



HAL
open science

On the Transfer Function Error of State-Space Filters in Fixed-Point Context

Thibault Hilaire

► **To cite this version:**

Thibault Hilaire. On the Transfer Function Error of State-Space Filters in Fixed-Point Context. IEEE Transactions on Circuits and Systems II: Express Briefs, 2009, 56 (12), pp.936-940. 10.1109/TC-SII.2009.2034193 . hal-01146515

HAL Id: hal-01146515

<https://hal.science/hal-01146515v1>

Submitted on 14 Dec 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

On the Transfer Function Error of State-Space Filters in Fixed-Point Context

Thibault Hilaire

Abstract—This paper presents a new measure used for the implementation of filters/controllers in state-space form. It investigates the transfer function deviation generated by the coefficient quantization. The classical L_2 -sensitivity measure is extended with precise consideration on their fixed-point representation, in order to make a more valid measure. By solving the related optimal realization problem, fixed-point accurate realizations in state-space form can be found.

Index Terms—Digital filter implementation, coefficient sensitivity, fixed-point implementation.

I. INTRODUCTION

The majority of control or signal processing systems is implemented in digital general purpose processors, DSPs¹, FPGAs², etc. Since these devices cannot compute with infinite precision and approximate real-number parameters with a finite binary representation, the numerical implementation of controllers (filters) leads to deterioration in characteristics and performance. This has two separate origins, corresponding to the quantization of the embedded coefficients and the roundoff errors occurring during the computations. They can be formalized as parametric errors and numerical noises, respectively. The focus of this paper are parametric errors, but one can refer to [1]–[4] for roundoff noises, where measures with fixed-point consideration already exist.

It is also well known that these Finite Word Length (FWL) effects depend on the structure of the realization. In state-space form, the realization depends on the choice of the basis of the state-vector. This motivates us to investigate the coefficient sensitivity minimization problem. It has been well studied with the L_2 -measure [1], [5]. However this measure only considers how sensitive to the coefficients the transfer function is, and does not investigate the coefficients' quantization, which depends on the fixed-point representation used. In [5], the transfer function error is exhibited for the first time, however, only for quantized coefficients with the same binary-point position.

This paper investigates the transfer function deviation generated by the coefficient quantization with precise consideration on their fixed-point representation. The classical L_2 -sensitivity analysis is shown in section II, whereas the new approach,

based on fixed-point consideration, is presented in section III. A comparison with the L_2 -sensitivity and some scaling considerations are provided. Finally, the optimal realization problem is solved in section IV and a numerical example is exhibited before conclusion.

II. L_2 -SENSITIVITY ANALYSIS

Let (A, b, c, d) be a stable, controllable and observable linear discrete time Single Input Single Output (SISO) state-space system, *i.e.*

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

where $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{b} \in \mathbb{R}^{n \times 1}$, $\mathbf{c} \in \mathbb{R}^{1 \times n}$ and $d \in \mathbb{R}$. $u(k)$ is the scalar input, $y(k)$ is the scalar output and $\mathbf{x}(k) \in \mathbb{R}^{n \times 1}$ is the state vector.

Its input-output relationship is given by the scalar transfer function $h : \mathbb{C} \rightarrow \mathbb{C}$ defined by:

$$h : z \mapsto \mathbf{c}(z\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{b} + d. \quad (2)$$

The quantization of the coefficients introduces some uncertainties to \mathbf{A} , \mathbf{b} , \mathbf{c} and d leading to $\mathbf{A} + \Delta\mathbf{A}$, $\mathbf{b} + \Delta\mathbf{b}$, $\mathbf{c} + \Delta\mathbf{c}$ and $d + \Delta d$ respectively. It is common to consider the sensitivity of the transfer function with respect to the coefficients, based on the following definitions.

Definition 1 (Transfer function sensitivity) Consider $\mathbf{X} \in \mathbb{R}^{m \times n}$ and $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{C}$ differentiable with respect to all the entries of \mathbf{X} .

The sensitivity of f with respect to \mathbf{X} is defined by the matrix $\mathbf{S}_{\mathbf{X}} \in \mathbb{R}^{m \times n}$:

$$\frac{\partial f}{\partial \mathbf{X}} \triangleq \mathbf{S}_{\mathbf{X}} \quad \text{with} \quad (\mathbf{S}_{\mathbf{X}})_{i,j} \triangleq \frac{\partial f}{\partial \mathbf{X}_{i,j}}. \quad (3)$$

Applied to a scalar transfer function h where $h(z)$ depends on a given matrix \mathbf{X} , $\frac{\partial h}{\partial \mathbf{X}}$ is a transfer function of a Multiple Inputs Multiple Outputs (MIMO) system.

