

# Determining Fixed-Point Formats for a Digital Filter Implementation using the Worst-Case Peak Gain measure

Anastasia Volkova, Thibault Hilaire, Christoph Lauter  
Sorbonne Universités, UPMC Univ Paris 06, UMR 7606, LIP6, France  
Email: first\_name.last\_name@lip6.fr

**Abstract**—In this article, we focus on the Fixed-Point implementation of Linear Time Invariant (LTI) filters in state-space representation. For that purpose, we give an algorithm to determine the Fixed-Point Formats of all the involved variables (states and outputs). For the sake of generality, the algorithm works in the case of Multiple Inputs Multiple Outputs filters.

The computational errors in the intermediate steps of the filter evaluation as well as their accumulation over time are fully taken into account. We handle several rounding modes (round to nearest and truncation) for two’s-complement-based Fixed-Point Arithmetic. Our approach is fully rigorous in the way that the output Fixed-Point formats are shown to be free of overflows and we do not use any (non-exhaustive) filter simulation steps but proofs.

## I. INTRODUCTION

The implementation of a digital filter can be considered as a path from filter specifications to a physical implementation (DSP device). This process includes generation of a transfer function, which defines the relation between the inputs and outputs of a filter. Further one needs to choose a filter structure to be eventually implemented, i.e. a computational scheme, a specific order of filter evaluation. Finally, there is software and hardware code generation.

However, on each step various issues arise. First of all, every transfer function must be discretized, which implies the modification of filter. Secondly, while in infinite precision filter structures (e.g. Direct Forms [1], State-Space [2], Wave [3] etc.) are equivalent, they are no longer the same objects in finite precision. Since the computational schemes are different, the round-off errors and consequently quality of structures vary. Finally, software and hardware implementations are performed under constraints, e.g. power consumption, area, output error, etc. And on each step of the filter implementation the degradation due to finite precision must be taken into account in order to produce a reliable implementation.

To unify the above-described process of filter implementation for any Linear Time-Invariant (LTI) filter and provide reliable implementation, we develop an automatized filter generator. Its work-flow is described in the Fig. 1.

In this paper we focus on one of the various steps: the Fixed-Point (FxP) Algorithm generation step and wordlength optimization. On this step we already know target specifications

This work has been sponsored by the *Agence nationale de la recherche* grant ANR-13-INSE-0007-02 MetaLibm.

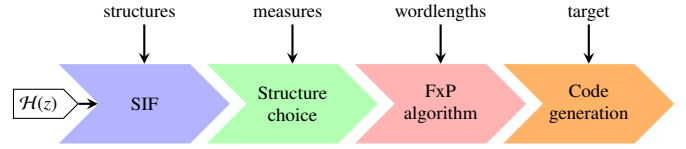


Fig. 1. Automatic filter generator flow.

and software implementation constraints. Therefore, during this optimization, we naturally look for a trade-off between area (for hardware implementation) and computational errors. If we take more bits, the implementation cost arises. If we take less bits, there exist a risk of an overflow. Depending on the application, one or even both criteria can be neglected. However, in the general case on each step of the optimization procedure we need to determine *rigorously* the Fixed-Point Formats (FxPF), i.e. the most and the least significant bits positions, for all variables in the filter implementation. These formats must allow no overflow and take into account computational errors.

There exist various approaches on determining FxPF, such as Affine Arithmetic [4], or the common approach of numerous simulations [5]. However, these methods are not suitable for our needs since they do not give any mathematical guarantee on the absence of overflow for any possible input and tend to overestimate. We, on the other hand, propose a rigorous algorithm completely based on mathematical proofs.

In Section II we describe the fixed-point arithmetic in question and theoretical background for our approach. In Section III we define the problem for a filter in a state-space representation and introduce a two-step algorithm. Error analysis of the algorithm is presented in Section IV. Finally, numerical results are provided before conclusion.

**Notation:** Throughout the article matrices are in uppercase boldface, vectors are in lowercase boldface, scalars are in lowercase. All matrix inequalities and absolute values are considered element-by-element. The matrix  $\mathbf{1}$  denotes a matrix of ones with the appropriate size.

## II. BASIC BRICKS

### A. The Worst-Case Peak Gain theorem

Let  $\mathcal{H} := (\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$  be a Bounded-Input Bounded-Output stable MIMO LTI filter in state-space representation:

$$\mathcal{H} \begin{cases} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k) &= \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{cases} \quad (1)$$

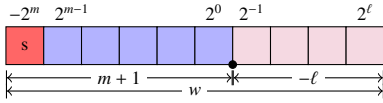


Fig. 2. Fixed-point representation (here,  $m = 5$  and  $\ell = -4$ ).

where  $\mathbf{u}(k) \in \mathbb{R}^{q \times 1}$  is the input vector,  $\mathbf{y}(k) \in \mathbb{R}^{p \times 1}$  the output vector,  $\mathbf{x}(k) \in \mathbb{R}^{n \times 1}$  the state vector and  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times q}$ ,  $\mathbf{C} \in \mathbb{R}^{p \times n}$  and  $\mathbf{D} \in \mathbb{R}^{p \times q}$  are the state-space matrices.

