

**Roundoff Noise Minimization
for 2-D Separable-Denominator Digital Filters
Using Jointly Optimal High-Order Error Feedback and Realization**

Takao Hinamoto

Hiroshima University
Higashi-Hiroshima, Japan

Akimitsu Doi

Hiroshima Institute of Technology
Hiroshima, Japan

Wu-Sheng Lu

University of Victoria
Victoria, Canada

May 31, 2017

Outline

- Early Work and Objectives
- Model
- Noise Gain and Its Minimization
- Numerical Example

1. Early Work and Objectives

Error Feedback (EF)

- Higgins and Munson, 1984.
- Vaidyanathan, 1985.
- Laakso and Hartimo, 1992.

EF Combined with State-Space Realization

- Hinamoto, Doi & Lu, 2003, 2005, 2013.

The Problem of Studies

Jointly optimizing *high-order* EF and realization for minimizing roundoff noise subject to l_2 -scaling constraints for 2-D separable denominator digital filters.

- ◇ Closed-form solutions for the controllability and observability Grammians;
- ◇ EF order is now made to be independent from the dimension of the state.

2. Model

$$\mathbf{x}_1(i, j) = \mathbf{A}\mathbf{x}_0(i, j) + \mathbf{b}u(i, j)$$

$$y(i, j) = \mathbf{c}\mathbf{x}_0(i, j) + du(i, j)$$

where the filter's separable denominator means that

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{0} & \mathbf{A}_4 \end{bmatrix}$$

- Taking error feedback and finite-word-length (FWL) effect into account, the model becomes

$$\begin{aligned} \tilde{\mathbf{x}}_1(i, j) &= \mathbf{A}Q[\tilde{\mathbf{x}}_0(i, j)] + \mathbf{b}u(i, j) \\ &\quad + \sum_{k=0}^{M-1} [\mathbf{D}_{1k} \oplus \mathbf{0}] \mathbf{e}_{-k}(i, j) + \sum_{l=0}^{N-1} [\mathbf{0} \oplus \mathbf{D}_{4l}] \mathbf{e}_{-l}(i, j) \\ \tilde{y}(i, j) &= \mathbf{c}Q[\tilde{\mathbf{x}}_0(i, j)] + du(i, j) + \mathbf{h}\mathbf{e}_0(i, j) \end{aligned}$$

where $\mathbf{h} \in R^{1 \times (m+n)}$ is an error-feedforward vector, $\mathbf{D}_{1k} \in R^{m \times m}$ and $\mathbf{D}_{4l} \in R^{n \times n}$ are high-order EF matrices, and

$$\mathbf{e}_k(i, j) = \tilde{\mathbf{x}}_k(i, j) - Q[\tilde{\mathbf{x}}_k(i, j)]$$

- From above equations, we have

$$\begin{aligned}\Delta \mathbf{x}_1(i, j) &= \mathbf{A} \Delta \mathbf{x}_0(i, j) + \mathbf{A} \mathbf{e}_0(i, j) - \sum_{k=0}^{M-1} [\mathbf{D}_{1k} \oplus \mathbf{0}] \mathbf{e}_{-k}(i, j) - \sum_{l=0}^{N-1} [\mathbf{0} \oplus \mathbf{D}_{4l}] \mathbf{e}_{-l}(i, j) \\ \Delta y(i, j) &= \mathbf{c} \Delta \mathbf{x}_0(i, j) + (\mathbf{c} - \mathbf{h}) \mathbf{e}_0(i, j)\end{aligned}$$

where $\Delta \mathbf{x}_k(i, j) = \mathbf{x}_k(i, j) - \tilde{\mathbf{x}}_k(i, j)$ and $\Delta y(i, j) = y(i, j) - \tilde{y}(i, j)$.

- Based on this, we obtain the frequency-domain model

$$\Delta Y(z_1, z_2) = \begin{bmatrix} \mathbf{H}_e^h(z_1) & \mathbf{H}_e^v(z_1, z_2) \end{bmatrix} \mathbf{E}_0(z_1, z_2)$$

where

$$\begin{aligned}\mathbf{H}_e^h(z_1) &= \mathbf{c}_1 \sum_{i=0}^{\infty} \left(\mathbf{A}_1^i - \sum_{k=0}^{M-1} \mathbf{A}_1^{i-k-1} \mathbf{D}_{1k} \right) z_1^{-i} - \mathbf{h}_1 \\ \mathbf{H}_e^v(z_1, z_2) &= \left[\mathbf{c}_2 + \mathbf{c}_1 (z_1 \mathbf{I}_m - \mathbf{A}_1)^{-1} \mathbf{A}_2 \right] \sum_{i=0}^{\infty} \left(\mathbf{A}_4^i - \sum_{k=0}^{N-1} \mathbf{A}_4^{i-k-1} \mathbf{D}_{4k} \right) z_2^{-i} - \mathbf{h}_2\end{aligned}$$