Definition 2 (L_2 -Norm) Let $\mathbf{H} : \mathbb{C} \rightarrow \mathbb{C}^{k \times l}$ be a function of the scalar complex variable z . $\|\mathbf{H}\|_2$ is the L_2 -norm of \mathbf{H} , defined by:

$$\|\mathbf{H}\|_2 \triangleq \sqrt{\frac{1}{2\pi} \int_0^{2\pi} \|\mathbf{H}(e^{j\omega})\|_F^2 d\omega} \quad (4)$$

where $\|\cdot\|_F$ is the Froebenius norm.

This work has been partially funded by the NFN SISE project (National Research Network "Signal and Information Processing in Science and Engineering") and the Institute of Communications and Radio-Frequency Engineering of Vienna University of Technology, France. T. Hilaire is now with the Laboratory of Computer Science (LIP6) of the University Pierre & Marie Curie of Paris, France.

¹Digital Signal Processors

²Field Programmable Gate-Array

Proposition 1 If H is the MIMO state-space system (K, L, M, N) , then its L_2 -norm can be computed by

$$\|H\|_2^2 = \text{tr}(NN^\top + MW_cM^\top) \quad (5)$$

$$= \text{tr}(N^\top N + L^\top W_oL) \quad (6)$$

where W_c and W_o are the controllability and observability Gramians, respectively. They are solutions of the Lyapunov equations

$$W_c = KW_cK^\top + LL^\top, \quad W_o = K^\top W_oK + M^\top M. \quad (7)$$

Proof: See [1]. ■

Gevers and Li [1] have proposed the L_2 -sensitivity measure to evaluate the coefficient roundoff errors. It is defined by

$$M_{L_2} \triangleq \left\| \frac{\partial h}{\partial A} \right\|_2^2 + \left\| \frac{\partial h}{\partial b} \right\|_2^2 + \left\| \frac{\partial h}{\partial c} \right\|_2^2 + \left\| \frac{\partial h}{\partial d} \right\|_2^2 \quad (8)$$

and can be computed by $\frac{\partial h}{\partial A}(z) = G^\top(z)F^\top(z)$, $\frac{\partial h}{\partial b}(z) = G^\top(z)$, $\frac{\partial h}{\partial c}(z) = F(z)$ and $\frac{\partial h}{\partial d}(z) = 1$, with

$$F(z) \triangleq (zI_n - A)^{-1}b, \quad G(z) \triangleq c(zI_n - A)^{-1}. \quad (9)$$

This measure is an extension of the more tractable but less natural L_1/L_2 sensitivity measure proposed by V. Tavşanoğlu and L. Thiele [6] ($\left\| \frac{\partial h}{\partial A} \right\|_1^2$ instead of $\left\| \frac{\partial h}{\partial A} \right\|_2^2$ in (8)).

Remark 1 To simplify the expressions, it is also possible to regroup all the coefficients in one unique matrix Z :

$$Z \triangleq \begin{pmatrix} A & b \\ c & d \end{pmatrix}. \quad (10)$$

Then, with L_2 -norm property, $M_{L_2} = \left\| \frac{\partial h}{\partial Z} \right\|_2^2$. From (9) and the associated state-spaces, the sensitivity transfer function $\frac{\partial h}{\partial Z}$ can be described by the MIMO state-space system $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D})$ with

$$\begin{aligned} \tilde{A} &\triangleq \begin{pmatrix} A & bc \\ \mathbf{0} & A \end{pmatrix}, \tilde{B} \triangleq \begin{pmatrix} \mathbf{0} & b \\ I_n & \mathbf{0} \end{pmatrix}, \\ \tilde{C} &\triangleq \begin{pmatrix} I_n & \mathbf{0} \\ \mathbf{0} & c \end{pmatrix}, \tilde{D} \triangleq \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix}. \end{aligned} \quad (11)$$

The proposition 1 is used to compute M_{L_2} . See [1] and [7] for more details.

Applying a coordinate transformation, defined by $\bar{x}(k) \triangleq \mathcal{U}^{-1}x(k)$ to the state-space system (A, b, c, d) , leads to a new equivalent realization $(\mathcal{U}^{-1}A\mathcal{U}, \mathcal{U}^{-1}b, c\mathcal{U}, d)$.