We will use the following approach on deducing the output interval for a filter: suppose all the inputs are guaranteed to be in a known interval:

$$\forall k, \quad |\mathbf{u}_i(k)| \leq \bar{u}_i, \quad i = 1, \dots, q. \quad (2)$$

Then, for all  $k$ , the output  $\mathbf{y}(k)$  lies in the interval

$$|\mathbf{y}(k)| \leq \langle\langle \mathcal{H} \rangle\rangle \bar{\mathbf{u}}, \quad (3)$$

where  $\langle\langle \mathcal{H} \rangle\rangle$  denotes the so-called Worst-Case Peak Gain (WCPG) matrix of the filter [6]. It can be computed as the  $\ell_1$ -norm of the impulse response of the filter.

It can be shown that the WCPG approach gives the smallest interval containing all the possible values of  $\mathbf{y}(k)$  and that the bound (3) can be attained. Moreover, in [7] the authors have published an algorithm to evaluate the WCPG at arbitrary precision.

### B. Fixed-Point arithmetic for filter implementation

The considered filters are implemented using signed FxP Arithmetic, with two's complement representation [8], [9]. Let  $t$  be such a FxP number. It can be written as  $t = -2^m t_m + \sum_{i=\ell}^{m-1} 2^i t_i$ , where its FxPF  $(m, \ell)$  gives the position of the *most* (MSB) significant and *least* (LSB) significant bits respectively, and  $t_i \in \mathbb{B} = \{0, 1\}$  is the  $i^{\text{th}}$  bit of  $t$  as shown in Fig. 2. The wordlength  $w$  is given by  $w = m - \ell + 1$ . The quantization step of the representation is  $2^\ell$ . Further in the article the LSB position  $\ell$  sometimes will be substituted directly with  $m - w + 1$ .

The real number  $t$  is represented in machine by the integer  $T$ , such that  $T = t \cdot 2^{-\ell}$ , where  $T \in [-2^{w-1}; 2^{w-1} - 1] \cap \mathbb{Z}$ . Note that this interval is not symmetric.

Let  $\mathbf{y}(k) \in \mathbb{R}$  be an output of a digital filter. Given the wordlength  $w_y$ , determining the Fixed-Point Formats for  $\mathbf{y}(k)$  means to find a MSB vector  $\mathbf{m}_y$  such that for each  $k$

$$\mathbf{y}_i(k) \in [-2^{m_{y_i}}; 2^{m_{y_i}} - 2^{m_{y_i} - w_{y_i} + 1}]. \quad (4)$$

Obviously, we are interested in the least possible MSB, which guarantees (4): any greater MSB would yield to a larger quantization step, hence round-off noise, at constant wordlength.

## III. DETERMINING THE FIXED-POINT FORMATS

### A. Exact filter

The problem of determining the FxPF for a filter  $\mathcal{H}$  implementation can be formulated as follows. Let  $\mathcal{H}$  be a filter in a state-space representation (1). Suppose all the inputs to be in an interval bounded by  $\bar{\mathbf{u}}$ .

Given the wordlength constraints vector  $\mathbf{w}_x$  for the state and  $\mathbf{w}_y$  for the output variables we look for a FxPF for  $\mathbf{x}(k)$  and

$\mathbf{y}(k)$  such that for any possible input no overflow occurs. We wish to determine the least MSB vectors  $\mathbf{m}_y$  and  $\mathbf{m}_x$  such that

$$\forall k, \quad \mathbf{y}(k) \in [-2^{-m_y}; 2^{m_y} - 2^{m_y - w_{y_i} + 1}], \quad (5)$$

$$\forall k, \quad \mathbf{x}(k) \in [-2^{-m_x}; 2^{m_x} - 2^{m_x - w_{x_i} + 1}]. \quad (6)$$

**Remark 1.** Since the filter  $\mathcal{H}$  is linear and input interval is centered at zero, the output interval is also centered in zero. Therefore, further we will seek to determine the least  $\mathbf{m}_y$  and  $\mathbf{m}_x$  such that

$$\forall k, \quad |\mathbf{y}(k)| \leq 2^{m_y} - 2^{m_y - w_{y_i} + 1}, \quad (7)$$

$$\forall k, \quad |\mathbf{x}(k)| \leq 2^{m_x} - 2^{m_x - w_{x_i} + 1}. \quad (8)$$

### B. Applying the WCPG to compute MSB positions

Applying the WCPG theorem on the filter  $\mathcal{H}$  yields a bound on the output interval:

$$|\mathbf{y}_i(k)| \leq \langle\langle \mathcal{H} \rangle\rangle \bar{u}_i, \quad i = 1, \dots, p. \quad (9)$$

Denote the bound vector  $\bar{\mathbf{y}} := \langle\langle \mathcal{H} \rangle\rangle \bar{\mathbf{u}}$ .

We can determine the FxPF for the output of a LTI filter  $\mathcal{H}$  utilizing the following lemma.

**Lemma 1.** Let  $\mathcal{H} = (\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$  be a BIBO-stable MIMO LTI filter and  $\bar{\mathbf{u}}$  be a bound on the input interval. Suppose the wordlengths  $w_y$  are known and  $w_{y_i} > 1$ ,  $i = 1, \dots, p$ .

