directly dealing with the l_2 -scaling constraints in (19) was found technically unfeasible. Instead, we consider a relaxation of (19) to a single constraint as

$$\text{tr}[\mathbf{T}^{-1}\mathbf{K}\mathbf{T}^{-T}] = \text{tr}[\mathbf{K}\mathbf{P}^{-1}] = p. \quad (20)$$

The relaxation proposed in (20) for the constraints in (19) has two advantages: First, rather than dealing with a total of p constraints in (19), with (20) we have only one single equality constraint that, as will be demonstrated below, can be handled by a Lagrange function with one additional parameter; second, from (20) it is seen that the constraint is now expressed in terms of matrix \mathbf{P} , therefore, once an optimal \mathbf{P} is identified, an optimal coordinate transformation matrix \mathbf{T} can be set to

$$\mathbf{T} = \mathbf{P}^{\frac{1}{2}}\mathbf{U} \quad (21)$$

with \mathbf{U} a $p \times p$ orthogonal matrix. Note that changing \mathbf{U} in (21) to a different orthogonal matrix does not affect the optimality of \mathbf{T} in (21) in the sense that, for any orthogonal \mathbf{U} , $\mathbf{P} = \mathbf{T}\mathbf{T}^T$ remains valid, hence it does not alter the optimal value of the sensitivity measure in (16b). Further notice that with (21), (19) becomes

$$(\mathbf{U}^T\mathbf{P}^{-\frac{1}{2}}\mathbf{K}\mathbf{P}^{-\frac{T}{2}}\mathbf{U})_{ii} = 1 \quad \text{for } i = 1, 2, \dots, p. \quad (22)$$

With \mathbf{K} fixed and \mathbf{P} determined, it is straightforward to find an orthogonal matrix \mathbf{U} so that the constraints in (22) (hence (19)) are satisfied. In words, via (21) a solution of the relaxed problem can be readily concerted to a solution that accurately satisfies the l_2 -scaling constraints in (19). The reader is referred to Section IV.A for a numerical example that illustrates the technique described above. This justifies the relaxation made in (20) and in this way, we now focus on the problem

$$\begin{aligned} &\text{minimize } J(\mathbf{P}) \text{ in (16b)} \\ &\text{subject to } \text{tr}[\mathbf{K}\mathbf{P}^{-1}] = p. \end{aligned} \quad (23)$$

To solve problem (23), we define the following Lagrange function of the problem:

$$\begin{aligned} J_o(\mathbf{P}, \lambda) &= \text{tr}[\mathbf{M}_A(\mathbf{P})\mathbf{P}] + \text{tr}[\mathbf{W}_B\mathbf{P}] \\ &\quad + \text{tr}[\mathbf{K}_C\mathbf{P}^{-1}] + \lambda(\text{tr}[\mathbf{K}\mathbf{P}^{-1}] - p) \end{aligned} \quad (24)$$

where λ is the Lagrange multiplier. Computing $J_o(\mathbf{P}, \lambda)/\partial\mathbf{P}$ by using [28, p.275]

$$\frac{d[\text{tr}(\mathbf{M}\mathbf{X})]}{d\mathbf{X}} = \mathbf{M}^T, \quad \frac{d[\text{tr}(\mathbf{M}\mathbf{X}^{-1})]}{d\mathbf{X}} = -(\mathbf{X}^{-1}\mathbf{M}\mathbf{X}^{-1})^T \quad (25)$$

and setting $\partial J_o(\mathbf{P}, \lambda)/\partial\mathbf{P} = \mathbf{0}$, it follows that

$$\mathbf{P}\mathbf{F}(\mathbf{P})\mathbf{P} = \mathbf{G}(\mathbf{P}, \lambda) \quad (26a)$$

where

$$\begin{aligned} \mathbf{F}(\mathbf{P}) &= \mathbf{M}_A(\mathbf{P}) + \mathbf{W}_B \\ \mathbf{G}(\mathbf{P}, \lambda) &= \mathbf{N}_A(\mathbf{P}) + \mathbf{K}_C + \lambda\mathbf{K} \end{aligned} \quad (26b)$$

with

$$\begin{aligned} \mathbf{N}_A(\mathbf{P}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} \mathbf{I}_p & \mathbf{0} \end{bmatrix} \mathbf{N}_{ij}^A \begin{bmatrix} \mathbf{I}_p \\ \mathbf{0} \end{bmatrix} \\ \mathbf{N}_{ij}^A &= \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2\mathbf{R}_{ij}\mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix} \mathbf{N}_{ij}^A \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2\mathbf{R}_{ij}\mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}^T \\ &\quad + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{P} \end{bmatrix}. \end{aligned} \quad (26c)$$

The equation in (26a) is highly nonlinear with respect to \mathbf{P} . An effective approach to solving the equation in (26a) is to relax it into the following recursive second-order matrix equation:

$$\mathbf{P}^{(k+1)} \mathbf{F}(\mathbf{P}^{(k)}) \mathbf{P}^{(k+1)} = \mathbf{G}(\mathbf{P}^{(k)}, \lambda^{(k+1)}) \quad (27)$$

where $\mathbf{P}^{(k)}$ is assumed to be known from the previous recursion and the solution $\mathbf{P}^{(k+1)}$ is given by

$$\begin{aligned} \mathbf{P}^{(k+1)} &= \mathbf{F}(\mathbf{P}^{(k)})^{-\frac{1}{2}} [\mathbf{F}(\mathbf{P}^{(k)})^{\frac{1}{2}} \\ &\quad \cdot \mathbf{G}(\mathbf{P}^{(k)}, \lambda^{(k+1)}) \mathbf{F}(\mathbf{P}^{(k)})^{\frac{1}{2}}]^{-\frac{1}{2}} \mathbf{F}(\mathbf{P}^{(k)})^{-\frac{1}{2}}. \end{aligned} \quad (28)$$