Since these two realizations are equivalent in infinite precision but are no more equivalent in finite precision (fixed point arithmetic, floating-point arithmetic, etc.), the L_2 -sensitivity then depends on \mathcal{U} , and is denoted $M_{L_2}(\mathcal{U})$. In this case, it is natural to define the following problem:

Problem 1 (optimal L_2 -sensitivity problem) Considering a state-space realization (A, b, c, d) , the optimal L_2 -sensitivity problem consists of finding the coordinate transformation \mathcal{U}_{opt} that minimizes M_{L_2} :

$$\mathcal{U}_{opt} = \arg \min_{\mathcal{U} \text{ invertible}} M_{L_2}(\mathcal{U}). \quad (12)$$

In [1], it is shown that the problem has one unique solution. Hence, for example, a gradient method can be used to solve it.

III. TRANSFER FUNCTION ERROR

A. Fixed-point implementation

In this paper, the notation (β, γ) is used for the fixed-point representation of a variable or coefficient (2's complement scheme), according to Figure 1. β is the total wordlength of the representation in bits, whereas γ is the wordlength of the fractional part (it determines the position of the binary-point). They are fixed for each variable (input, states, output) and each coefficient, and implicit (unlike the floating-point representation). β and γ will be suffixed by the variable/coefficient they refer to. These parameters could be scalars, vectors or matrices, according to the variables they refer to.

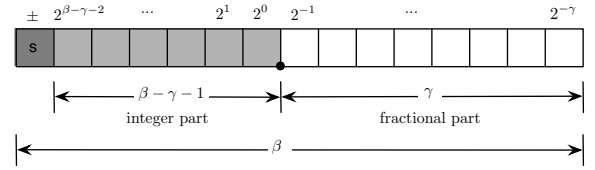


Fig. 1. Fixed-point representation

Let us suppose that the wordlength of the coefficients β_Z is given³. Then, the coefficients Z_{ij} are represented in fixed-point by $(\beta_{Z_{ij}}, \gamma_{Z_{ij}})$ with:

$$\gamma_{Z_{ij}} = \beta_{Z_{ij}} - 2 - \lfloor \log_2 |Z_{ij}| \rfloor. \quad (13)$$

where the $\lfloor a \rfloor$ operation rounds a to the nearest integer less or equal to a (for positive numbers $\lfloor a \rfloor$ is the integer part).

Remark 2 The binary point position is not defined for null coefficients, however this is no problem because these coefficients will not be represented in the final algorithm (the null multiplications are removed).

So, in order to consider coefficients that will be quantized without error, we introduced a *weighting* matrix δ_Z such that

$$(\delta_Z)_{ij} \triangleq \begin{cases} 0 & \text{if } Z_{ij} \text{ is exactly implemented} \\ 1 & \text{otherwise.} \end{cases} \quad (14)$$

The exactly implemented coefficients are 0, ± 1 and also positive and negative coefficients of power of 2.

Remark 3 In some specific computational cases the fixed-point representation chosen for the coefficients is not always the best one as defined in (13). For example, in the *Roundoff Before Multiplication* scheme, some extra quantizations are added to the coefficients, in order to avoid shift operations after multiplications [2]. Only the classical case (corresponding to the *Roundoff After Multiplication*) is considered here, as defined by equation (13).

Remark 4 It is also possible to choose the same fixed-point representation for all the coefficients (determined by the

³In FPGA or ASIC, it is of interest to consider the wordlength as optimization variables, in order to find hardware realizations that minimize hardware criteria like power consumption or surface, under certain numerical accuracy constraints, like L_2 -sensitivity ones [8]. This is not considered in this paper.

coefficient with the highest magnitude). But in that case, the lowest coefficients (in magnitude) do not have an appropriate representation. They are coded with less meaningful bits and have a higher relative error. When the ratio between the greatest and lowest magnitude is too high, then underflows occurs for the lowest coefficients that cannot be represented. For example, this is common for the Direct Form realizations with high (or low) L_2 -gain.