If for  $i = 1, \dots, p$  the MSBs are computed with

$$\mathbf{m}_{y_i} = \left\lceil \log_2(\bar{y}_i) - \log_2(1 - 2^{1 - w_{y_i}}) \right\rceil \quad (10)$$

and the LSBs are computed with  $\ell_{y_i} = \mathbf{m}_{y_i} + 1 - w_{y_i}$ , then for all  $k$   $|\mathbf{y}_i(k)| \leq 2^{m_{y_i}} - 2^{m_{y_i} - w_{y_i} + 1}$  and  $\mathbf{m}_{y_i}$  is the least.

*Proof.* We look for the least  $\mathbf{m}_y$  such that (7) holds. Since the bound  $\bar{\mathbf{y}}$  can be reached, it is sufficient to require:

$$\bar{y}_i \leq 2^{m_{y_i}} - 2^{m_{y_i} - w_{y_i} + 1} \quad (11)$$

Solving this inequality for  $\mathbf{m}_{y_i}$  we obtain that the smallest integer, which satisfies the above inequality is

$$\mathbf{m}_{y_i} = \left\lceil \log_2(\bar{y}_i) - \log_2(1 - 2^{1 - w_{y_i}}) \right\rceil. \quad (12)$$

■

### C. Modification of filter $\mathcal{H}$ to determine bounds on the state variable $\mathbf{x}$

Using Lemma 1 we can determine the FxPF for the output of a filter. In order to determine the FxPF for the state variable we modify the filter  $\mathcal{H}$  by vertically concatenating the state vector and the output and include necessary changes into the state matrices.

Denote vector  $\zeta(k) := \begin{pmatrix} \mathbf{x}(k) \\ \mathbf{y}(k) \end{pmatrix}$  to be the new output vector. Then the state-space relationship (1) takes the form:

$$\mathcal{H}_\zeta \begin{cases} \mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \zeta(k) = \begin{pmatrix} \mathbf{I} \\ \mathbf{C} \end{pmatrix} \mathbf{x}(k) + \begin{pmatrix} \mathbf{0} \\ \mathbf{D} \end{pmatrix} \mathbf{u}(k) \end{cases} \quad (13)$$

Hence the problem is to find the least MSB vector  $\mathbf{m}_\zeta$  such that (element-by-element)

$$\forall k, \quad |\zeta_i(k)| \leq 2^{m_{\zeta_i}} - 2^{m_{\zeta_i} - w_{\zeta_i} + 1}. \quad (14)$$

Now, applying the WCPG theorem on the filter  $\mathcal{H}_\zeta$  and using Lemma 1, we can deduce the MSB positions of the state and output vectors for an implementation of the filter  $\mathcal{H}$ .

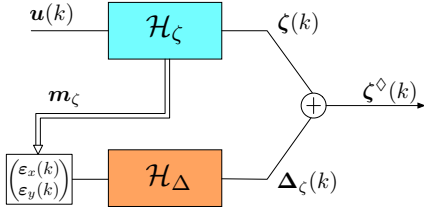


Fig. 3. Implemented filter decomposition.

#### D. Taking rounding errors into account

However, due to the finite-precision degradation what we actually compute is not the exact filter  $\mathcal{H}_\zeta$  but an implemented filter  $\mathcal{H}_\zeta^\diamond$ :

$$\mathcal{H}_\zeta^\diamond \begin{cases} \mathbf{x}^\diamond(k+1) = \diamond_{\ell_x} (\mathbf{A}\mathbf{x}^\diamond(k) + \mathbf{B}\mathbf{u}(k)) \\ \zeta^\diamond(k) = \diamond_{\ell_\zeta} \left( \begin{pmatrix} \mathbf{I} \\ \mathbf{C} \end{pmatrix} \mathbf{x}^\diamond(k) + \begin{pmatrix} \mathbf{0} \\ \mathbf{D} \end{pmatrix} \mathbf{u}(k) \right) \end{cases} \quad (15)$$

where the Sums-of-Products (accumulation of scalar products on the right side) are computed with some rounding operator  $\diamond_\ell$ . Suppose, this operator ensures faithful rounding [10]:

$$|\diamond_\ell(x) - x| < 2^\ell \quad (16)$$

In [11], [12] it was shown that such an operator can be implemented using some extra guard bits for the accumulation.

Denote the errors due to operator  $\diamond_\ell$  as  $\boldsymbol{\varepsilon}_x(k)$  and  $\boldsymbol{\varepsilon}_y(k)$  for the state and output vectors, respectively. Essentially, the vectors  $\boldsymbol{\varepsilon}_x(k)$  and  $\boldsymbol{\varepsilon}_y(k)$  may be associated with the noise which is induced by the filter implementation. Then the implemented filter can be rewritten as

$$\mathcal{H}_\zeta^\diamond \begin{cases} \mathbf{x}^\diamond(k+1) = \mathbf{A}\mathbf{x}^\diamond(k) + \mathbf{B}\mathbf{u}(k) + \boldsymbol{\varepsilon}_x(k) \\ \zeta^\diamond(k) = \begin{pmatrix} \mathbf{I} \\ \mathbf{C} \end{pmatrix} \mathbf{x}^\diamond(k) + \begin{pmatrix} \mathbf{0} \\ \mathbf{D} \end{pmatrix} \mathbf{u}(k) + \begin{pmatrix} \mathbf{0} \\ \mathbf{I} \end{pmatrix} \boldsymbol{\varepsilon}_y(k) \end{cases}, \quad (17)$$

where

$$|\boldsymbol{\varepsilon}_x(k)| < 2^{\ell_x}, \quad |\boldsymbol{\varepsilon}_y(k)| < 2^{\ell_y}.$$