Here, the Lagrange multiplier $\lambda^{(k+1)}$ can be efficiently obtained using a bisection method [29] so that

$$f(\lambda^{(k+1)}) = |p - \text{tr}[\tilde{\mathbf{K}}^{(k)} \tilde{\mathbf{G}}^{(k)}(\lambda^{(k+1)})]| < \varepsilon \quad (29a)$$

where

$$\begin{aligned} \tilde{\mathbf{G}}^{(k)}(\lambda^{(k+1)}) &= [\mathbf{F}(\mathbf{P}^{(k)})^{\frac{1}{2}} \mathbf{G}(\mathbf{P}^{(k)}, \lambda^{(k+1)}) \mathbf{F}(\mathbf{P}^{(k)})^{\frac{1}{2}}]^{-\frac{1}{2}} \\ \tilde{\mathbf{K}}^{(k)} &= \mathbf{F}(\mathbf{P}^{(k)})^{\frac{1}{2}} \mathbf{K} \mathbf{F}(\mathbf{P}^{(k)})^{\frac{1}{2}}. \end{aligned} \quad (29b)$$

This iteration process continues until

$$|J(\mathbf{P}^{(k)}, \lambda^{(k+1)}) - J(\mathbf{P}^{(k-1)}, \lambda^{(k)})| < \varepsilon \quad (30)$$

is satisfied for a prescribed tolerance $\varepsilon > 0$. If the iteration is terminated at step k , $\mathbf{P}^{(k)}$ is claimed to be a solution point.

B. An Unconstrained Optimization Method

Defining

$$\hat{\mathbf{T}} = \mathbf{T}^T \mathbf{K}^{-\frac{1}{2}} \quad (31)$$

the l_2 -scaling constraints in (19) can be written as

$$(\hat{\mathbf{T}}^{-T} \hat{\mathbf{T}}^{-1})_{ii} = 1 \quad \text{for } i = 1, 2, \dots, p. \quad (32)$$

It is obvious that the conditions in (32) are always satisfied by choosing $\hat{\mathbf{T}}^{-1}$ as

$$\hat{\mathbf{T}}^{-1} = \left[\frac{t_1}{\|t_1\|}, \frac{t_2}{\|t_2\|}, \dots, \frac{t_p}{\|t_p\|} \right]. \quad (33)$$

By writing $J(\mathbf{P})$ in (16b) as

$$J(\mathbf{T}) = \text{tr}[\mathbf{T}^T \mathbf{M}_A (\mathbf{T} \mathbf{T}^T) \mathbf{T}] + \text{tr}[\mathbf{T}^T \mathbf{W}_B \mathbf{T}] + \text{tr}[\mathbf{T}^{-1} \mathbf{K}_C \mathbf{T}^{-T}] \quad (34)$$

and then using (31), the l_2 -sensitivity measure can be expressed as

$$J(\mathbf{x}) = \text{tr}[\hat{\mathbf{T}} \hat{\mathbf{M}}_A (\hat{\mathbf{T}} \hat{\mathbf{T}}^T) \hat{\mathbf{T}}] + \text{tr}[\hat{\mathbf{T}} \hat{\mathbf{W}}_B \hat{\mathbf{T}}^T] + \text{tr}[\hat{\mathbf{T}}^{-T} \hat{\mathbf{K}}_C \hat{\mathbf{T}}^{-1}] \quad (35a)$$

where $\mathbf{x} = (t_1^T, t_2^T, \dots, t_p^T)^T$ and

$$\begin{aligned} \hat{\mathbf{M}}_A(\hat{\mathbf{T}}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} \mathbf{0} & \mathbf{K}^{\frac{1}{2}} \end{bmatrix} \hat{\mathbf{M}}_{ij}^A \begin{bmatrix} \mathbf{0} \\ \mathbf{K}^{\frac{1}{2}} \end{bmatrix} \\ \hat{\mathbf{M}}_{ij}^A &= \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}^T \hat{\mathbf{M}}_{ij}^A \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix} \\ &\quad + \begin{bmatrix} \mathbf{K}^{-\frac{1}{2}} \hat{\mathbf{T}}^{-1} \hat{\mathbf{T}}^{-T} \mathbf{K}^{-\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \hat{\mathbf{W}}_B &= \mathbf{K}^{\frac{1}{2}} \mathbf{W}_B \mathbf{K}^{\frac{1}{2}}, \quad \hat{\mathbf{K}}_C = \mathbf{K}^{-\frac{1}{2}} \mathbf{K}_C \mathbf{K}^{-\frac{1}{2}}. \end{aligned} \quad (35b)$$

In this way, the original constrained optimization problem can be converted into an unconstrained optimization problem of obtaining a $p^2 \times 1$ vector \mathbf{x} which minimizes $J(\mathbf{x})$ in (35a).

Applying a quasi-Newton (known as the Broyden-Fletcher-Goldfarb-Shanno (BFGS)) algorithm [30] to minimize $J(\mathbf{x})$ in (35a), in the k th iteration the most recent point \mathbf{x}_k is updated to point \mathbf{x}_{k+1} as

$$\mathbf{x}_{k+1} = \mathbf{x} + \alpha_k \mathbf{d}_k \quad (36a)$$

where

$$\begin{aligned} \mathbf{d}_k &= -\mathbf{S}_k \nabla J(\mathbf{x}_k), \quad \alpha_k = \arg \min_{\alpha} J(\mathbf{x}_k + \alpha \mathbf{d}_k) \\ \mathbf{S}_{k+1} &= \mathbf{S}_k + \left(1 + \frac{\gamma_k^T \mathbf{S}_k \gamma_k}{\gamma_k^T \delta_k} \right) \frac{\delta_k \delta_k^T}{\gamma_k^T \delta_k} - \frac{\delta_k \gamma_k^T \mathbf{S}_k + \mathbf{S}_k \gamma_k \delta_k^T}{\gamma_k^T \delta_k} \\ \mathbf{S}_0 &= \mathbf{I}_{p^2}, \quad \delta_k = \mathbf{x}_{k+1} - \mathbf{x}_k, \quad \gamma_k = \nabla J(\mathbf{x}_{k+1}) - \nabla J(\mathbf{x}_k). \end{aligned} \quad (36b)$$