During the quantization process, the coefficients are changed from \mathbf{Z} into $\mathbf{Z}^\dagger \triangleq \mathbf{Z} + \Delta\mathbf{Z}$. For a best-roundoff quantization, the $\{\Delta\mathbf{Z}_{i,j}\}$ are independent centered random variables uniformly distributed [9] within the ranges $-2^{-\gamma\mathbf{z}_{ij}-1} \leq \mathbf{Z}_{i,j} < 2^{-\gamma\mathbf{z}_{ij}-1}$, so their second-order moments are given by

$$\sigma_{\mathbf{Z}_{i,j}}^2 \triangleq E\{(\Delta\mathbf{Z}_{i,j})^2\} \quad (15)$$

$$= \frac{2^{-2\gamma\mathbf{z}_{ij}}}{12} \delta_{\mathbf{Z}_{i,j}}, \quad (16)$$

where $E\{\cdot\}$ is the mean operator.

B. Transfer function error

Due to the quantization of the coefficients, the transfer function is changed from h to $h^\dagger \triangleq h + \Delta h$. This degradation can be evaluated in a statistical way with the following definition.

Definition 3 (Transfer function error) A measure of the transfer function error can be statistically defined by [5]

$$\sigma_{\Delta h}^2 \triangleq \frac{1}{2\pi} \int_0^{2\pi} E\{|\Delta h(e^{j\omega})|^2\} d\omega. \quad (17)$$

The transfer function error is a tractable measure that can be evaluated with the following proposition.

Proposition 2 The transfer function error is given by:

$$\sigma_{\Delta h}^2 = \left\| \frac{\partial h}{\partial \mathbf{Z}} \times \Xi_{\mathbf{Z}} \right\|_2^2 \quad (18)$$

where \times is the Schur product, $\Xi_{\mathbf{Z}} \in \mathbb{R}^{(n+1) \times (n+1)}$ defined by:

$$(\Xi_{\mathbf{Z}})_{i,j} \triangleq \begin{cases} \frac{2^{-\beta\mathbf{z}_{ij}+1}}{\sqrt{3}} \lfloor \mathbf{Z}_{ij} \rfloor_2 (\delta_{\mathbf{Z}})_{ij} & \text{if } \mathbf{Z}_{ij} \neq 0 \\ 0 & \text{if } \mathbf{Z}_{ij} = 0 \end{cases} \quad (19)$$

and $\lfloor x \rfloor_2$ is nearest power of 2 lower than $|x|$:

$$\lfloor x \rfloor_2 \triangleq 2^{\lfloor \log_2 |x| \rfloor}. \quad (20)$$

Proof: A first order approximation gives

$$\Delta h(z) = \sum_{i,j} \frac{\partial h}{\partial \mathbf{Z}_{i,j}}(z) \Delta \mathbf{Z}_{i,j}, \quad \forall z \in \mathbb{C}. \quad (21)$$

Hence

$$E\{|\Delta h(e^{j\omega})|^2\} = \sum_{i,j} \left| \frac{\partial h(e^{j\omega})}{\partial \mathbf{Z}_{i,j}} \right|^2 \sigma_{\mathbf{Z}_{i,j}}^2 \quad (22)$$

because the random variables $\Delta\mathbf{Z}_{i,j}$ are independent.

Considering (13) and (16) for non-null coefficients, we get

$$\sigma_{\mathbf{Z}_{i,j}}^2 = \frac{4}{3} 2^{-2\beta\mathbf{z}_{ij}} \lfloor \mathbf{Z}_{i,j} \rfloor_2^2 (\delta_{\mathbf{Z}})_{ij} \quad (23)$$

and

$$\sigma_{\Delta h}^2 = \sum_{i,j} \left\| (\Xi)_{i,j} \frac{\partial h}{\partial \mathbf{Z}_{i,j}} \right\|_2^2 \quad (24)$$

Then, with $(\frac{\partial h}{\partial \mathbf{Z}})_{i,j} = \frac{\partial h}{\partial \mathbf{Z}_{i,j}}$ and (4), eq. (18) holds. ■

Remark 5 In the classical case where the wordlength of the coefficients are all the same (equal to β), we can define a normalized transfer error $\bar{\sigma}_{\Delta h}^2$ defined by:

$$\bar{\sigma}_{\Delta h}^2 \triangleq \frac{3\sigma_{\Delta h}^2}{2^{-2\beta+2}}. \quad (25)$$

This measure is now independent of the wordlength, and can be used for some comparisons.

C. Comparison with the classical M_{L_2} measure

It is of interest to remark the relationship with the classical M_{L_2} measure. In [5] where the transfer function error appears for the first time, the coefficients are supposed to have the same fixed-point representation, so their second-order moment ($\sigma_{\mathbf{Z}_{i,j}}^2$) are all equal and denoted σ_0^2 . The M_{L_2} satisfies then

$$M_{L_2} = \frac{\sigma_{\Delta h}^2}{\sigma_0^2}. \quad (26)$$

Here, the transfer function error $\sigma_{\Delta h}^2$ can be seen as an extension of the M_{L_2} measure with fixed-point considerations. The sensitivity is weighted according to the variance of the quantization noise of each coefficient.