3. Noise Gain and Its Minimization

3.A Noise Gain

- Based on above analysis, the ***normalized noise gain*** which is a scaled noise variance at the filter output is found to be

$$J_0(\mathbf{h}, \mathbf{D}_1, \mathbf{D}_4) = J_0^h(\mathbf{h}_1, \mathbf{D}_1) + J_0^v(\mathbf{h}_2, \mathbf{D}_4)$$

where

$$J_0^h(\mathbf{h}_1, \mathbf{D}_1) = \text{trace} \left[\frac{1}{2\pi j} \oint_{|z_1|=1} \mathbf{H}_e^h(z_1)^* \mathbf{H}_e^h(z_1) \frac{dz_1}{z_1} \right]$$

$$J_0^v(\mathbf{h}_2, \mathbf{D}_4) = \text{trace} \left[\frac{1}{(2\pi j)^2} \oint_{|z_1|=1} \oint_{|z_2|=1} \mathbf{H}_e^v(z_1, z_2)^* \mathbf{H}_e^v(z_1, z_2) \frac{dz_1 dz_2}{z_1 z_2} \right]$$

with

$$\mathbf{D}_1 = \begin{bmatrix} \mathbf{D}_{10} & \mathbf{D}_{11} & \cdots & \mathbf{D}_{1,M-1} \end{bmatrix}, \quad \mathbf{D}_4 = \begin{bmatrix} \mathbf{D}_{40} & \mathbf{D}_{41} & \cdots & \mathbf{D}_{4,N-1} \end{bmatrix}$$

- Utilizing the expression of $\mathbf{H}_e^h(z_1)$ and $\mathbf{H}_e^v(z_1, z_2)$ in terms of the state-space model parameters, the noise gain can be expressed explicitly as follows:

$$\begin{aligned}
J_0^h(\mathbf{h}_1, \mathbf{D}_1) = & \text{trace} \left[\mathbf{W}^h - \sum_{k=0}^{M-1} \left\{ \mathbf{D}_{1k}^T \mathbf{W}^h \mathbf{A}_1^{k+1} + (\mathbf{A}_1^T)^{k+1} \mathbf{W}^h \mathbf{D}_{1k} \right\} \right. \\
& + \sum_{k=0}^{M-1} \sum_{l=0}^{M-1} \mathbf{D}_{1k}^T \left\{ \mathbf{W}^h \mathbf{A}_1^{k-l} + (\mathbf{A}_1^T)^{l-k} \mathbf{W}^h \right\} \mathbf{D}_{1l} \\
& \left. - \sum_{k=0}^{M-1} \mathbf{D}_{1k}^T \mathbf{W}^h \mathbf{D}_{1k} - 2\mathbf{h}_1^T \mathbf{c}_1 + \mathbf{h}_1^T \mathbf{h}_1 \right]
\end{aligned}$$

and

$$\begin{aligned}
J_0^v(\mathbf{h}_2, \mathbf{D}_4) = & \text{trace} \left[\mathbf{W}^v - \sum_{k=0}^{N-1} \left\{ \mathbf{D}_{4k}^T \mathbf{W}^v \mathbf{A}_4^{k+1} + (\mathbf{A}_4^T)^{k+1} \mathbf{W}^v \mathbf{D}_{4k} \right\} \right. \\
& + \sum_{k=0}^{N-1} \sum_{l=0}^{N-1} \mathbf{D}_{4k}^T \left\{ \mathbf{W}^v \mathbf{A}_4^{k-l} + (\mathbf{A}_4^T)^{l-k} \mathbf{W}^v \right\} \mathbf{D}_{4l} \\
& \left. - \sum_{k=0}^{N-1} \mathbf{D}_{4k}^T \mathbf{W}^v \mathbf{D}_{4k} - 2\mathbf{h}_2^T \mathbf{c}_2 + \mathbf{h}_2^T \mathbf{h}_2 \right]
\end{aligned}$$

where \mathbf{W}^h and \mathbf{W}^v are the horizontal and vertical observability Grammians of the state-space filter which can be found by solving the Lyapunov equations:

$$\mathbf{W}^h = \mathbf{A}_1^T \mathbf{W}^h \mathbf{A}_1 + \mathbf{c}_1^T \mathbf{c}_1$$

and

$$\mathbf{W}^v = \mathbf{A}_4^T \mathbf{W}^v \mathbf{A}_4 + \mathbf{A}_2^T \mathbf{W}^h \mathbf{A}_2 + \mathbf{c}_2^T \mathbf{c}_2$$

