

# NEW $L_2$ -DYNAMIC-RANGE-SCALING CONSTRAINTS FOR LOW PARAMETRIC SENSITIVITY REALIZATIONS

Thibault Hilaire

Vienna University of Technology, Austria  
Institute of Communications and Radio-Frequency Engineering  
thibault.hilaire@nt.tuwien.ac.at

## ABSTRACT

This paper presents a new dynamic-range scaling for the implementation of filters/controllers in state-space form. Specific fixed-point considerations allow us to relax the classical  $L_2$ -scaling constraints while still preserving the implementation from overflows. It gives more degrees of freedom for the optimal  $L_2$ -parametric sensitivity problem. The underlying constrained problem is converted into an unconstrained problem for which a solution can be provided. This leads to realizations which are still scaled but less sensitive.

## 1. INTRODUCTION

The majority of signal processing (or control) systems is implemented in digital general purpose processors, DSPs<sup>1</sup>, FPGAs<sup>2</sup>, 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 [3, 6, 10, 14] for roundoff noises.

It is also well known that these Finite Word Length (FWL) effects depend on the structure of the realization. This motivates to investigate the coefficient sensitivity minimization problem. It has been widely studied since Thiele published [16, 17] and the definition of a tractable input-output sensitivity norm (the  $L_1/L_2$ -sensitivity). This work has been extended with a more natural and reasonable measure, the  $L_2$ -sensitivity ([3, 8]).

The dynamic-range-scaling constraints have been introduced in [11] and [9] to prevent overflow and underflow during the evaluation of the state-vector, and as well as the state and criteria normalization. These constraints have to be considered in the  $L_2$ -sensitivity minimization problem, for which [7] proposes an efficient quasi-Newton algorithm to solve it.

This paper investigates the  $L_2$ -dynamic-range-scaling problem by considering concrete fixed-point implementation of state-space realizations. It reveals that the classical  $L_2$ -scaling is only a sufficient condition to prevent overflows and thus it can be slightly relaxed in order to extend the degrees of freedom for the optimization process. New *relaxed*- $L_2$ -dynamic-range-scalings are then presented with respect to the described computational scheme. Finally, the  $L_2$ -sensitivity minimization problem with relaxed  $L_2$ -scaling constraints is solved. A numerical example illustrates that the proposed constraints can offer reduced  $L_2$ -sensitivity with overflow protection.

## 2. $L_2$ -SENSITIVITY ANALYSIS

Let  $(\mathbf{A}, \mathbf{b}, \mathbf{c}, d)$  be a stable, controllable and observable linear discrete time SISO<sup>3</sup> 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 uncertainty 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 of interest to consider the sensitivity of the transfer function with respect to the coefficients, based on the following definitions.

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

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

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

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

$$\|\mathbf{H}\|_p \triangleq \left( \frac{1}{2\pi} \int_0^{2\pi} \|\mathbf{H}(e^{j\omega})\|_F^p d\omega \right)^{\frac{1}{p}} \quad (4)$$

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

Gevers and Li [3] have proposed the  $L_2$ -sensitivity measure to evaluate the coefficient roundoff errors. It is defined by

$$M_{L_2} \triangleq \left\| \frac{\partial h}{\partial \mathbf{A}} \right\|_2^2 + \left\| \frac{\partial h}{\partial \mathbf{b}} \right\|_2^2 + \left\| \frac{\partial h}{\partial \mathbf{c}} \right\|_2^2 + \left\| \frac{\partial h}{\partial d} \right\|_2^2 \quad (5)$$

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

$$\mathbf{F}(z) \triangleq (z\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{b}, \quad \mathbf{G}(z) \triangleq \mathbf{c}(z\mathbf{I}_n - \mathbf{A})^{-1}. \quad (6)$$

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

[16] ( $\left\| \frac{\partial h}{\partial \mathbf{A}} \right\|_1^2$  instead of  $\left\| \frac{\partial h}{\partial \mathbf{A}} \right\|_2^2$  in (5)).

This work has been funded in parts by the NFN SISE project (National Research Network "Signal and Information Processing in Science and Engineering").