The M_{L_2} measure considers each coefficient with the same weight, even if the quantization of one particular coefficient induces a small or big modification of this coefficient. To avoid this problem, a *normalization* has been created. Since a coordinate transformation $\mathbf{U} = \alpha \mathbf{I}_n$ does not change \mathbf{A} , but only multiplies the coefficients of \mathbf{b} by $\frac{1}{\alpha}$, and those of \mathbf{c} by α , it was of interest to set a condition on \mathbf{b} or \mathbf{c} for normalization.

The L_2 -dynamic-range-scaling constraints have been introduced by Jackson in [10] and Hwang in [11]. It consists in scaling the state-variable vector so as to prevent overflows or underflows during its evaluation and it imposes a condition on \mathbf{b} (it avoids big values for \mathbf{c} and small for \mathbf{b} , and *vice versa*).

Definition 4 (L_2 -scaling) A state-space realization $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)$ is said to be *L_2 -scaled* if the transfer functions from the input to each state have unitary L_2 -norms, i.e.:

$$\|e_i^\top (z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b}\|_2 = 1, \quad \forall 1 \leq i \leq n \quad (27)$$

where e_i is the column vector of appropriate dimension and with all elements being 0 except for the i^{th} element which is 1.

With proposition 1 applied to the system $(\mathbf{A}, \mathbf{b}, e_i^\top, 0)$, the L_2 -scaling constraints (27) can be expressed as:

$$(\mathbf{W}_c)_{i,i} = 1, \quad \forall 1 \leq i \leq n \quad (28)$$

where \mathbf{W}_c is the controllability Gramian of the state-space system $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)$.

In addition, it has been shown in [12] that the L_2 -scaling is not necessary to implement a realization without overflows. Two equivalent choices are possible for the implementation:

- define a binary-point position for each state (for example, the same as the input), and apply a scaling to them in order to adapt the peak values of each state to the chosen binary-point position. This scaling can also be based on a *relaxed* L_2 -scaling, as in [12].
- or set the binary-point position for each state, according to (29), to make sure that the fixed-point representation of the states avoids state-overflows:

$$\gamma_{x_i} = \beta_{x_i} - 2 - \left\lfloor \log_2 \overset{\text{up}}{x_i} \right\rfloor. \quad (29)$$

where $\overset{\text{up}}{x_i}$ is an upper-bound of the i^{th} state, given by a L_1 -norm:

$$\overset{\text{up}}{x_i} = \|e_i^\top (z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b}\|_1^{\max} u, \quad (30)$$

or estimated by a L_2 -norm [12], [13]:

$$\overset{\text{up}}{x_i} \simeq \delta \|e_i^\top (z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b}\|_2^{\max} u. \quad (31)$$

The L_2 -scaling constraints have to be applied in the first case, whereas no other constraint has to be applied in the second case..

Contrary to the L_2 -sensitivity measure M_{L_2} , the transfer function error $\sigma_{\Delta h}^2$ can be applied in a general cases and be meaningful. In the very particular case where all the coefficients have the same fixed-point representation, these two measure are equivalent. But for other cases, only the $\sigma_{\Delta h}^2$ measure is a meaningful measure of the degradation of the transfer function due to the quantization of the coefficients.

D. Scaling

Let us consider a scaling of the states: $x(k)$ is changed in $\mathbf{U}^{-1}x(k)$ with \mathbf{U} an invertible diagonal matrix. The realization \mathbf{Z}_0 is changed into \mathbf{Z}_1 :

$$\mathbf{Z}_1 = \mathbf{T}^{-1} \mathbf{Z}_0 \mathbf{T}, \quad \text{with } \mathbf{T} = \begin{pmatrix} \mathbf{U}^{-1} & \\ & \mathbf{I}_n \end{pmatrix} \quad (32)$$

Proposition 3 (Invariance to scaling) *A scaling with powers of 2 (\mathbf{U} diagonal with $\mathbf{U}_{ii} = 2^{p_i}$, $p_i \in \mathbb{Z}$, $1 \leq i \leq n$) does not change the transfer function error $\sigma_{\Delta h}^2$.*