- The noise gain can be considerably simplified if the element EF matrices in \mathbf{D}_1 and \mathbf{D}_4 are diagonal:

$$J^h(\mathbf{h}_1, \mathbf{D}_1) = \text{tr} \left[\mathbf{W}^h - \mathbf{c}_1^T \mathbf{c}_1 - 2 \sum_{k=0}^{M-1} \mathbf{W}^h \mathbf{A}_1^{k+1} \mathbf{D}_{1k} + \sum_{k=0}^{M-1} \sum_{l=0}^{M-1} \mathbf{W}^h \mathbf{A}_1^{|k-l|} \mathbf{D}_{1k} \mathbf{D}_{1l} \right] + (\mathbf{c}_1 - \mathbf{h}_1)(\mathbf{c}_1 - \mathbf{h}_1)^T$$

and

$$J^v(\mathbf{h}_2, \mathbf{D}_4) = \text{tr} \left[\mathbf{W}^v - \mathbf{c}_2^T \mathbf{c}_2 - 2 \sum_{k=0}^{N-1} \mathbf{W}^v \mathbf{A}_4^{k+1} \mathbf{D}_{4k} + \sum_{k=0}^{N-1} \sum_{l=0}^{N-1} \mathbf{W}^v \mathbf{A}_4^{|k-l|} \mathbf{D}_{4k} \mathbf{D}_{4l} \right] + (\mathbf{c}_2 - \mathbf{h}_2)(\mathbf{c}_2 - \mathbf{h}_2)^T$$

3.B Minimization of Noise Gain

We now seek to find a state-space realization of a given state-space filter, $(\bar{A}, \bar{b}, \bar{c}, d)$ where

$$\bar{A} = T^{-1}AT, \quad \bar{b} = T^{-1}b, \quad \bar{c} = cT$$

with $T = T_1 \oplus T_4$, such that its noise gain is minimized subject to the l_2 -scaling constraints

$$(\bar{K})_{ii} = (T^{-1}KT)_{ii} = 1 \quad \text{for } 1 \leq i \leq m+n$$

where $K = K^h \oplus K^v$ is the controllability Grammians which can be found by solving the Laypunov equations

$$K^v = A_4 K^v A_4^T + b_2 b_2^T$$

and

$$K^h = A_1 K^h A_1^T + A_2 K^v A_2^T + b_1 b_1^T$$

- Evidently, the variables involved in this problem are coordinate transformation matrix T as well as diagonal EF matrices in D_1 and D_4 which are jointly optimized by minimizing the objective

$$\bar{J}(T, D_1, D_4) = \bar{J}^h(T_1, D_1) + \bar{J}^v(T_4, D_4)$$

where

$$\begin{aligned} \bar{J}^h(\mathbf{T}_1, \mathbf{D}_1) = & \text{tr} \left[\mathbf{T}_1^T (\mathbf{W}^h - \mathbf{c}_1^T \mathbf{c}_1) \mathbf{T}_1 - 2 \sum_{k=0}^{M-1} \mathbf{T}_1^T \mathbf{W}^h \mathbf{A}_1^{k+1} \mathbf{T}_1 \mathbf{D}_{1k} \right. \\ & \left. + \sum_{k=0}^{M-1} \sum_{l=0}^{M-1} \mathbf{T}_1^T \mathbf{W}^h \mathbf{A}_1^{|k-l|} \mathbf{T}_1 \mathbf{D}_{1k} \mathbf{D}_{1l} \right] \end{aligned}$$

and

$$\begin{aligned} \bar{J}^v(\mathbf{T}_4, \mathbf{D}_4) = & \text{tr} \left[\mathbf{T}_4^T (\mathbf{W}^v - \mathbf{c}_2^T \mathbf{c}_2) \mathbf{T}_4 - 2 \sum_{k=0}^{N-1} \mathbf{T}_4^T \mathbf{W}^v \mathbf{A}_4^{k+1} \mathbf{T}_4 \mathbf{D}_{4k} \right. \\ & \left. + \sum_{k=0}^{N-1} \sum_{l=0}^{N-1} \mathbf{T}_4^T \mathbf{W}^v \mathbf{A}_4^{|k-l|} \mathbf{T}_4 \mathbf{D}_{4k} \mathbf{D}_{4l} \right] \end{aligned}$$