It should be remarked that since the operator  $\diamond_\ell$  is applied  $\boldsymbol{\varepsilon}_x(k) \neq \mathbf{x}(k) - \mathbf{x}^\diamond(k)$  and  $\boldsymbol{\varepsilon}_y(k) \neq \mathbf{y}(k) - \mathbf{y}^\diamond(k)$ . As the rounding also affects the filter state, the  $\mathbf{x}^\diamond(k)$  drifts away from  $\mathbf{x}(k)$  over time, whereas with  $\boldsymbol{\varepsilon}_x(k)$  we consider the error due to one step only.

It can be observed that at each instance of time the state and output vectors are computed out of  $\mathbf{u}(k)$  and error-vectors, which can be considered as inputs as well. Thanks to the linearity of the filters, we can decompose the actually implemented filter into a sum of the exact filter and an "error-filter"  $\mathcal{H}_\Delta$  as shown in Fig. 3. Note that this "error-filter" is an artificial one; it is not required to be "implemented" by itself and serves exclusively for error-analysis purposes.

The filter  $\mathcal{H}_\Delta$  is obtained by computing the difference between  $\mathcal{H}_\zeta^\diamond$  and  $\mathcal{H}_\zeta$ . This filter takes the rounding errors  $\boldsymbol{\varepsilon}(k) := \begin{pmatrix} \boldsymbol{\varepsilon}_x(k) \\ \boldsymbol{\varepsilon}_y(k) \end{pmatrix}$  as input and returns the result of their propagation through the filter:

$$\mathcal{H}_\Delta \begin{cases} \Delta_x(k+1) = \mathbf{A}\Delta_x(k) + \begin{pmatrix} \mathbf{I} \\ \mathbf{0} \end{pmatrix} \boldsymbol{\varepsilon}(k) \\ \Delta_\zeta(k) = \begin{pmatrix} \mathbf{I} \\ \mathbf{C} \end{pmatrix} \Delta_x(k) + \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \boldsymbol{\varepsilon}(k) \end{cases}, \quad (18)$$

where, the vector  $\boldsymbol{\varepsilon}(k)$  is guaranteed to be in the interval bounded by  $\bar{\boldsymbol{\varepsilon}} := 2^{\ell_\varepsilon}$ .

Once the decomposition is done, we can apply the WCPG theorem on the "error-filter"  $\mathcal{H}_\Delta$  and deduce the output interval of the computational errors propagated through filter:

$$\forall k, \quad |\Delta_\zeta(k)| \leq \langle \langle \mathcal{H}_\Delta \rangle \rangle \cdot \bar{\boldsymbol{\varepsilon}}. \quad (19)$$

Hence, the output of the implemented filter is bounded with

$$|\zeta^\diamond(k)| \leq |\zeta(k)| + |\Delta_\zeta(k)|. \quad (20)$$

Applying Lemma 1 on the implemented filter and using (21) we obtain that the MSB vector  $\mathbf{m}_\zeta^\diamond$  can be upper bounded by

$$\mathbf{m}_{\zeta_i}^\diamond = \left[ \log_2 \left( \left( \langle \langle \mathcal{H}_\zeta \rangle \rangle \cdot \bar{\mathbf{u}} \right)_i + \left( \langle \langle \mathcal{H}_\Delta \rangle \rangle \cdot \bar{\boldsymbol{\varepsilon}} \right)_i \right) - \log_2 \left( 1 - 2^{1-w_{\zeta_i}} \right) \right]. \quad (21)$$

Therefore, the FxPF  $(\mathbf{m}_\zeta^\diamond, \ell_\zeta^\diamond)$  guarantee that no overflows occur for the implemented filter.

Since the input of the error filter  $\mathcal{H}_\Delta$  depends on the FxPF chosen for implementation, we cannot directly use (22). The idea is to first compute the FxPF of the variables in the exact filter  $\mathcal{H}$ , where computational errors are not taken into account, and use it as an initial guess for implemented filter  $\mathcal{H}_\zeta^\diamond$ . Hence, we obtain the following two-step algorithm:

- Step 1: Determine the FxPF  $(\mathbf{m}_\zeta, \ell_\zeta)$  for the exact filter  $\mathcal{H}_\zeta$
- Step 2: Construct the "error-filter"  $\mathcal{H}_\Delta$ , which gives the propagation of the computational errors induced by format  $(\mathbf{m}_\zeta, \ell_\zeta)$ ; then, compute the FxPF  $(\mathbf{m}_\zeta^\diamond, \ell_\zeta^\diamond)$  of the actually implemented filter  $\mathcal{H}_\zeta^\diamond$  using (22).