<sup>1</sup>Digital Signal Processors

<sup>2</sup>Field Programmable Gate-Array

<sup>3</sup>Single Input Single Output

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

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  $T$ , and is denoted  $M_{L_2}(T)$ . 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  $T_{opt}$  that minimizes  $M_{L_2}$ :*

$$T_{opt} = \arg \min_{T \text{ invertible}} M_{L_2}(T). \quad (7)$$

[3] shows that the problem has one unique solution. Hence, for example, a gradient method can be used to solve it.

### 3. $L_p$ -DYNAMIC-RANGE SCALING

The  $L_p$ -dynamic-range-scaling constraints have been introduced by Jackson in [11] and Hwang in [9]. It consists in scaling the state-variable vector such that overflows or underflows during its evaluation are prevented.

**Definition 3 ( $L_p$ -scaling)** *A state-space realization  $(A, b, c, d)$  is said to be  $L_p$ -scaled if the  $L_p$ -norms of the transfer functions from the input to each state are set to 1, i.e.:*

$$\left\| e_i^\top (zI_n - A)^{-1} b \right\|_p = 1, \quad \forall 1 \leq i \leq n \quad (8)$$

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

Let  $\bar{u}^{\max}$  denote the maximum value of the input  $u$ :

$$\bar{u}^{\max} \triangleq \max_{k \in \mathbb{N}} |u(k)|. \quad (9)$$

The  $L_1$ -scaling guarantees that the dynamic of each state  $x_i$  is lower than  $\bar{u}^{\max}$ , whereas the  $L_2$ -scaling guarantees that the variance of each state is unitary for a unit-variance centered white noise input.  $L_2$ -scaling doesn't completely prevent overflow as does  $L_1$ , but it is less conservative and more realistic, so it is widely used [15].

**Proposition 1** *The  $L_2$ -norms of the transfer functions from the input to each state are given by*

$$\left\| e_i^\top (zI_n - A)^{-1} b \right\|_2 = \sqrt{(\mathbf{W}_c)_{i,i}} \quad (10)$$

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

$$\mathbf{W}_c = A\mathbf{W}_cA^\top + bb^\top. \quad (11)$$

*Proof:*

A classical result on  $L_2$ -norm gives the  $L_2$ -norm of a state-space system  $H := (K, L, M, N)$ :

$$\|H\|_2^2 = \text{tr}(NN^\top + M\mathbf{W}_cM^\top) \quad (12)$$

where  $\text{tr}(\cdot)$  is the trace operator and  $\mathbf{W}_c$  is the controllability Gramian of the system.

This is here applied to the system  $(A, b, e_i^\top, 0)$ . ■

### Problem 2 (sensitivity problem with $L_2$ -scaling constraints)

*The optimal  $L_2$ -sensitivity problem with  $L_2$ -norm dynamic-range-scaling constraints can be formulated as the optimization problem 1, subject to the constraints*

$$(\mathbf{W}_c)_{i,i} = 1, \quad \forall 1 \leq i \leq n. \quad (13)$$

This constrained problem can be transformed into an unconstrained problem, and an optimization algorithm can be used to solve it [7].

## 4. FIXED-POINT IMPLEMENTATION

### 4.1 Fixed-point representation

In this paper, the notation  $(\beta, \gamma)$  is used for the fixed-point representation of a variable or coefficient ( $2$ 's complement scheme), according to Figure 1.  $\beta$  is the total wordlength of the representation in bits, whereas  $\gamma$  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).  $\beta$  and  $\gamma$  will be suffixed by the variable/coefficient they refer to.

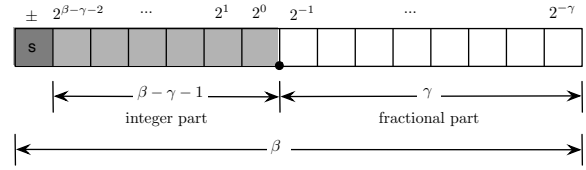


Figure 1: Fixed-point representation

To represent a value  $x$  without overflow, a fixed-point representation  $(\beta_x, \gamma_x)$  may satisfy:

$$\beta_x - \gamma_x - 1 \geq \lfloor \log_2 |x| \rfloor + 1 \quad (14)$$

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).