In the above, $\nabla J(\mathbf{x})$ is the gradient of $J(\mathbf{x})$ with respect to \mathbf{x} , and \mathbf{S}_k is a positive-definite approximation of the inverse Hessian matrix of $J(\mathbf{x})$. The algorithm starts with a trivial initial point \mathbf{x}_0 obtained from an initial assignment $\mathbf{T} = \mathbf{I}_p$, and this iteration process continues until

$$|J(\mathbf{x}_{k+1}) - J(\mathbf{x}_k)| < \epsilon \quad (37)$$

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

The BFGS algorithm is a well-known descent algorithm meaning that a twice continuously differentiable objective function at the iterates generated by the BFGS algorithm is monotonically decreasing. In theory, it was shown [31] that if the starting point is sufficiently close to a local minimizer and the initial Hessian approximation is sufficiently close to the Hessian at that minimizer, then the BFGS iterates will converge to the minimizer.

The implementation of (36a) requires the gradient of $J(\mathbf{x})$, which can be efficiently evaluated using closed-form expressions, see Appendix IV.

IV. A CASE STUDY

We now present a case study to demonstrate the effectiveness of the two algorithms developed in Section III. The case study was carried out using MATLAB on a PC with an Intel Core i-2500 CPU at 3.3 GHz.

Consider a stable 3-D separable-denominator digital filter in (2a) and (2b) specified by

$$\begin{aligned} \Delta_0 &= 10^{-2} \begin{bmatrix} 0.00730 & 0.34297 & -0.09594 & 0.20541 \\ 3.33408 & -5.73707 & 3.94939 & -1.61598 \\ -1.46081 & 2.66051 & -1.68094 & 0.68022 \\ 1.12651 & -1.62192 & 1.24735 & -0.55781 \end{bmatrix} \\ \Delta_1 &= 10^{-2} \begin{bmatrix} 2.81318 & -5.00467 & 3.46926 & -0.84798 \\ -5.29980 & 9.24831 & -6.29206 & 2.80791 \\ 4.95232 & -8.39641 & 5.73329 & -1.62170 \\ 0.72029 & -1.34272 & 0.95941 & 0.54827 \end{bmatrix} \\ \Delta_2 &= 10^{-2} \begin{bmatrix} -0.69409 & 1.54874 & -0.94779 & 0.39116 \\ 3.93785 & -6.79910 & 4.66564 & -1.96344 \\ -2.37995 & 4.20737 & -2.75482 & 0.95329 \\ 0.70545 & -0.90615 & 0.73168 & -0.55633 \end{bmatrix} \end{aligned}$$

$$\Delta_3 = 10^{-2} \begin{bmatrix} 1.67681 & -2.69078 & 1.98218 & -0.33567 \\ -0.59937 & 1.11289 & -0.71981 & 0.43504 \\ 1.82472 & -2.93685 & 2.11591 & -0.43417 \\ 1.28875 & -2.01749 & 1.51782 & -0.09016 \end{bmatrix}$$

$$\begin{bmatrix} b_{11} & b_{12} & b_{13} \end{bmatrix} = \begin{bmatrix} b_{31} & b_{32} & b_{33} \end{bmatrix}$$

$$= \begin{bmatrix} -1.81600 & 1.23756 & -0.31382 \end{bmatrix}$$

$$\begin{bmatrix} b_{21} & b_{22} & b_{23} \end{bmatrix} = \begin{bmatrix} -1.81611 & 1.23775 & -0.31391 \end{bmatrix}.$$

The 3-D filter can be realized by a minimal state-space model in (3) with

$$A_2 = \begin{bmatrix} 0.00000 & -0.19089 & 0.29060 \\ 0.74393 & -86.40470 & 133.71075 \\ -0.27211 & -57.01643 & 88.22081 \end{bmatrix}$$

$$B_2 = 10^3 \begin{bmatrix} 0.00602 & -0.00921 & 0.00699 & -0.00095 \\ -1.10247 & 1.68622 & -1.27902 & 0.17267 \\ -0.71455 & 1.09291 & -0.82977 & 0.11192 \end{bmatrix}$$

$$C_2 = \begin{bmatrix} 0.07236 & 0.06711 & -0.10298 \\ 0.01930 & 0.01789 & -0.02745 \\ 0.05887 & 0.05460 & -0.08378 \\ 0.07079 & 0.06565 & -0.10073 \end{bmatrix}.$$