Proof: Let $\mathcal{F}_2(x)$ denotes the fractional value of $\log_2 |x|$:

$$\mathcal{F}_2(x) \triangleq \log_2 |x| - \lfloor \log_2 |x| \rfloor \quad (33)$$

Then the operator $[\cdot]_2$ satisfies

$$[ab]_2 = [a]_2 [b]_2 2^{\lfloor \mathcal{F}_2(a) + \mathcal{F}_2(b) \rfloor} \quad (34)$$

and hence

$$\left[(\mathbf{Z}_1)_{ij} \right]_2 = \left[\mathbf{T}_{ii}^{-1} \right]_2 \left[(\mathbf{Z}_0)_{ij} \right]_2 \left[\mathbf{T}_{jj} \right]_2 \Phi_{ij} \quad (35)$$

with $\Phi_{ij} \triangleq 2^{\lfloor \mathcal{F}_2(\mathbf{T}_{ii}^{-1}) + \mathcal{F}_2(\mathbf{T}_{jj}) + \mathcal{F}_2((\mathbf{Z}_0)_{ij}) \rfloor}$. So, $\Xi_{\mathbf{Z}}|_{\mathbf{Z}_1}$ is deduced from $\Xi_{\mathbf{Z}}|_{\mathbf{Z}_0}$ by

$$\left(\Xi|_{\mathbf{Z}_1} \right)_{ij} = \left(\Xi|_{\mathbf{Z}_0} \right)_{ij} \left[\mathbf{T}_{ii}^{-1} \right]_2 \left[\mathbf{T}_{jj} \right]_2 \Phi_{ij}. \quad (36)$$

The similarity on \mathbf{Z}_0 changes the sensitivity such that

$$\frac{\partial h}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_1} = \mathbf{T}^\top \frac{\partial h}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_0} \mathbf{T}^{-1}, \quad (37)$$

then

$$\left(\frac{\partial h}{\partial \mathbf{Z}_{ij}} \Xi_{ij} \right) \Big|_{\mathbf{Z}_1} = \left(\frac{\partial h}{\partial \mathbf{Z}_{ij}} \Xi_{ij} \right) \Big|_{\mathbf{Z}_0} \frac{\left[\mathbf{T}_{ii}^{-1} \right]_2 \left[\mathbf{T}_{jj} \right]_2}{\mathbf{T}_{ii}^{-1} \mathbf{T}_{jj}} \Phi_{ij}. \quad (38)$$

Now we can remark that $\Phi_{ij} \in \{1, 2, 4\}$ and $\Phi_{ij} = 1$ if the power of 2 are used for the scaling. Also $\frac{\lfloor a \rfloor_2}{a} = 1$ if a is a power of 2. ■

IV. OPTIMAL REALIZATIONS

It is now possible to consider the optimal transfer function error problem. It consists of finding the optimal coordinate transformation \mathbf{U}_{opt} :

$$\mathbf{U}_{opt} = \arg \min_{\mathbf{U} \text{ invertible}} \sigma_{\Delta h}^2(\mathbf{U}). \quad (39)$$

Since $\sigma_{\Delta h}^2$ is invariant to power-of-2 scaling, this optimization problem has an infinite number of solutions. So it could be of interest to *normalize* all the coordinate transforms with regard to an extra consideration. For example, this could be a L_2 -scaling constraints, even if it is not necessary here.

One possible *normalization* is to apply on a realization the power-of-2 scaling \mathcal{V} :

$$\mathcal{V}_{ii} = \left\lfloor \sqrt{(\mathbf{W}_c)_{ii}} \right\rfloor_2, \quad 1 \leq i \leq n \quad (40)$$

where \mathbf{W}_c is the controllability Gramian of this realization.

Then, the *normalized* realization has the following property:

Proposition 4 *A normalized realization satisfies the relaxed- L_2 -scaling constraints, i.e.:*

$$1 \leq (\mathbf{W}_c)_{ii} < 4, \quad 1 \leq i \leq n. \quad (41)$$

This relaxed-constraints were proposed in [12] as an extension of the strict- L_2 -scaling constraints (28) that still prevents the implementation from overflows.

Proof: After the *normalization*, the new realization will have a controllability Gramian changed into $\mathbf{W}'_c = \mathcal{V}^{-1} \mathbf{W}_c \mathcal{V}^{-\top}$, hence

$$(\mathbf{W}'_c)_{ii} = \left(\frac{\sqrt{(\mathbf{W}_c)_{ii}}}{\left\lfloor \sqrt{(\mathbf{W}_c)_{ii}} \right\rfloor_2} \right)^2. \quad (42)$$

We can remark that $\forall x, 1 \leq \frac{x}{\lfloor x \rfloor_2} < 2$, so the *normalized* realization satisfies (41). ■

Of course, any other normalization is possible, but this one allows us to use the existing L_2 -scaling or relaxed L_2 -scaling methods [12], [14].