An important fixed-point issue is to find a valid fixed-point representation, such that (14) is satisfied for all values which  $x$  can assume during the execution of the algorithm.

### 4.2 State-overflow

**Definition 4 (State-overflow)** *The overflow of the state variables  $(x_i)_{1 \leq i \leq n}$  can be strictly avoided if all the values  $x_i(k)$  are within the range allowed by the fixed-point representations, i.e. ( $1 \leq i \leq n$ )*

$$\forall k, \quad -2^{\beta_{x_i} - \gamma_{x_i} - 1} \leq x_i(k) < 2^{\beta_{x_i} - \gamma_{x_i} - 1}. \quad (15)$$

*The overflows are avoided if the binary-point position of each state is carefully chosen, such that*

$$\gamma_{x_i} = \beta_{x_i} - 2 - \left\lfloor \log_2 \bar{x}_i^{\max} \right\rfloor, \quad (16)$$

where  $\bar{x}_i^{\max}$  is the maximum magnitude for the  $i^{\text{th}}$  state:

$$\bar{x}_i^{\max} \triangleq \max_{k \in \mathbb{N}} |x_i(k)|. \quad (17)$$

**Remark 1** At least,  $\gamma_{x_i}$  should satisfy  $\gamma_{x_i} \leq \beta_{x_i} - 2 - \left\lfloor \log_2 \bar{x}_i^{\max} \right\rfloor$ , but the greater  $\gamma_{x_i}$  is, the more accurate is the fixed-point format.

However only upper bounds can be computed. A first upper bound  $\bar{x}_i^{\text{up}}$  can be obtained by an  $L_1$ -norm:

$$\bar{x}_i^{\text{up}} = \left\| e_i^\top (z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b} \right\|_1^{\max} u, \quad (18)$$

and a second one can be estimated by an  $L_2$ -norm [15]:

$$\bar{x}_i^{\text{up}} \simeq \delta \left\| e_i^\top (z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b} \right\|_2^{\max} u. \quad (19)$$

Here, the parameter  $\delta$  can be interpreted as a representation of the number of standard deviations of  $x_i$ , if the input is unit-variance white centered noise ( $\delta \geq 1$ ). Since the  $L_2$ -norm in (19) doesn't give a strict bound (contrary to (18)),  $\delta$  can be seen as a *safety* parameter [15].

Finally, these upper bounds are used to define the binary-point positions:

$$\gamma_{x_i} = \beta_{x_i} - 2 - \left\lfloor \log_2 \bar{x}_i^{\text{up}} \right\rfloor. \quad (20)$$

In general, the  $L_1$  and  $L_2$  estimations of  $\bar{x}_i^{\text{up}}$  approximately leads to the same binary-point position, with 1 or 2 bits deviation. However, since the  $L_2$ -norm is more tractable (with proposition 1) and the  $L_1$ -norm too conservative ( $\bar{x}_i^{\text{max}} \ll \bar{x}_i^{\text{up}}$ ), in practice (19) is used, with  $\delta = 1$ . After implementation, a simulation-based estimation like in [1] or [12] can also be used to verify *in situ* the peak values and the binary point positions, according to the inputs.

### 4.3 Computational scheme

In order to implement a realization without overflows, two equivalent choices are possible:

- set the binary-point position for each state, according to (20), to make sure that the fixed-point representation of the states avoids state-overflows;
- or define a binary-point position for each state, and apply a scaling to them in order to adapt the peak values of each state to the chosen binary-point position.

Here, we here focus on the 2<sup>nd</sup> choice, referring to dynamic-range-scaling constraints.

Let us consider in detail the fixed-point implementation of the system given in (1). It leads to  $(n+1)$  scalar products to be evaluated, of the form:

$$S = \sum_{i=1}^N p_i q_i \quad (21)$$

where the  $(p_i)$  are given coefficients and  $(q_i)$  are bounded variables. To avoid bit-shift operations between each addition in the evaluation of eq. (21), the binary-point positions of each partial product of the sum should be equal.