The above algorithm takes into account the filter implementation errors. However, the algorithm itself is implemented in finite-precision and can suffer from rounding errors, which influence the output result. All operations in the MSB computation will induce errors, so what we actually compute are only floating-point approximations  $\widehat{\mathbf{m}}_\zeta$  and  $\widehat{\mathbf{m}}_\zeta^\diamond$ . In what follows, we propose error-analysis of the floating-point evaluation of the MSB positions via (10) and (22).

#### IV. ERROR ANALYSIS OF THE MSB COMPUTATION FORMULA

Let us consider case of  $\widehat{\mathbf{m}}_{\zeta_i}^\diamond$  and show afterwards that  $\widehat{\mathbf{m}}_{\zeta_i}^\diamond$  is its special case. To reduce the size of expressions, denote

$$\mathfrak{m} := \log_2 \left( \left( \langle \langle \mathcal{H}_\zeta \rangle \rangle \cdot \bar{\mathbf{u}} \right)_i + \left( \langle \langle \mathcal{H}_\Delta \rangle \rangle \cdot \bar{\boldsymbol{\varepsilon}} \right)_i \right) - \log_2 \left( 1 - 2^{1-w_{\zeta_i}} \right).$$

Handling floating-point analysis of multiplications and additions in (22) is trivial using approach by Higham [13]. The difficulty comes from the WCPG matrices, which cannot be computed exactly. However both approximations  $\widehat{\langle \langle \mathcal{H}_\zeta \rangle \rangle}$  and  $\widehat{\langle \langle \mathcal{H}_\Delta \rangle \rangle}$ , even if computed with arbitrary precision, bear some errors  $\boldsymbol{\varepsilon}_{\text{WCPG}_\zeta}$  and  $\boldsymbol{\varepsilon}_{\text{WCPG}_\Delta}$  that satisfy

$$0 \leq \widehat{\langle \langle \mathcal{H}_\Delta \rangle \rangle} - \langle \langle \mathcal{H}_\Delta \rangle \rangle \leq \boldsymbol{\varepsilon}_{\text{WCPG}_\Delta} \cdot \mathbf{1} \quad (22)$$

$$0 \leq \widehat{\langle \langle \mathcal{H}_\zeta \rangle \rangle} - \langle \langle \mathcal{H}_\zeta \rangle \rangle \leq \boldsymbol{\varepsilon}_{\text{WCPG}_\zeta} \cdot \mathbf{1} \quad (23)$$

Introducing the errors on the WCPG computations into the formula (22) we obtain that what we actually compute is

$$\widehat{\mathbf{m}}_{\zeta_i}^\diamond \leq \left\lceil \mathfrak{m} + \log_2 \left( 1 + \frac{\boldsymbol{\varepsilon}_{\text{WCPG}_\zeta} \sum_{j=1}^q \bar{\mathbf{u}}_j + \boldsymbol{\varepsilon}_{\text{WCPG}_\Delta} \sum_{j=1}^{n+p} \bar{\boldsymbol{\varepsilon}}_j}{\left( \langle \langle \mathcal{H}_\zeta \rangle \rangle \bar{\mathbf{u}} \right)_i + \left( \langle \langle \mathcal{H}_\Delta \rangle \rangle \bar{\boldsymbol{\varepsilon}} \right)_i} \right) \right\rceil. \quad (24)$$

The error term in (25) cannot be zero (apart from trivial case with zero  $\bar{\mathbf{u}}$ ). However, since we can control the accuracy of the WCPG matrices, we can deduce conditions for the approximation  $\widehat{\mathbf{m}}_{\zeta_i}^\diamond$  to be off by at most one.

**Lemma 2.** *If the WCPG matrices  $\langle\langle\mathcal{H}_\zeta\rangle\rangle$  and  $\langle\langle\mathcal{H}_\Delta\rangle\rangle$  are computed such that (23) and (24) hold with*

$$\varepsilon_{\text{WCPG}_\Delta} < \frac{1}{2} \frac{\langle\langle\mathcal{H}_\Delta\rangle\rangle \cdot \bar{\mathbf{e}}}{\sum_{j=1}^{p+n} \bar{\mathbf{e}}_j} \quad (25)$$

$$\varepsilon_{\text{WCPG}_\zeta} < \frac{1}{2} \frac{\langle\langle\mathcal{H}_\zeta\rangle\rangle \cdot \bar{\mathbf{u}}}{\sum_{j=1}^q \bar{\mathbf{u}}_j}, \quad (26)$$

then

$$0 \leq \widehat{\mathbf{m}}_{\zeta_i}^\diamond - \mathbf{m}_{\zeta_i}^\diamond \leq 1. \quad (27)$$