Applying a coordinate transformation defined by

$$T = 10^3 \text{diag}\{0.01077, 2.58588, 1.68384\}$$

and using (14b) with truncation $(0,0) \leq (i,j) \leq (100,100)$ and (18) to evaluate the Gramians $M_A(I_3)$, W_B , K_C , N_{Δ_0} , W_1 , K_3 , K_0 and K yields

$$M_A(I_3) = 10^7 \begin{bmatrix} 0.00021 & 0.03478 & -0.03471 \\ 0.03478 & 5.83540 & -5.82400 \\ -0.03471 & -5.82400 & 5.81261 \end{bmatrix}$$

$$W_B = 10^8 \begin{bmatrix} 0.00014 & 0.02404 & -0.02399 \\ 0.02404 & 4.36159 & -4.35410 \\ -0.02399 & -4.35410 & 4.34663 \end{bmatrix}$$

$$K_C = 10 \begin{bmatrix} 5.70413 & -4.79529 & -4.78664 \\ -4.79529 & 5.70413 & 5.70410 \\ -4.78664 & 5.70410 & 5.70413 \end{bmatrix}$$

$$N_{\Delta_0} = 10 \begin{bmatrix} 8.61482 & 8.61482 & 8.61482 & 8.61482 \\ 8.61482 & 8.61482 & 8.61482 & 8.61482 \\ 8.61482 & 8.61482 & 8.61482 & 8.61482 \\ 8.61482 & 8.61482 & 8.61482 & 8.61482 \end{bmatrix}$$

$$W_1 = 10^3 \begin{bmatrix} 2.68584 & 2.59589 & 2.34321 \\ 2.59589 & 2.68584 & 2.59589 \\ 2.34321 & 2.59589 & 2.68584 \end{bmatrix}$$

$$K_3 = 10^3 \begin{bmatrix} 2.52585 & 2.43584 & 2.18528 \\ 2.43584 & 2.52585 & 2.43584 \\ 2.18528 & 2.43584 & 2.52585 \end{bmatrix}$$

$$K_0 = 10 \begin{bmatrix} 1.42603 & 1.29740 & 0.99842 \\ 1.29740 & 1.42603 & 1.29740 \\ 0.99842 & 1.29740 & 1.42603 \end{bmatrix}$$

$$K = \begin{bmatrix} 1.00000 & -0.84067 & -0.83915 \\ -0.84067 & 1.00000 & 0.99999 \\ -0.83915 & 0.99999 & 1.00000 \end{bmatrix}.$$

The total l_2 -sensitivity defined by (14a) was found to be

$$S = J(I_3) + 1.59797 \times 10^4$$

where $J(I_3)$ defined by (16b) was found to be

$$J(I_3) = 9.87319 \times 10^8.$$

We see that the sensitivity associated with $H_2(z_2)$, namely $J(I_3)$, was indeed dominating the filter's sensitivity.

A. Application of the Lagrange method

Choosing $P^{(0)} = I_3$ in (28) as an initial estimate and a tolerance $\epsilon = 10^{-8}$ in the bisection method of (29a) as well as in (30), it took the Lagrange-based algorithm 7 iterations to converge to the solution

$$P^{opt} = \begin{bmatrix} 2.26693 & -2.29606 & -2.28843 \\ -2.29606 & 3.27352 & 3.26762 \\ -2.28843 & 3.26762 & 3.26176 \end{bmatrix}.$$

With K given above and $P = P^{opt}$, an orthogonal matrix U satisfying (22) was found to be

$$U = \begin{bmatrix} 0.91281 & 0.01760 & -0.40801 \\ 0.37727 & 0.34617 & 0.85897 \\ -0.15636 & 0.93801 & -0.30934 \end{bmatrix}$$

which in conjunction with (21) yielded the optimal coordinate transformation matrix as

$$T^{opt} = \begin{bmatrix} 0.95919 & -0.79515 & -0.84535 \\ -0.31867 & 1.52484 & 0.92024 \\ -0.31499 & 1.52445 & 0.91574 \end{bmatrix}$$

and the optimal state-space coefficient matrices in (15) were constructed as

$$\bar{A}_2 = \begin{bmatrix} 0.62960 & -0.10557 & -0.29017 \\ -0.03842 & 0.56172 & -0.29573 \\ 0.40269 & 0.09232 & 0.62479 \end{bmatrix}$$

$$\bar{B}_2 = \begin{bmatrix} 0.36268 & -0.55634 & 0.50511 & -0.05885 \\ -0.12255 & 0.18883 & -0.23407 & 0.02035 \\ -0.13464 & 0.20306 & 0.02528 & 0.01847 \end{bmatrix}$$

$$\bar{C}_2 = \begin{bmatrix} 0.06690 & -0.34372 & 0.24644 \\ 0.01682 & -0.08626 & 0.06888 \\ 0.05258 & -0.27075 & 0.20585 \\ 0.06038 & -0.31154 & 0.25590 \end{bmatrix}.$$

The controllability Gramian associated with the optimized state-space model was found to be

$$\bar{K} = \begin{bmatrix} 1.00041 & -0.63416 & 0.63416 \\ -0.63416 & 1.00041 & -0.89767 \\ 0.63416 & -0.89767 & 1.00041 \end{bmatrix}$$

showing that the l_2 -scaling constraints in (19) were practically satisfied. Then the l_2 -sensitivity measure in (24) was computed as

$$J_o(P^{opt}, \lambda) = 3.24252 \times 10^3$$

where $\lambda = 8.42179 \times 10^2$. The profiles of the l_2 -sensitivity measure $J_o(P, \lambda)$ and the Lagrange multiplier λ during the first 7 iterations are shown in Fig. 2, from which it is observed that with a tolerance $\epsilon = 10^{-8}$ the algorithm converges in 7 iterations.

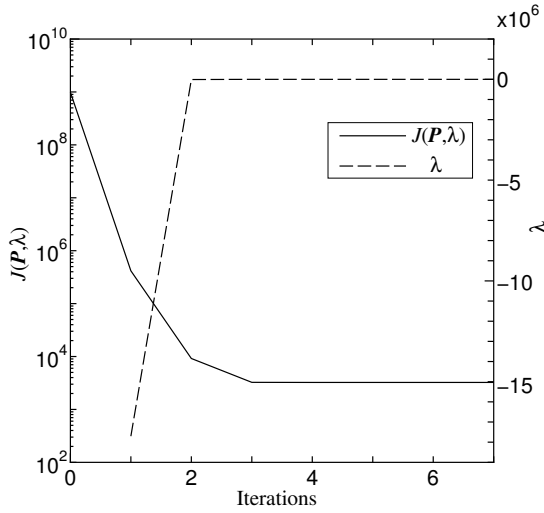


Fig. 2. Profiles of l_2 -sensitivity and λ performances during the first 7 iterations of the constrained optimization method.