- First, the l_2 constraints $(\mathbf{T}^{-1} \mathbf{K} \mathbf{T})_{ii} = 1$ are eliminated in two simple steps:
 (1) Variable change $\hat{\mathbf{T}} = \hat{\mathbf{T}}_1 \oplus \hat{\mathbf{T}}_4 = (\mathbf{T}_1 \oplus \mathbf{T}_4)^T (\mathbf{K}^h \oplus \mathbf{K}^v)^{-1/2}$ simplifies the constraints to $(\hat{\mathbf{T}}^{-T} \hat{\mathbf{T}}^{-1})_{ii} = 1$;
 (2) The above constraints are automatically satisfied by assuming

$$\hat{\mathbf{T}}_1^{-1} = \begin{bmatrix} \frac{\mathbf{t}_{11}}{\|\mathbf{t}_{11}\|} & \frac{\mathbf{t}_{12}}{\|\mathbf{t}_{12}\|} & \dots & \frac{\mathbf{t}_{1m}}{\|\mathbf{t}_{1m}\|} \end{bmatrix} \quad \text{and} \quad \hat{\mathbf{T}}_4^{-1} = \begin{bmatrix} \frac{\mathbf{t}_{41}}{\|\mathbf{t}_{41}\|} & \frac{\mathbf{t}_{42}}{\|\mathbf{t}_{42}\|} & \dots & \frac{\mathbf{t}_{4n}}{\|\mathbf{t}_{4n}\|} \end{bmatrix}$$

- Now the design variables in the gain minimization problem are $t_{11}, \dots, t_{1m}, t_{41}, \dots, t_{4n}, D_{10}, \dots, D_{1,M-1}, D_{40}, \dots, D_{4,N-1}$. These variables are arranged to form a column vector \mathbf{x} , so the noise gain is a function of \mathbf{x} , hence is denoted by $J(\mathbf{x})$.
- A quasi-Newton algorithm with BFGS updates was applied to minimize $J(\mathbf{x})$ iteratively.

◊ It starts with a trivial initial point \mathbf{x}_0 which is obtained by assigning T_1, T_4 , and all D_{1k} and D_{4l} to be identity matrices.

◊ It updates iterate \mathbf{x}_k to \mathbf{x}_{k+1} using $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ where

$$\mathbf{d}_k = -\mathbf{S}_k \nabla J(\mathbf{x}_k), \quad \alpha_k = \arg \left[\min_{\alpha} J(\mathbf{x}_k + \alpha \mathbf{d}_k) \right]$$

$$\mathbf{S}_{k+1} = \mathbf{S}_k + \left(1 + \frac{\gamma_k^T \mathbf{S}_k \gamma_k}{\gamma_k^T \delta_k} \right) \frac{\delta_k \delta_k^T}{\gamma_k^T \delta_k} - \frac{\delta_k \gamma_k^T \mathbf{S}_k + \mathbf{S}_k \gamma_k \delta_k^T}{\gamma_k^T \delta_k}$$

$$\mathbf{S}_0 = \mathbf{I}, \quad \delta_k = \mathbf{x}_{k+1} - \mathbf{x}_k, \quad \gamma_k = \nabla J(\mathbf{x}_{k+1}) - \nabla J(\mathbf{x}_k)$$

◊ Closed-form formula for $\nabla J(\mathbf{x}_k)$ for fast and accurate computation has been derived .

4. Numerical Example

We consider a 2-D separable-denominator digital filter $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)_{3+3}$ with

$$\mathbf{A}_1^T = \begin{bmatrix} 0 & 0 & 0.599655 \\ 1 & 0 & -1.836929 \\ 0 & 1 & 2.173645 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 0.064564 & 0.033034 & 0.012881 \\ 0.091213 & 0.110512 & 0.102759 \\ 0.097256 & 0.151864 & 0.172460 \end{bmatrix}$$

$$\mathbf{A}_4 = \begin{bmatrix} 0 & 0 & 0.564961 \\ 1 & 0 & -1.887939 \\ 0 & 1 & 2.280029 \end{bmatrix}, \quad \mathbf{b}_1 = \begin{bmatrix} 0.047053 \\ 0.062274 \\ 0.060436 \end{bmatrix}, \quad \mathbf{b}_2 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{c}_1^T = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{c}_2 = [0.016556 \quad 0.012550 \quad 0.008243], \quad d = 0.019421.$$