Then, two computational schemes are possible: the *Roundoff After Multiplication* scheme, where shifts are added after each product to align the operands of the sum ( $p_i q_i$  is implemented as  $(p_i' * q_i') \gg d_i$ ) and the *Roundoff Before Multiplication* scheme, where the required shifts are reported into the coefficients ( $p_i q_i$  is implemented as  $(p_i' \gg d_i) * q_i'$ ). See Figure 2.

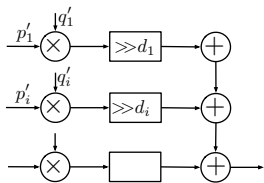


Figure 2: Scalar product with *Roundoff After Multiplication* scheme

The main idea of the scaling is to scale each variable  $(q_i)$  such that the shifts ( $d_i = 0, \forall i$ ) are prevented. In fixed-point representation, the scaling only implies that all the  $(q_i)$  have a common

format, and so have the  $(p_i)$ . See [2, 6] for more details on implementation schemes.

Applied to the state-space realization (1), this yields that all the states must have the same binary-point position as the input and the coefficients  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$  and  $d$ .

Besides, since they have the same fractional part  $\gamma$ , their quantization's errors  $\Delta\mathbf{A}$ ,  $\Delta\mathbf{b}$ ,  $\Delta\mathbf{c}$  and  $\Delta d$  have the same magnitude  $2^{-\gamma-1}$ , and the  $L_2$ -sensitivity measure represents a meaningful bound on the transfer function error  $\Delta h$  (considering that  $h$  is shifted in  $h + \Delta h$  by the coefficients' quantization):

$$\|\Delta h\|_2^2 \leq \left\| \frac{\partial h}{\partial \mathbf{A}} \times \Delta \mathbf{A} \right\|_2^2 + \left\| \frac{\partial h}{\partial \mathbf{b}} \times \Delta \mathbf{b} \right\|_2^2 + \left\| \frac{\partial h}{\partial \mathbf{c}} \times \Delta \mathbf{c} \right\|_2^2 + \left\| \frac{\partial h}{\partial d} \times \Delta d \right\|_2^2 \quad (22)$$

$$\leq 2^{-2(\gamma+1)} M_{L_2} \quad (23)$$

### 4.4 New $L_2$ -scaling constraints

Taken this into consideration, the overflows will be avoided by setting the same binary-point position for the states and the input, and by applying an appropriate scaling on the states such that the constraints (16) are satisfied.

Compared to strict  $L_2$ -scaling where the states must satisfy  $\bar{x}_i^{\text{max}} = \bar{u}^{\text{max}}$ , here, the constraints are relaxed (but still restrictive enough to guarantee the protection against overflow) and replaced by  $\gamma_{x_i} = \gamma_u$ .

**Proposition 2 (relaxed- $L_2$ -scaling constraints)** *Since the input and the states may have the same binary-point position, the  $L_2$ -scaling constraints (13) are now transformed into*

$$\frac{2^{2\alpha_i}}{\delta^2} \leq (\mathbf{W}_c)_{i,i} < 4 \frac{2^{2\alpha_i}}{\delta^2}, \quad \forall 1 \leq i \leq n \quad (24)$$

where

$$\alpha_i \triangleq \beta_{x_i} - \beta_u - \mathcal{F}_2 \left( \frac{\bar{u}^{\text{max}}}{u} \right) \quad (25)$$

and  $\mathcal{F}_2(x)$  is defined as the fractional value of  $\log_2(x)$ :

$$\mathcal{F}_2(x) \triangleq \log_2(x) - \lfloor \log_2(x) \rfloor \quad (26)$$

*Proof:*

The binary-point position of the input is set to  $\gamma_u = \beta_u - 2 - \left\lfloor \log_2 \frac{\bar{u}^{\text{max}}}{u} \right\rfloor$ . Hence, with (20), the constraints  $\gamma_u = \gamma_{x_i}$  lead to

$$\beta_u - \left\lfloor \log_2 \frac{\bar{u}^{\text{max}}}{u} \right\rfloor = \beta_{x_i} - \left\lfloor \log_2 \left( \delta \left\| e_i^\top (z\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{b} \right\|_2^{\max} u \right) \right\rfloor$$

and

$$\left\lfloor \log_2 \left( \delta \sqrt{(\mathbf{W}_c)_{i,i}} \right) + \mathcal{F}_2 \left( \frac{\bar{u}^{\text{max}}}{u} \right) \right\rfloor = \beta_{x_i} - \beta_u \quad (27)$$