B. Application of the quasi-Newton method

By choosing $\hat{T} = K^{-\frac{1}{2}}$ with $T = I_3$ as an initial estimate in (36a) and a tolerance $\epsilon = 10^{-8}$ in (37), the quasi-Newton algorithm took 37 iterations to converge to the solution

$$\hat{T}^{opt} = \begin{bmatrix} 1.03339 & 1.41754 & -0.13472 \\ -0.29813 & 2.46592 & -0.04995 \\ 0.90070 & -1.17579 & 0.96141 \end{bmatrix}$$

or equivalently,

$$T^{opt} = \begin{bmatrix} 0.34976 & -1.20354 & 0.83486 \\ 0.42464 & 1.68766 & -0.49626 \\ 0.42613 & 1.68528 & -0.49116 \end{bmatrix}.$$

The l_2 -sensitivity measure in (35a) was then computed as

$$J(x) = 3.24356 \times 10^3$$

and the optimal state-space coefficient matrices in (15) were constructed as

$$\begin{aligned} \bar{A}_2 &= \begin{bmatrix} 0.61448 & -0.17729 & 0.32325 \\ -0.03747 & 0.57287 & 0.17307 \\ -0.43314 & -0.03545 & 0.62876 \end{bmatrix} \\ \bar{B}_2 &= \begin{bmatrix} 0.24441 & -0.37603 & 0.40596 & -0.04083 \\ -0.25567 & 0.39234 & -0.37591 & 0.04130 \\ 0.19877 & -0.30151 & 0.06567 & -0.02906 \end{bmatrix} \\ \bar{C}_2 &= \begin{bmatrix} 0.07181 & -0.29561 & -0.30148 \\ 0.02071 & -0.07290 & -0.08203 \\ 0.06091 & -0.23067 & -0.24815 \\ 0.07718 & -0.26188 & -0.30285 \end{bmatrix}. \end{aligned}$$

The profile of the l_2 -sensitivity measure $J(x)$ during the first 37 iterations is shown in Fig. 3, from which it is observed that with a tolerance $\epsilon = 10^{-8}$ the algorithm converges in 37 iterations.

We see that as far as the example discrete system is concerned, the two algorithms offer practically the same performance in terms of sensitivity minimization. Concerning

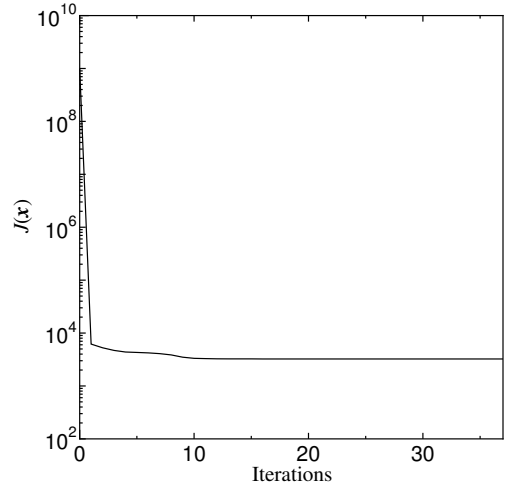


Fig. 3. Profile of l_2 -sensitivity performance during the first 37 iterations of the unconstrained optimization method.

the computational complexity, the 7 iterations required by the constrained algorithm consumed 46.01250 seconds of CPU time versus 7.15158×10^3 seconds for the unconstrained algorithm to finish its 37 iterations. Obviously the constrained method is considerably more efficient than its unconstrained counterpart. However, for discrete system of large scales, it is quite likely that the algorithms' complexity will change due to the fact that the constrained algorithm heavily involves expensive matrix manipulations such as computing square root as well as inverse of matrices. For this reason, we believe that together the two algorithms can serve for sensitivity minimization for a wide range of system scales.

C. Comparisons with the work in [25]

The Lagrange method was also applied to the 3-D filter investigated in Section 4 of [25]. We shall omit the intermediate details but focus on comparisons of the results obtained. The total l_2 -sensitivity defined by (14a) with (16b) after l_2 -scaling was found to be $S = J(I_3) + 21.89505 \times 10^2 = 112.71349 \times 10^2$ where $J(I_3) = 90.81844 \times 10^2$. After 6 iterations, the Lagrange method converges to a solution P^{opt} with $J(P^{opt}) = 9.14487 \times 10^2$, which implies a total l_2 -sensitivity $S = 31.03992 \times 10^2$. At this solution, the dynamic-range constraints $(\bar{K})_{ii} = 1$ for $i = 1, 2, 3$ are perfectly satisfied. The total number of nontrivial parameters in its realization was 55. Alternatively, the algorithm proposed in [25] was applied, which led to a minimized total l_2 -sensitivity 3.12274×10^2 . The Gramian \bar{K} associated with the optimal state-space filter was found to be

$$\bar{K} = \begin{bmatrix} 1.76972 & 1.14927 & 0.77867 \\ 1.14927 & 1.21704 & 0.94657 \\ 0.77867 & 0.94657 & 2.25271 \end{bmatrix}.$$

Clearly, the dynamic-range constraints $(\bar{K})_{ii} = 1$ are violated for all $i = 1, 2, 3$ in this case. The number of nontrivial parameters involved in this realization was 73. From above results, we see two options for 3-D state-space filters with separable

denominators: [25] offers techniques for minimum sensitivity but with more system parameters and possible violation of dynamic-range constraints, while the present paper provides techniques for low sensitivity, reduced number of parameters (hence faster implementation), and guaranteed freedom of internal overflow.