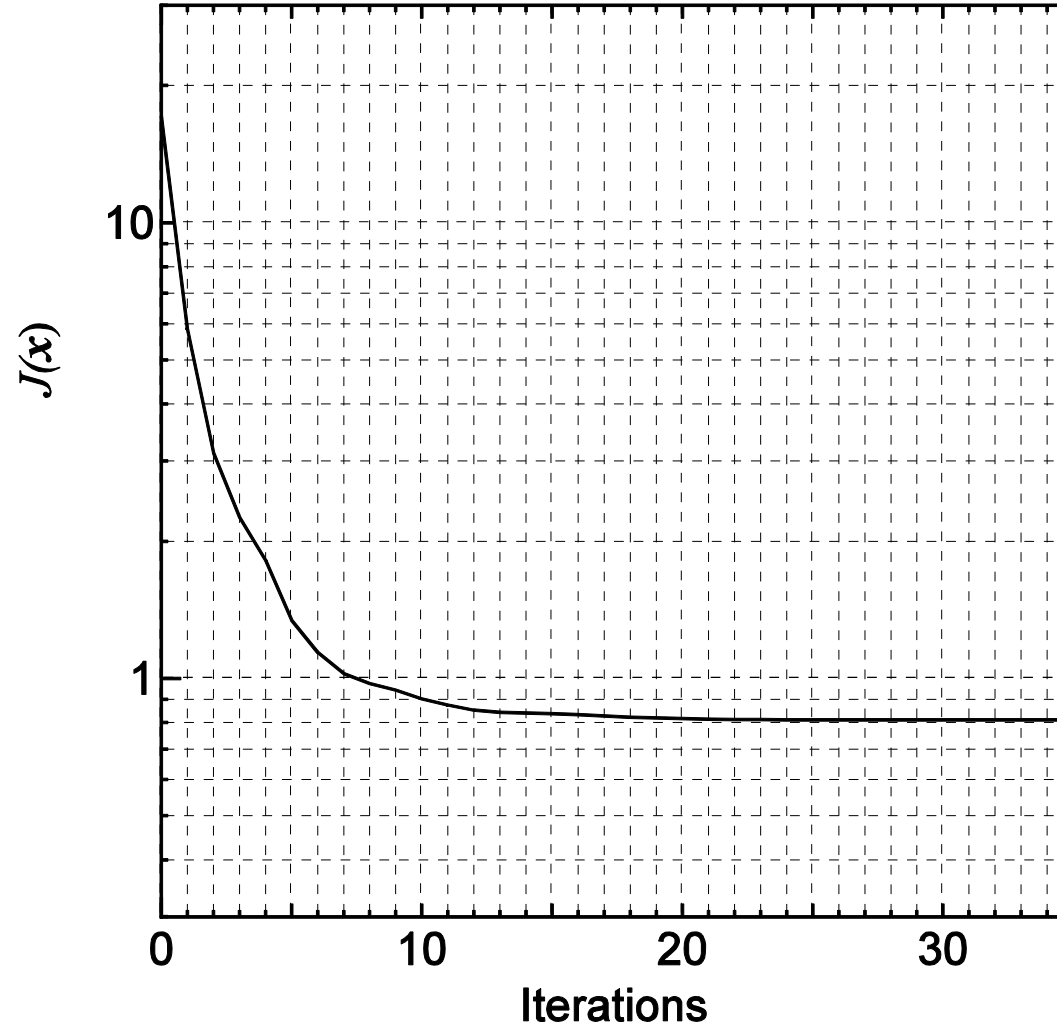
Case 1: Minimum noise gain without error feedforward and feedback :

$$\hat{J}_{\min}(\hat{T}, \mathbf{0}, \mathbf{0}) = 7.936533$$

Case 2: Minimum noise gain with error feedforward and feedback:

$(\underline{M}, \underline{N})$	<u>Infinite Precision</u>	<u>3-Bit Quantization</u>	<u>Integer Quantization</u>
(1, 1)	1.3004	1.3518	3.0832
(2, 1)	1.0217	1.0672	2.2677
(1, 2)	1.0899	1.1349	3.0247
(2, 2)	0.8112	0.8502	2.2091

- The figure below depicts the profile of the noise gain $J(\mathbf{x}_k)$ in first 35 iterations for the case of $(M, N) = (2, 2)$.



Thank you.

Q & A