Then, the optimal problem can be defined as

Problem 2 (normalized $\sigma_{\Delta h}^2$ -optimal realization)

Considering a state-space realization (A, b, c, d) , the $\sigma_{\Delta h}^2$ -optimal realization can be found by solving the following optimization problem:

$$\mathbf{U}_{opt} = \arg \min_{\mathbf{U} \text{ invertible}} \sigma_{\Delta h}^2(\mathbf{U}\mathbf{V}), \quad (43)$$

where \mathbf{V} is a diagonal scaling matrix such that

$$\mathbf{V}_{ii} = \left[\sqrt{(\mathbf{U}^{-1}\mathbf{W}_c\mathbf{U}^{-\top})_{ii}} \right]_2 \quad (44)$$

Proof: \mathbf{V} is built in such way that the coordinate transformation $\mathcal{T} = \mathbf{U}\mathbf{V}$ applied on the realization performs the normalization. ■

Since the $\sigma_{\Delta h}^2$ measure is non smooth, this optimization problem can be solved with a global optimization method such as the Adaptive Simulated Algorithm (ASA) [15], [16]. A gradient-base method such as the quasi-Newton algorithm leads to local optima, which in practice are however not too far from the global optimum.

The FWR Toolbox⁴ was used for the numerical examples, and few minutes of computation were here required on a desktop computer.

V. NUMERICAL EXAMPLE

Let us consider the filter with coefficients given by the Matlab command `butter(4, 0.05)`, and some equivalent (in infinite precision) realizations:

\mathcal{Z}_1 : the Direct Form II,

\mathcal{Z}_2 : the balanced realization,

\mathcal{Z}_3 : the normalized $\bar{\sigma}_{\Delta h}^2$ -optimal realization (obtained with ASA and proposition 2).

The following table gives the $\bar{\sigma}_{\Delta h}^2$ measure of these different realizations. This could be compared to the *a posteriori* difference between the transfer function h and the transfer function with quantized coefficients (denoted h^\dagger). The wordlengths used are 16, 14 and 10 bits (but 10 bits are not enough to preserve the stability of the Direct Form II, that cannot be used with so few bits):

realization	$\bar{\sigma}_{\Delta h}^2$	$\ h - h^\dagger\ _2$		
		16 bits	14 bits	10 bits
\mathcal{Z}_1	$5.690e + 6$	$2.055e - 2$	0.1578	N.A.
\mathcal{Z}_2	3.693	$3.678e - 5$	$1.6994e - 4$	$3.0375e - 3$
\mathcal{Z}_3	1.439	$2.189e - 5$	$2.3148e - 5$	$1.8358e - 4$

The realization \mathcal{Z}_3 has of course the lowest transfer function error. Even with few bits, the degradation of realization \mathcal{Z}_3 remain low, compared to \mathcal{Z}_1 and \mathcal{Z}_2 .

Even if it is non-optimal, \mathcal{Z}_2 performs quite well with 16-bit implementation, but is less accurate than \mathcal{Z}_3 with less bits.

As we can remark, the coefficients of these three realizations have different magnitudes, and that legitimates the use of $\bar{\sigma}_{\Delta h}^2$. Here are the 16-bit fixed-point coefficients of \mathcal{Z}_3^\dagger (each one has a different binary-point position):

$$\mathbf{z}_3^\dagger = \begin{pmatrix} +29648.2^{-15} & +27141.2^{-18} & +20820.2^{-20} & -30467.2^{-19} & -32227.2^{-19} \\ +24569.2^{-20} & +29679.2^{-15} & +22295.2^{-17} & -31725.2^{-20} & +19083.2^{-22} \\ -31503.2^{-20} & -31152.2^{-19} & +29148.2^{-15} & +30424.2^{-22} & -32633.2^{-15} \\ +22733.2^{-17} & +21076.2^{-20} & -32727.2^{-21} & +29154.2^{-15} & -26416.2^{-31} \\ +28776.2^{-24} & -32739.2^{-22} & -25371.2^{-26} & -32767.2^{-18} & +16771.2^{-29} \end{pmatrix}$$

⁴sources available at <http://fwrtoolbox.gforge.inria.fr>

VI. CONCLUSION

This paper has presented the optimal realization problem in fixed-point context and exhibited a new measure to evaluate the coefficients roundoff errors. Compared to the classical L_2 -sensitivity measure, the transfer function error is a meaningful measure for every realization, even for non- L_2 -scaled ones.