V. CONCLUSION

A minimal state-space realization technique for 3-D separable-denominator digital filters has been explored and the l_2 -sensitivity of the minimal state-space model realized from a given 3-D separable-denominator digital filter has been analyzed more precisely by taking into account 0 and 1 elements contained in the model. Two iterative methods have been developed to minimize the l_2 -sensitivity of the filter subject to a fixed number of nontrivial parameters and dynamic-range l_2 -scaling constraints. Computer simulation results have demonstrated the validity and effectiveness of the proposed techniques.

A problem that might be worthwhile for future studies is to minimize coefficient sensitivity subject to filter's robust stability under l_2 -scaling constraints. In such a problem, the results reported in [24] are expected to play a significant role.

APPENDIX I

MINIMAL REALIZATION OF $H_2(z_2)$

We write $H_2(z_2)$ in (2b) as

$$H_2(z_2) = \Delta_0 + \frac{(\Delta_1 - b_{21}\Delta_0)z_2^{-1} + \dots + (\Delta_{N_2} - b_{2N_2}\Delta_0)z_2^{-N_2}}{D_2(z_2)} \quad (\text{A.1})$$

which can be realized with a multivariable observable canonic form $(A_0, B_0, C_0, \Delta_0)_{N_2(N_1+1)}$ as

$$A_0 = \begin{bmatrix} \mathbf{0} & \cdots & \mathbf{0} & -b_{2N_2}\mathbf{I}_{N_1+1} \\ \mathbf{I}_{N_1+1} & \cdots & \mathbf{0} & \vdots \\ \vdots & \ddots & \vdots & -b_{22}\mathbf{I}_{N_1+1} \\ \mathbf{0} & \cdots & \mathbf{I}_{N_1+1} & -b_{21}\mathbf{I}_{N_1+1} \end{bmatrix}$$

$$B_0 = \begin{bmatrix} \Delta_{N_2} - b_{2N_2}\Delta_0 \\ \vdots \\ \Delta_2 - b_{22}\Delta_0 \\ \Delta_1 - b_{21}\Delta_0 \end{bmatrix}, \quad C_0 = [\mathbf{0} \quad \cdots \quad \mathbf{0} \quad \mathbf{I}_{N_1+1}]. \quad (\text{A.2})$$

Although this realization is always observable, it is uncontrollable unless $\text{rank } V_{N_2(N_1+1)-r} = N_2(N_1+1)$, i.e. full rank, where $V_i = [B_0, A_0 B_0, \dots, A_0^i B_0]$ and r is the rank of B_0 [27]. Suppose that

$$\text{rank } V_{N_2(N_1+1)-1} = p \quad (\text{A.3})$$

and v_1, v_2, \dots, v_p are p linearly independent column vectors from matrix $V_{N_2(N_1+1)-1}$. By defining $M = [v_1, v_2, \dots, v_p]$, it is shown [32] that $H_2(z_2)$ can be realized with a minimal state-space model $(A_2, B_2, C_2, \Delta_0)_p$ as

$$A_2 = (M^T M)^{-1} M^T A_0 M, \quad B_2 = (M^T M)^{-1} M^T B_0$$

$$C_2 = C_0 M. \quad (\text{A.4})$$

APPENDIX II

ERROR IN TRANSFER FUNCTION DUE TO PARAMETER VARIATIONS

Consider a transfer function $H(z_1, z_2, z_3)$ that contains N parameters $\{p_1, p_2, \dots, p_N\}$. Let $\{\tilde{p}_i\}$ be the FWL version of $\{p_i\}$, where $\tilde{p}_i = p_i + \Delta p_i$ with Δp_i the parameter perturbation, and $\tilde{H}(z_1, z_2, z_3)$ be the transfer function associated with perturbed parameters $\{\tilde{p}_i\}$. The first-order approximation of $\tilde{H}(z_1, z_2, z_3)$ then gives

$$\tilde{H}(z_1, z_2, z_3) = H(z_1, z_2, z_3) + \Delta H(z_1, z_2, z_3) \quad (\text{A.5a})$$

where $\Delta H(z_1, z_2, z_3)$ will be

$$\Delta H(z_1, z_2, z_3) = \sum_{i=1}^N \frac{\partial H(z_1, z_2, z_3)}{\partial p_i} \Delta p_i. \quad (\text{A.5b})$$

For a fixed-point implementation of B bits, the parameter perturbations are considered to be independent uniformly distributed random variables within the range $[-2^{-B-1}, 2^{-B-1}]$. Then, a measure of the transfer function error can statistically be defined as

$$\sigma_{\Delta H}^2 = \frac{1}{(2\pi j)^3} \oint_{|z_1|=1} \oint_{|z_2|=1} \oint_{|z_3|=1} E[|\Delta H(z_1, z_2, z_3)|^2] \cdot \frac{dz_1}{z_1} \frac{dz_2}{z_2} \frac{dz_3}{z_3} \quad (\text{A.6})$$

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

$$E[|\Delta H(z_1, z_2, z_3)|^2] = \sum_{i=1}^N \left| \frac{\partial H(z_1, z_2, z_3)}{\partial p_i} \right|^2 \sigma^2 \quad (\text{A.7})$$

where $\sigma^2 = E[(\Delta p_i)^2] = 2^{-2B}/12$. Eq. (A.7) establishes an analytic relationship between variations in the transfer function induced by an FWL realization and parameter sensitivity.