And finally

$$2^{\alpha_i} \leq \delta \sqrt{(\mathbf{W}_c)_{i,i}} < 2^{\alpha_i+1} \quad (28)$$

**Remark 2** For microcontroller or DSP implementations (contrary to FPGA or some ASIC implementations), the wordlength of all variables is most of the time equal, i.e.  $\beta_u = \beta_{x_i}$  ( $1 \leq i \leq n$ ).  $\delta$  is set to unity (as for classical  $L_2$ -scaling constraints).

In best case,  $\frac{\bar{u}^{\text{max}}}{u}$  is equal to a power of 2 ( $2^p$ ,  $p \in \mathbb{Z}$ ), and the relaxed- $L_2$ -scaling constraints (24) become:

$$1 \leq (\mathbf{W}_c)_{i,i} < 4, \quad \forall 1 \leq i \leq n, \quad (29)$$

and in worst case,  $u^{\max}$  is equal to the representable value immediately lower than a power of 2 ( $2^p - 2^{\beta_u - 2 - p}$ ), and the constraints become:

$$\frac{1}{4} \leq (\mathbf{W}_c)_{i,i} < 1, \quad \forall 1 \leq i \leq n \quad (30)$$

(in the first case,  $p + 1$  bits are used for the integer part of the states, whereas  $p$  bits only are used in the second case).

It is important to remark that these new constraints allow more freedom for the scaling and introduce a new degree of freedom for the search for optimal realizations. Even though not considered in this paper, moreover it could give more freedom for the minimization of the roundoff noise power.

### 5. OPTIMAL $L_2$ -SENSITIVITY REALIZATION WITH RELAXED $L_2$ -NORM DYNAMIC-RANGE-SCALING CONSTRAINTS

Then, these relaxed constraints can be applied to a new sensitivity problem:

**Problem 3 (relaxed sensitivity problem)** *The optimal  $L_2$ -sensitivity problem with relaxed  $L_2$ -norm dynamic-range-scaling constraints can be expressed in the form of the constrained problem 2 subject to constraints in (24).*

This constrained problem can be solved by scaling the system after each transformation in order to ensure that the constraints are met:

**Proposition 3 (a posteriori relaxed scaling)** *Considering a state-space realization, it is possible to a posteriori scale it with a diagonal transformation matrix  $\mathbf{T}$  given by*

$$\mathbf{T}_{i,i} = \delta \sqrt{(\mathbf{W}_c)_{i,i}} 2^{-\mathcal{F}_2(\delta \sqrt{(\mathbf{W}_c)_{i,i}}) - \alpha_i}, \quad (31)$$

such that the constraints (24) are satisfied. Moreover, it is possible to build all the transformation matrices that meet the constraints (24): Let us consider an invertible matrix  $\mathbf{U} \in \mathbb{R}^{n \times n}$ , then the transformation matrix  $\mathbf{T} = \mathbf{U}\mathbf{V}$  with  $\mathbf{V}$  diagonal such that:

$$\mathbf{V}_{i,i} = \delta \sqrt{(\mathbf{U}^{-1} \mathbf{W}_c \mathbf{U}^{-\top})_{i,i}} 2^{-\mathcal{F}_2(\delta \sqrt{(\mathbf{U}^{-1} \mathbf{W}_c \mathbf{U}^{-\top})_{i,i}}) - \alpha_i} \quad (32)$$

produces the relaxed- $L_2$ -scaling.