*Proof.* Proof by construction, we reason as follows: since the error-term caused by the WCPG floating-point evaluation is positive and the ceil function is increasing, then

$$\widehat{\mathbf{m}}_{\zeta_i}^\diamond - \mathbf{m}_{\zeta_i}^\diamond \geq 0, \quad (28)$$

i.e. the floating-point approximation  $\widehat{\mathbf{m}}_{\zeta_i}^\diamond$  is guaranteed to never be underestimated. However, it can overestimate the MSB position by

$$\widehat{\mathbf{m}}_{\zeta_i}^\diamond - \mathbf{m}_{\zeta_i}^\diamond \leq \left\lceil \mathbf{m} - \left\lfloor \mathbf{m} \right\rfloor \right\rceil_{-1 < \leq 0} + \log_2 \left( 1 + \frac{\varepsilon_{\text{WCPG}_\zeta} \sum_{j=1}^q \bar{\mathbf{u}}_j + \varepsilon_{\text{WCPG}_\Delta} \sum_{j=1}^{n+p} \bar{\mathbf{e}}_j}{\langle\langle\mathcal{H}_\zeta\rangle\rangle \bar{\mathbf{u}} + \langle\langle\mathcal{H}_\Delta\rangle\rangle \bar{\mathbf{e}}} \right). \quad (29)$$

The approximation  $\widehat{\mathbf{m}}_{\zeta_i}^\diamond$  overestimates at most by one bit if and only if the error term is contained in the interval  $[0, 1)$ , i.e. if

$$0 \leq \log_2 \left( 1 + \frac{\varepsilon_{\text{WCPG}_\zeta} \sum_{j=1}^q \bar{\mathbf{u}}_j + \varepsilon_{\text{WCPG}_\Delta} \sum_{j=1}^{n+p} \bar{\mathbf{e}}_j}{\langle\langle\mathcal{H}_\zeta\rangle\rangle \bar{\mathbf{u}} + \langle\langle\mathcal{H}_\Delta\rangle\rangle \bar{\mathbf{e}}} \right) < 1. \quad (30)$$

Hence, using the above condition we can deduce the upper bounds on the  $\varepsilon_{\text{WCPG}_\zeta}$  and  $\varepsilon_{\text{WCPG}_\Delta}$ :

$$0 \leq \frac{\varepsilon_{\text{WCPG}_\zeta} \sum_{j=1}^q \bar{\mathbf{u}}_j + \varepsilon_{\text{WCPG}_\Delta} \sum_{j=1}^{n+p} \bar{\mathbf{e}}_j}{\langle\langle\mathcal{H}_\zeta\rangle\rangle \bar{\mathbf{u}} + \langle\langle\mathcal{H}_\Delta\rangle\rangle \bar{\mathbf{e}}} < 1. \quad (31)$$

Since all the terms are positive, the left inequality is always true. The right inequality in (32) is satisfied for instance if

$$\frac{\varepsilon_{\text{WCPG}_\zeta} \sum_{j=1}^q \bar{\mathbf{u}}_j}{\langle\langle\mathcal{H}_\zeta\rangle\rangle \cdot \bar{\mathbf{u}}_i} < \frac{1}{2} \quad \frac{\varepsilon_{\text{WCPG}_\Delta} \sum_{j=1}^{n+p} \bar{\mathbf{e}}_j}{\langle\langle\mathcal{H}_\Delta\rangle\rangle \cdot \bar{\mathbf{e}}_i} < \frac{1}{2}. \quad (32)$$

Rearranging terms we obtain following inequalities on the WCPG computation with error:

$$\varepsilon_{\text{WCPG}_\zeta} < \frac{1}{2} \cdot \frac{\langle\langle\mathcal{H}_\zeta\rangle\rangle \bar{\mathbf{u}}_i}{\sum_{j=1}^q \bar{\mathbf{u}}_j} \quad \varepsilon_{\text{WCPG}_\Delta} < \frac{1}{2} \cdot \frac{\langle\langle\mathcal{H}_\Delta\rangle\rangle \bar{\mathbf{e}}_i}{\sum_{j=1}^{n+p} \bar{\mathbf{e}}_j}. \quad (33)$$

Unfortunately, the above results cannot be used in practice, since they depend themselves on the exact WCPG matrices.

Denote  $\langle\langle\mathcal{H}\rangle\rangle := |\mathbf{D}| + |\mathbf{CB}| + |\mathbf{CAB}|$ . It can be shown that it is a lower bound of the WCPG matrix. We can compute this matrix exactly. Obviously,

$$\frac{\langle\langle\mathcal{H}_\Delta\rangle\rangle \cdot \bar{\mathbf{e}}}{\sum_{j=1}^{p+n} \bar{\mathbf{e}}_j} \leq \frac{\langle\langle\mathcal{H}_\Delta\rangle\rangle \cdot \bar{\mathbf{e}}_i}{\sum_{j=1}^{p+n} \bar{\mathbf{e}}_j} \quad (34)$$

$$\frac{\langle\langle\mathcal{H}_\zeta\rangle\rangle \cdot \bar{\mathbf{u}}}{\sum_{j=1}^q \bar{\mathbf{u}}_j} \leq \frac{\langle\langle\mathcal{H}_\zeta\rangle\rangle \cdot \bar{\mathbf{u}}_i}{\sum_{j=1}^q \bar{\mathbf{u}}_j}. \quad (35)$$

Hence, if the WCPG matrices in the right sides of (34) are substituted with their lower bounds, the condition (32) stays satisfied and we obtain bounds (26) and (27).  $\blacksquare$

Analogously, Lemma 2 can be applied to the computation of  $\widehat{\mathbf{m}}_{\zeta_i}$  with the terms concerning filter  $\mathcal{H}_\Delta$  set to zero.