A normalization procedure used to solve the optimal realization problem has been presented with a numerical example, but there is still further work to be done, specially to develop an *ad hoc* optimization algorithm.

This measure could be easily adapted to the Specialized Implicit Framework [17] that allows to encompass all the existing structures (direct forms, cascade/parallel decompositions, lattices, δ or ρ -operator based realizations, etc.). Some other measures, like the pole-sensitivity measure, could be extended to consider fixed-point coefficients.

REFERENCES

- [1] M. Gevers and G. Li, *Parametrizations in Control, Estimation and Filtering Problems*. Springer-Verlag, 1993.
- [2] T. Hilaire, D. Ménard, and O. Sentieys, "Bit accurate roundoff noise analysis of fixed-point linear controllers," in *Proc. IEEE Int. Symposium on Computer-Aided Control System Design (CACSD'08)*, Sept. 2008.
- [3] S. Hwang, "Minimum uncorrelated unit noise in state-space digital filtering," *IEEE Trans. on Acoust., Speech, and Signal Processing*, vol. 25, no. 4, pp. 273–281, August 1977.
- [4] C. Mullis and R. Roberts, "Synthesis of minimum roundoff noise fixed point digital filters," in *IEEE Transactions on Circuits and Systems*, vol. CAS-23, no. 9, September 1976.
- [5] T. Hinamoto, S. Yokoyama, T. Inoue, W. Zeng, and W. Lu, "Analysis and minimization of L_2 -sensitivity for linear systems and two-dimensional state-space filters using general controllability and observability gramians," in *IEEE Transactions on Circuits and Systems, Fundamental Theory and Applications*, vol. 49, no. 9, september 2002.
- [6] V. Tavşanoğlu and L. Thiele, "Optimal design of state-space digital filters by simultaneous minimization of sensibility and roundoff noise," in *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. CAS-31, October 1984.
- [7] T. Hilaire and P. Chevrel, "On the compact formulation of the derivation of a transfer matrix with respect to another matrix," INRIA, Tech. Rep. RR-6760, 2008.
- [8] R. Rocher, D. Ménard, N. Hervé, and O. Sentieys, "Fixed-point configurable hardware components," *EURASIP Signal of Embedded Systems*, no. 1, p. 20, January 2006.
- [9] A. Sripad and D. Snyder, "A necessary and sufficient condition for quantization error to be uniform and white," in *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. 5, no. 10, 1977, pp. 442–448.
- [10] L. Jackson, "Roundoff-noise analysis for fixed-point digital filters realized in cascade or parallel form," *Audio and Electroacoustics, IEEE Transactions on*, vol. 18, no. 2, pp. 107–122, June 1970.
- [11] S. Hwang, "Dynamic range constraint in state-space digital filtering," in *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. 23, 1975, pp. 591–593.
- [12] T. Hilaire, "Low parametric sensitivity realizations with relaxed l_2 -dynamic-range-scaling constraints," *IEEE Trans. on Circuits & Systems II*, vol. 56, no. 7, pp. 590–594, July 2009.
- [13] K. Parhi, *VLSI Digital Signal Processing Systems: Design and Implementation of Digital Controllers*. John Wiley & Sons, 1999.
- [14] T. Hinamoto, H. Ohnishi, and W.-S. Lu, "Minimization of l_2 sensitivity of one- and two dimensional state-space digital filters subject to l_2 -dynamic-range-scaling constraints," *IEEE Trans. on Circuits and Systems-II*, vol. 52, no. 10, pp. 641–645, October 2005.
- [15] L. Ingber, "Adaptive simulated annealing (ASA): Lessons learned," *Control and Cybernetics*, vol. 25, no. 1, pp. 33–54, 1996.
- [16] S. Chen and B. Luk, "Adaptive Simulated Annealing for optimization in signal processing applications," *Signal Processing*, vol. 79, pp. 117–128, 1999.
- [17] T. Hilaire, P. Chevrel, and J. Whidborne, "A unifying framework for finite wordlength realizations," *IEEE Trans. on Circuits and Systems*, vol. 8, no. 54, pp. 1765–1774, August 2007.