APPENDIX III

l_2 -SCALING CONSTRAINTS

Suppose that the input and the output of a 2-D filter $f(z_2, z_3)$ are denoted by $x(k, r)$ and $u(k, r)$, respectively. It follows from (9b) and (5a) that

$$x(k, r) = \sum_{p=1}^k \sum_{q=0}^r h(p, q) u(k-p, r-q) \quad (\text{A.8a})$$

where

$$h(p, q) = A_2^{p-1} B_2 C_3 A_3^{q-1} \hat{e}_{N_3} \quad \text{for } p, q \geq 1$$

$$h(p, 0) = A_2^{p-1} B_2 \hat{e}_1 \quad \text{for } p \geq 1. \quad (\text{A.8b})$$

Let e_i be the i th column of the identity matrix I_p , an upper bound of the i th component of the local state vector $x(k, r)$

in (A.8a) can be obtained as

$$\begin{aligned} |e_i^T \mathbf{x}(k, r)|^2 &= \left[\sum_{p=1}^k \sum_{q=0}^r e_i^T \mathbf{h}(p, q) u(k-p, r-q) \right]^2 \\ &\leq e_i^T \left(\sum_{p=1}^k \sum_{q=0}^r \mathbf{h}(p, q) \mathbf{h}^T(p, q) \right) e_i \\ &\quad \cdot \sum_{p=1}^k \sum_{q=0}^r u^2(k-p, r-q) \leq e_i^T \mathbf{K} e_i \|u\|^2 \end{aligned} \quad (\text{A.9a})$$

where

$$\mathbf{K} = \sum_{p=1}^{\infty} \sum_{q=0}^{\infty} \mathbf{h}(p, q) \mathbf{h}^T(p, q). \quad (\text{A.9b})$$

Note that $e_i^T \mathbf{K} e_i$ is the i th diagonal element of \mathbf{K} , thus if all the diagonal elements of \mathbf{K} are equal to unity, (A.9a) implies that the amplitude of each component of vector $\mathbf{x}(k, r)$ will not exceed $\|u\|$. Therefore, the dynamic-range constraints on the state-variables may be imposed as

$$(\mathbf{K})_{ii} = 1 \quad \text{for } i = 1, 2, \dots, p. \quad (\text{A.10})$$

The above constraints correspond to (19) in the new realization and matrix \mathbf{K} in (A.9b) coincides with that derived from (18).

APPENDIX IV

GRADIENT OF $J(\mathbf{x})$

The gradient of $J(\mathbf{x})$ can be evaluated using closed-form expressions as

$$\frac{\partial J(\hat{\mathbf{T}})}{\partial t_{ij}} = \lim_{\Delta \rightarrow 0} \frac{J(\hat{\mathbf{T}}_{ij}) - J(\hat{\mathbf{T}})}{\Delta} = 2\beta_1 - 2\beta_2 + 2\beta_3 - 2\beta_4 \quad (\text{A.11a})$$

where $\hat{\mathbf{T}}_{ij}$ is the matrix obtained from $\hat{\mathbf{T}}$ with a perturbed (i, j) th component, which is given by [33, p.655]

$$\begin{aligned} \hat{\mathbf{T}}_{ij} &= \hat{\mathbf{T}} + \frac{\Delta \hat{\mathbf{T}} \mathbf{g}_{ij} e_j^T \hat{\mathbf{T}}}{1 - \Delta e_j^T \hat{\mathbf{T}} \mathbf{g}_{ij}}, \quad \hat{\mathbf{T}}_{ij}^{-1} = \hat{\mathbf{T}}^{-1} - \Delta \mathbf{g}_{ij} e_j^T \\ \mathbf{g}_{ij} &= \partial \left\{ \frac{t_j}{\|t_j\|} \right\} / \partial t_{ij} = \frac{1}{\|t_j\|^3} (t_{ij} t_j - \|t_j\|^2 e_i) \\ \beta_1 &= e_j^T \hat{\mathbf{T}} \mathbf{M}_A(\hat{\mathbf{T}}) \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij}, \quad \beta_3 = e_j^T \hat{\mathbf{T}} \mathbf{W}_B \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{g}_{ij} \\ \beta_2 &= e_j^T \hat{\mathbf{T}}^{-T} \hat{\mathbf{N}}_A(\hat{\mathbf{T}}) \mathbf{g}_{ij}, \quad \beta_4 = e_j^T \hat{\mathbf{T}}^{-T} \hat{\mathbf{K}}_C \mathbf{g}_{ij} \\ \hat{\mathbf{N}}_A(\hat{\mathbf{T}}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \begin{bmatrix} \mathbf{K}^{-\frac{1}{2}} & \mathbf{0} \end{bmatrix} \hat{\mathbf{N}}_{ij}^A \begin{bmatrix} \mathbf{K}^{-\frac{1}{2}} \\ \mathbf{0} \end{bmatrix} \\ \hat{\mathbf{N}}_{ij}^A &= \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix} \hat{\mathbf{N}}_{ij}^A \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \mathbf{R}_{ij} \mathbf{C}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}^T \\ &\quad + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}^{\frac{1}{2}} \hat{\mathbf{T}}^T \hat{\mathbf{T}} \mathbf{K}^{\frac{1}{2}} \end{bmatrix}. \end{aligned} \quad (\text{A.11b})$$