In order to guarantee the output MSB position to be optimal, an instance of the Table Maker's Dilemma [10] must be solved. This is due to the simplification by using the triangle inequality in (21), which basically requires both filters  $\mathcal{H}_\zeta$  and  $\mathcal{H}_\Delta$  reach the worst-case bound with the same input sequence.

## V. COMPLETE ALGORITHM

The two-step algorithm, presented in subsection III-D takes into account accumulation of computational errors in a filter over time and Lemma 2 presents error-analysis of the MSB position computation procedure. However, one additional fact has not been taken into account.

In most cases the MSB vectors  $\widehat{\mathbf{m}}_\zeta$  (computed on Step 1) and  $\widehat{\mathbf{m}}_\zeta^\diamond$  (computed on Step 2) are the same. However, in some cases they are not, which can happen due to one of the following reasons:

- the accumulated rounding error due to the FxPF ( $\widehat{\mathbf{m}}_\zeta, \widehat{\ell}_\zeta$ ) makes output of the actually implemented filter pass over to the next binade; or
- the floating-point approximation  $\widehat{\mathbf{m}}_\zeta^\diamond$  is off by one.

Moreover, if the MSB position after Step 2 of the algorithm is increased, the LSB position moves along and increases the error. Therefore, the modified format must be re-checked. Hence, the FxPF determination algorithm gets transformed into the following three-step procedure:

- Step 1: Determine the FxPF ( $\widehat{\mathbf{m}}_\zeta, \widehat{\ell}_\zeta$ ) for the exact filter  $\mathcal{H}_\zeta$ ;
- Step 2: Construct the "error-filter"  $\mathcal{H}_\Delta$ , which shows the propagation of the computational errors induced by format ( $\widehat{\mathbf{m}}_\zeta, \widehat{\ell}_\zeta$ ); then, compute the FxPF ( $\widehat{\mathbf{m}}_\zeta^\diamond, \widehat{\ell}_\zeta^\diamond$ ) of the actually implemented filter  $\mathcal{H}_\zeta^\diamond$ ;
- Step 3: If  $\widehat{\mathbf{m}}_\zeta^\diamond == \widehat{\mathbf{m}}_\zeta$ , then return ( $\widehat{\mathbf{m}}_\zeta^\diamond, \widehat{\ell}_\zeta^\diamond$ ); otherwise  $\widehat{\mathbf{m}}_{\zeta_i} \leftarrow \widehat{\mathbf{m}}_{\zeta_i} + 1$  and go to Step 2.

## VI. NUMERICAL RESULTS

The above described algorithm was implemented as a C library, using GNU MPFR<sup>1</sup> [14] version 3.1.12, GNU MPFI<sup>2</sup> version 1.5.1 and the WCPG library [7].

<sup>1</sup><http://www.mpfr.org/>

<sup>2</sup><https://gforge.inria.fr/projects/mpfi/>

Consider following example: let  $\mathcal{H}$  be a random stable filter with 3 states, 1 input and 1 output. Suppose the inputs are in an interval absolutely bounded by  $\bar{u} = 5.125$ . In this example we set all the wordlength constraints to 7 bits, however our library supports the multiple wordlength paradigm.

The results of the work of our algorithm can be observed in Table I. On the Step 1 we deduce the MSB positions for the filter  $\mathcal{H}$ , which does not take computational errors into account. These MSBs serve as an initial guess for the Step 2. We compute the vector  $\bar{\epsilon}_z$  and construct the "error-filter"  $\mathcal{H}_\Delta$ . Taking into account the error propagation yields changes in the MSB positions: the rounding errors force the third state  $x_3(k)$  to pass over to the next binade. However, moving the MSB yields larger quantization step and an addition step is required. Step 3 verifies that the FxPF deduced on the Step 2 guarantees no overflow.

We can verify the result by tracing the state and output vectors in two cases: the exact vectors  $\mathbf{y}(k)$  and  $\mathbf{x}(k)$ ; and the quantized vectors  $\mathbf{y}^\diamond(k)$  and  $\mathbf{x}^\diamond(k)$ . The quantized vectors are the result of an implementation of the filter  $\mathcal{H}$  with FxPF deduced on the Step 1. For each state  $x_i(k)$  and output  $y(k)$  we take an input sequence which makes the state respectively output attain the theoretical bound. However, this sequence is not guaranteed to maximize the quantized error.

In Fig. 4 we can observe that the exact vector  $\mathbf{y}(k)$  attains the bound computed with the WCPG theorem. The quantized output stays within the bound  $\bar{y}^\diamond$  for the implemented filter but passes over the bound of the exact filter. However, as the bound  $\bar{y}$  is far from the next binade, the MSB position is not changed and the computational errors do not result in overflow.