*Proof:*

$\mathcal{F}_2$  acts as a modulo operator. For  $x \in \mathbb{R}$ ,  $\bar{x} \triangleq 2^{\mathcal{F}_2(x) + a}$  is such that  $2^a \leq \bar{x} < 2^{a+1}$ .

Since the constraints (24) are equal to

$$2^{\alpha_i} \leq \delta \sqrt{(\mathbf{W}_c)_{i,i}} < 2^{\alpha_i + 1} \quad (33)$$

and  $\mathbf{T}$  transforms  $(\mathbf{W}_c)_{i,i}$  into  $\mathbf{T}_{i,i}^{-2} (\mathbf{W}_c)_{i,i}$ , then  $\mathbf{T}_{i,i}$  has to be of the form:

$$\delta \mathbf{T}_{i,i}^{-1} \sqrt{(\mathbf{W}_c)_{i,i}} = 2^{\mathcal{F}_2(\delta \sqrt{(\mathbf{W}_c)_{i,i}}) + \alpha_i} \quad (34)$$

Thus, the optimization problem is given by

$$\mathbf{U}_{opt} = \arg \min_{\substack{\mathbf{U} \text{ invertible} \\ \mathbf{V} \text{ defined by (32)}}} M_{L_2}(\mathbf{U}\mathbf{V}). \quad (35)$$

This was implemented in the FWR toolbox<sup>4</sup> for Matlab, with `fminsearch`, and `fminunc` functions, and they both give same results with similar numbers of iterations.

Of course, the use of the matrice  $\mathbf{V}$ , that is merely used to eliminate the constraints and solve an unconstrained minimization problem, increases the degree of non-linearity for the objective function to minimize. However, this seems not to be a problem, since in our tests, the optimal realizations found seem to be global optima.

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

## 6. EXAMPLE

Let us consider the following digital filter, given by its transfer function<sup>5</sup>:

$$H(z) = \frac{4.297e - 3z^4 - 8.595e - 3z^2 + 4.297e - 3}{z^4 - 3.803z^3 + 5.550z^2 - 3.678z + 0.9355} \quad (36)$$

and its multiple equivalent (in infinite precision) realizations:

- $\mathcal{R}_1$  is the Direct Form II given by (37).
- $\mathcal{R}_2$  is the optimal  $L_2$ -scaled realization (solution of problem 2). The numerical values are given by (38).
- $\mathcal{R}_3$  is the optimal relaxed- $L_2$ -scaled realization (problem 3), with  $u^{\max}$  a power of 2, and  $\delta = 1$ . It is obtained with proposition 3. The numerical values are given by (39).

The balanced realization (known to be close to the unscaled  $M_{L_2}$ -optimal realization) was chosen as the starting point for the optimization. The numerical results are presented in the appendix and the pseudo-code associated to the realization  $\mathcal{R}_3$  is presented in algorithm 1 (16 bits are considered for the input, output, states and coefficients, 32 bits for the accumulator (no guard bits), and  $u^{\max} = 16$ ).

The following table gives the  $M_{L_2}$  sensitivities of these different realizations:

realization	$M_{L_2}$ sensitivity
$\mathcal{R}_1$	1.690e + 09
$\mathcal{R}_2$	5247.9
$\mathcal{R}_3$	5222.1

In this example, the relaxed  $L_2$ -scaled realization  $\mathcal{R}_3$  achieves lower sensitivity than the strict  $L_2$ -scaled optimal realization  $\mathcal{R}_2$  while protecting implementation from overflows.

But it is not always the case : if we consider the example in [7], the optimal relaxed- $L_2$ -scaled realization satisfies  $(\mathbf{W}_c)_{i,i} = 1$  and is then also a strict  $L_2$ -scaled realization. This depends on the diagonal terms of the controllability Gramians of the (non scaled) optimal realization.

It is also interesting to notice that a good estimation of  $u^{\max}$  (if it is not a power of 2) can allow to achieve lower sensitivity by moving the constraints (it could also be the case for the example in [7]).

## 7. CONCLUSION

This paper has presented the  $L_2$ -sensitivity minimization problem and the associated  $L_2$ -scaling constraints. These constraints that prevent from overflows have been considered with concrete fixed-point implementation schemes. Novel  $L_2$ -dynamic-range constraints have been exhibited.