REFERENCES

- [1] 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.
- [2] G. Li and M. Gevers, "Optimal synthetic FWL design of state-space digital filters," in *Proc. 1992 IEEE Int. Conf. Acoust., Speech, Signal Processing*, vol. 4, pp. 429-432.
- [3] M. Gevers and G. Li, *Parameterizations in Control, Estimation and Filtering Problems: Accuracy Aspects*, Springer-Verlag, 1993.
- [4] C. Xiao, "Improved L_2 -sensitivity for state-space digital system," *IEEE Trans. Signal Processing*, vol. 45, pp. 837-840, Apr. 1997.
- [5] T. Hinamoto, S. Yokoyama, T. Inoue, W. Zeng and W.-S. Lu, "Analysis and minimization of L_2 -sensitivity for linear systems and two-dimensional state-space filters using general controllability and observability Gramians," *IEEE Trans. Circuits Syst. I*, vol. 49, pp. 1279-1289, Sept. 2002.
- [6] S. Yamaki, M. Abe and M. Kawamata, "A closed form solution to L_2 -sensitivity minimization of second-order state-space digital filters," in *Proc. 2006 IEEE Int. Symp. Circuits Syst.*, pp. 5223-5226.
- [7] 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.
- [8] G. Li, "Two-dimensional system optimal realizations with L_2 -sensitivity minimization," *IEEE Trans. Signal Processing*, vol. 46, pp. 809-813, Mar. 1998.
- [9] T. Hinamoto and Y. Sugie, " L_2 -sensitivity analysis and minimization of 2-D separable-denominator state-space digital filters," *IEEE Trans. Signal Processing*, vol. 50, pp. 3107-3114, Dec. 2002.
- [10] T. Hinamoto, H. Ohnishi and W.-S. Lu, "Minimization of L_2 -sensitivity for state-space digital filters subject to L_2 -dynamic-range scaling constraints," *IEEE Trans. Circuits Syst.-II*, vol. 52, pp. 641-645, Oct. 2005.
- [11] T. Hinamoto, K. Iwata and W.-S. Lu, " L_2 -sensitivity Minimization of one- and two-dimensional state-space digital filters subject to L_2 -scaling constraints," *IEEE Trans. Signal Processing*, vol. 54, pp. 1804-1812, May 2006.
- [12] S. Yamaki, M. Abe and M. Kawamata, "A novel approach to L_2 -sensitivity minimization of digital filters subject to L_2 -scaling constraints," in *Proc. 2006 IEEE Int. Symp. Circuits Syst.*, pp. 5219-5222.
- [13] T. Hinamoto, T. Oumi, O. I. Omoifo, and W.-S. Lu, "Minimization of frequency-weighted l_2 -sensitivity subject to l_2 -scaling constraints for two-dimensional state-space filters," *IEEE Trans. Signal Processing*, vol. 56, pp. 5157-5168, Oct. 2008.
- [14] C. T. Mullis and R. A. Roberts, "Synthesis of minimum roundoff noise fixed-point digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-23, pp. 551-562, Sept. 1976.
- [15] S. Y. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-25, pp. 273-281, Aug. 1977.
- [16] K. Hirano, M. Sakane and M. Z. Mulk, "Design of three-dimensional recursive digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-31, pp. 550-561, June 1984.
- [17] H. Mutluay and M. M. Fahmy, "Frequency-domain design of N -D digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-32, pp. 1226-1233, Dec. 1985.
- [18] M. E. Zervakis and A. N. Venetsanopoulos, "Design of three-dimensional digital filters using two-dimensional rotated filters," *IEEE Trans. Circuits Syst.*, vol. CAS-34, pp. 1452-1469, Dec. 1987.
- [19] H. K. Kwan and C. L. Chan, "Design of multidimensional spherically symmetric and constant group delay recursive digital filters with sum of power-of-two coefficients," *IEEE Trans. Circuits Syst.*, vol. CAS-37, pp. 1027-1035, Aug. 1990.
- [20] P. K. Rajan and M. N. S. Swamy, "Quadrantal symmetry associated with two-dimensional digital transfer functions," *IEEE Trans. Circuits Syst.*, vol. 25, pp. 340-343, June 1983.
- [21] L. Xu, H. Fan, Z. Lin and N. K. Bose, "A direct-construction approach to multidimensional realization and LFR uncertainty modeling," *Multidimensional Systems and Signal Processing*, vol. 19, pp. 323-359, Dec. 2008.
- [22] B. Sumanasena and P. H. Bauer, "Realization using the Roesser model for implementations in distributed grid sensor networks," *Multidimensional Systems and Signal Processing*, vol. 22, pp. 131-146, Mar. 2011.
- [23] T. Zhou, "Boundedness of multidimensional filters over a prescribed frequency domain," *IEEE Trans. Signal Processing*, vol. 56, pp. 5487-5499, Nov. 2008.
- [24] T. Zhou, "Stability and stability margin for a two-dimensional system," *IEEE Trans. Signal Processing*, vol. 54, pp. 3483-3488, Sept. 2006.
- [25] T. Hinamoto, Y. Sugie, A. Doi and M. Muneyasu, "Synthesis of 3-D separable-denominator state-space digital filters with minimum L_2 -sensitivity," *Multidimensional Systems and Signal Processing*, vol. 15, pp. 147-167, Apr. 2004.

- [26] R. J. Ober, X. Lai, Z. Lin and E. S. Ward, "State-space realization of a three-dimensional image set with application to noise reduction of fluorescence microscopy images of cells," *Multidimensional Systems and Signal Processing*, vol. 16, pp. 7-48, Jan. 2005.
- [27] C. T. Chen, *Introduction to Linear System Theory*, Holt, Rinehart and Winston, Inc., 1970.
- [28] L. L. Scharf, *Statistical Signal Processing*, Reading, MA: Addison-Wesley, 1991.
- [29] H. Togawa, *Handbook of Numerical Methods*, Saiensu-sha, Tokyo, 1992.
- [30] R. Fletcher, *Practical Methods of Optimization*, 2nd ed., Wiley, New York, 1987.
- [31] J. E. Dennis and J. J. More, "Quasi-Newton methods, motivation and theory," *SIAM Review*, vol. 19, pp. 46-89, 1977.
- [32] R. E. Kalman, "Mathematical description of linear dynamical systems," *SIAM J. Control (ser. A)*, vol. 1, pp. 152-192, 1963.
- [33] T. Kailath, *Linear System*, Englewood Cliffs, N.J.: Prentice-Hall, 1980.