For the state  $x_3(k)$  the bound  $\bar{x}_3$  is very close to the  $2^{m_{x_3}} - 2^{\ell_{x_3}}$  and taking into account the filter computational errors yields passing to the next binade. It can be clearly observed in Fig. 5 that the quantized state passes over the value  $2^{m_{x_3}} - 2^{\ell_{x_3}}$ . If the computational errors were not taken into account, the initial format (deduced on the Step 1) would result in overflow.

Therefore, we observed that deducing the FxPF without taking into account the computational errors can result in overflow but our approach guarantees a reliable implementation without overestimation.

	states			output
	$x_1(k)$	$x_2(k)$	$x_3(k)$	$y(k)$
<b>Step 1</b>	6	7	5	6
<b>Step 2</b>	6	7	6	6
<b>Step 3</b>	6	7	6	6

TABLE I

EVOLUTION OF THE MSB POSITIONS VECTOR THROUGH ALGORITHM ITERATIONS

## VII. CONCLUSIONS

We give an algorithm to determine the FxPF for all variables of a LTI filter, where the computation errors are taken into account. We ensure, by construction, that no overflow occurs. However, the computed MSB positions can overestimate at most by one bit. In order to guarantee no overestimation at all an instance of the TMD must be solved. Our algorithm can be extended to any LTI filter realization (Direct Forms, Lattice filters, etc.) using the Specialized Implicit Framework (SIF) tool

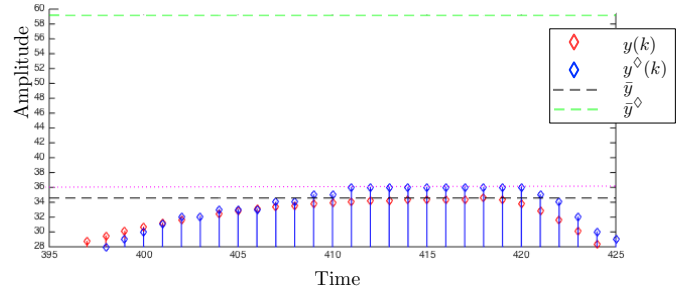


Fig. 4. The exact and quantized outputs of the example. Quantized output does not pass over to the next binade.

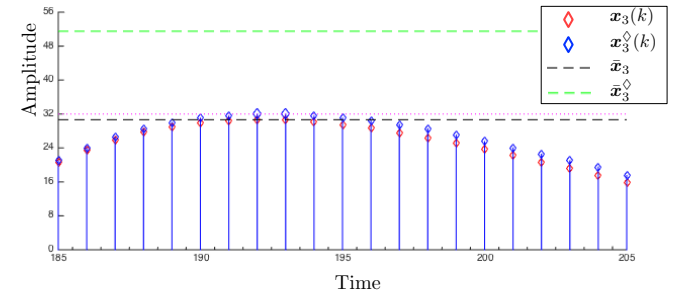


Fig. 5. The exact and quantized third state of the example. Quantized output passes over to the next binade.

[6]. SIF enables unification of various LTI filter realizations for analysis, comparison and FxP code generation. Therefore, this work is an essential basic brick for the automatic filter-to-code conversion.

## REFERENCES

- [1] A. V. Oppenheim and R. W. Schaffer, *Digital Signal Processing*. Englewood Cliffs, NJ: Prentice-Hall, 1975.
- [2] B. Friedland, *Control Systems Design: An Introduction to State-Space Methods*. McGraw-Hill Higher Education, 1985.
- [3] A. Fettweis, "Wave digital filters: Theory and practice," *Proc. of the IEEE*, vol. 74, no. 2, 1986.
- [4] J. Lopez, C. Carreras, and O. Nieto-Taladriz, "Improved interval-based characterization of fixed-point LTI systems with feedback loops," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, vol. 26, no. 11, pp. 1923–1933, 2007.
- [5] 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*, vol. 45, pp. 1455–1464, November 1998.
- [6] T. Hilaire and B. Lopez, "Reliable implementation of linear filters with fixed-point arithmetic," in *Proc. SiPS*, 2013.
- [7] A. Volkova, T. Hilaire, and C. Lauter, "Reliable evaluation of the worst-case peak gain matrix in multiple precision," in *Proc. 22nd IEEE Symposium on Computer Arithmetic (ARITH22)*, 2015.
- [8] J. v. Neumann, "First draft of a report on the edvac," tech. rep., 1945.
- [9] T. Finley, "Two's complement." Cornell University lecture notes, 2000.
- [10] J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres, *Handbook of Floating-Point Arithmetic*. Birkhäuser Boston, 2010.
- [11] B. Lopez, T. Hilaire, and L. S. Didier, "Formatting bits to better implement signal processing algorithms," in *Proc. PECCS, Portugal*, pp. 104–111, 2014.
- [12] F. de Dinechin, M. Istoan, and A. Massouri, "Sum-of-product architectures computing just right," in *Application-Specific Systems, Architectures and Processors (ASAP)*, IEEE, 2014.
- [13] N. J. Higham, *Accuracy and stability of numerical algorithms (2. ed.)*. SIAM, 2002.
- [14] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier, and P. Zimmermann, "MPFR: A multiple-precision binary floating-point library with correct rounding," *ACM Transactions on Mathematical Software*, vol. 33, no. 2, pp. 13:1–13:15, 2007.