Even if the goal of this paper is not a detailed optimization algorithm like in [7], a proposition to solve the constrained optimization problem has been exhibited and applied on a numerical example.

These relaxed constraints could also be very important for some other realizations, like  $\delta$ -operator state-space, the  $\rho$ -Direct Form II transposed [13] and also classical filtering structures. For these realizations where a parameter  $\Delta$  should be used to achieve the  $L_2$ -scaling, a relaxed- $L_2$ -scaling permits to fix this parameter as a power of 2, in order to decrease the amount of computations.

To apply this work to other classical structures, it will be soon extended to the Specialized Implicit Framework [5] that allows to encompass existing structures in an implicit state-space form.

## REFERENCES

- [1] P. Belanovic and M. Rupp. Automated floating-point to fixed-point conversion with the fixify environment. *Rapid System Prototyping, 2005. (RSP 2005). The 16th IEEE International Workshop on*, pages 172–178, June 2005.

<sup>5</sup>Due to a lack of space, only 4 digits are given, but more may be required to completely define the system.

- [2] G. Constantinides, P. Cheung, and W. Luk. *Synthesis And Optimization Of DSP Algorithms*. Kluwer Academic Publishers, 2004.
- [3] M. Gevers and G. Li. *Parametrizations in Control, Estimation and Filtering Problems*. Springer-Verlag, 1993.
- [4] T. Hilaire and P. Chevrel. On the compact formulation of the derivation of a transfer matrix with respect to another matrix. Technical Report RR-6760, INRIA, 2008.
- [5] T. Hilaire, P. Chevrel, and J. Whidborne. A unifying framework for finite wordlength realizations. *IEEE Trans. on Circuits and Systems*, 8(54), August 2007.
- [6] T. Hilaire, D. Ménard, and O. Sentieys. Bit accurate roundoff noise analysis of fixed-point linear controllers. In *Proc. IEEE International Symposium on Computer-Aided Control System Design (CACSD'08)*, September 2008.
- [7] T. Hinamoto, H. Ohnishi, and W.-S. Lu. Minimization of l2 sensitivity of one- and two dimensional state-space digital filters subject to l2-dynamic-range-scaling constraints. *IEEE Trans. on Circuits and Systems-II*, 52(10):641–645, October 2005.
- [8] T. Hinamoto and Y. Sugie. L2-sensitivity analysis and minimization of 2-d separable-denominator state-space digital filters. *Signal Processing, IEEE Transactions on*, 50(12):3107–3114, Dec 2002.
- [9] S. Hwang. Dynamic range constraint in state-space digital filtering. In *IEEE Trans. on Acoustics, Speech and Signal Processing*, volume 23, pages 591–593, 1975.
- [10] S. Hwang. Minimum uncorrelated unit noise in state-space digital filtering. *IEEE Trans. on Acoust., Speech, and Signal Processing*, 25(4):273–281, August 1977.
- [11] L. Jackson. Roundoff-noise analysis for fixed-point digital filters realized in cascade or parallel form. *Audio and Electroacoustics, IEEE Transactions on*, 18(2):107–122, June 1970.
- [12] S. Kim, K. Kum, and W. Sung. Fixed-point optimization utility for C and C++ based digital signal processing programs. *IEEE Transactions on Circuits and Systems*, 45(11):1455–1464, November 1998.
- [13] G. Li and Z. Zhao. On the generalized DFII structure and its state-space realization in digital filter implementation. *IEEE Trans. on Circuits and Systems*, 51(4):769–778, April 2004.
- [14] C. Mullis and R. Roberts. Synthesis of minimum roundoff noise fixed point digital filters. In *IEEE Transactions on Circuits and Systems*, volume CAS-23, September 1976.
- [15] K. Parhi. *VLSI Digital Signal Processing Systems: Design and Implementation of Digital Controllers*. Number ISBN 0-471-24186-5. John Wiley & Sons, 1999.
- [16] 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*, volume CAS-31, October 1984.
- [17] L. Thiele. Design of sensitivity and round-off noise optimal state-space discrete systems. *Int. J. Circuit Theory Appl.*, 12:39–46, 1984.

$$\mathbf{A}_2 = \begin{pmatrix} 0.9386 & -0.2477 & -0.0427 & 0.0077 \\ 0.2525 & 0.9675 & -0.0010 & -0.0003 \\ 0.0228 & -0.0031 & 0.9301 & -0.2524 \\ 0.0039 & -3.071e-05 & 0.2517 & 0.9669 \end{pmatrix}, \quad \mathbf{b}_2 = \begin{pmatrix} -0.0652 \\ -0.0193 \\ 0.3130 \\ 0.0471 \end{pmatrix},$$

$$\mathbf{c}_2 = (-1.1151 \quad 0.2002 \quad -0.1789 \quad 0.0729), \quad d_2 = 4.2974e-3 \quad (38)$$

$$\mathbf{A}_3 = \begin{pmatrix} 0.9300 & -0.2504 & 0.0290 & -0.0030 \\ 0.2528 & 0.9671 & 0.0027 & -0.0059 \\ -0.0483 & 0.0091 & 0.9385 & -0.2466 \\ 0.0001 & 0.0056 & 0.2527 & 0.9674 \end{pmatrix}, \quad \mathbf{b}_3 = \begin{pmatrix} 0.6289 \\ 0.0665 \\ -0.1326 \\ -0.0440 \end{pmatrix},$$

$$\mathbf{c}_3 = (-0.0872 \quad 0.0335 \quad -0.5567 \quad 0.1104), \quad d_3 = 4.2974e-3 \quad (39)$$

**Input:**  $u$ : 16 bits integer  
**Output:**  $y$ : 16 bits integer  
**Data:**  $xn, xnp$ : array [1..4] of 16 bits integers  
**Data:**  $Acc$ : 32 bits integer  
**begin**

```

// Intermediate variables
Acc ← (xn(1) * 31209);
Acc ← Acc + (xn(2) * -8303);
Acc ← Acc + (xn(3) * -600);
Acc ← Acc + (xn(4) * 78);
Acc ← Acc + (u * -806);
xnp(1) ← Acc >> 15;
Acc ← (xn(1) * 8304);
Acc ← Acc + (xn(2) * 31690);
Acc ← Acc + (xn(3) * -77);
Acc ← Acc + (xn(4) * 10);
Acc ← Acc + (u * -104);
xnp(2) ← Acc >> 15;
Acc ← (xn(1) * 2398);
Acc ← Acc + (xn(2) * -315);
Acc ← Acc + (xn(3) * 30009);
Acc ← Acc + (xn(4) * -8149);
Acc ← Acc + (u * 1523);
xnp(3) ← Acc >> 15;
Acc ← (xn(1) * 313);
Acc ← Acc + (xn(2) * -41);
Acc ← Acc + (xn(3) * 8147);
Acc ← Acc + (xn(4) * 31711);
Acc ← Acc + (u * 198);
xnp(4) ← Acc >> 15;
// Outputs
Acc ← (xn(1) * -25782);
Acc ← Acc + (xn(2) * 3381);
Acc ← Acc + (xn(3) * -12182);
Acc ← Acc + (xn(4) * 1575);
Acc ← Acc + (u * 18);
y ← Acc >> 15;
// Permutations
xn ← xnp;

```

**end**

**Algorithm 1:** Numerical fixed-point algorithm of realization  $\mathcal{R}_3$

## A. NUMERICAL RESULTS

$$\mathbf{A}_1 = \begin{pmatrix} 3.8031 & -1.3875 & 0.9196 & -0.4678 \\ 4 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0.5 & 0 \end{pmatrix}, \quad \mathbf{b}_1 = \begin{pmatrix} 0.1250 \\ 0 \\ 0 \\ 0 \end{pmatrix},$$

$$\mathbf{c}_1 = (0.1307 \quad -0.0649 \quad 0.0316 \quad 0.0011), \quad d_1 = 4.2974e-3 \quad (37)$$