



# Reliable fixed-point implementation of linear data-flows

Thibault Hilaire, Anastasia Volkova, Maminionja Ravoson

► **To cite this version:**

Thibault Hilaire, Anastasia Volkova, Maminionja Ravoson. Reliable fixed-point implementation of linear data-flows. IEEE International Workshop on Signal Processing Systems (SiPS), Oct 2016, Dallas, United States. <hal-01347637v2>

**HAL Id: hal-01347637**

**<http://hal.upmc.fr/hal-01347637v2>**

Submitted on 7 Sep 2016

**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.

# Reliable fixed-point implementation of linear data-flows

Thibault Hilaire, Anastasia Volkova, Maminionja Ravoson  
Sorbonne Universités, UPMC Univ Paris 06, UMR 7606, LIP6, F-75005, Paris, France  
Email: first\_name.last\_name@lip6.fr

**Abstract**—In this article, we propose a complete methodology to implement a signal processing or control-engineering algorithm described with a linear data-flow into numerical code using fixed-point arithmetic. Our approach is based on a reliable determination of the Worst-Case Peak gain of a filter as well as on rigorous error analysis of roundoff error propagation. It guarantees that no overflow will occur and that the output error due to the finite precision implementation is less than a given bound.

Without loss of generality, we consider the linear data-flows given in the form of Simulink block diagram. It is first transposed into an internal matrix-based representation and then the reliable evaluation of the magnitudes of each internal variable is performed. Our approach allows determining the minimum word-length required to achieve a given accuracy. Finally, the methodology is illustrated with numerical examples.

## I. INTRODUCTION

The implementation of embedded algorithms is a tedious task for many reasons. Due to the limited resources available (power, memory, computing units, processing time, etc.), the transformation from a mathematical model (whatever its form: equations, transfer function, data-flow graph, etc.) into embedded code is not straightforward. One principal pitfall concerns the accuracy part: generally finite precision arithmetic, like the Fixed-Point (FxP) arithmetic is used, and the computational errors induced by limited bit-width may be significant.

Moreover, in FxP arithmetic, the developer has to deal with the binary-point alignment, fixed-point conversion, etc., while in floating-point arithmetic this is done by the hardware.

The two main bottlenecks for the Fixed-Point implementation are the setting of the FxP formats (they depend on the magnitude of each variable) and the roundoff error analysis. We are interested in determining the impact of roundoff errors on the output(s) of the implemented system, therefore we need to rigorously analyse the difference between the exact system and the implemented one. These two issues may be not so easy to deal reliably with, especially in the cases where the system has an internal feedback, like for the Infinite Impulse Response (IIR) filters.

In practice, the FxP conversion is done using simulations, typically using Matlab/Simulink tools<sup>1</sup>:

- a) run the system with some inputs (using floating-point double precision arithmetic) and determine the magnitudes of all variables;

This work has been partially sponsored by french ANR agency under grant No ANR-13-INSE-0007-04 MetaLibm.

<sup>1</sup><http://www.mathworks.com>

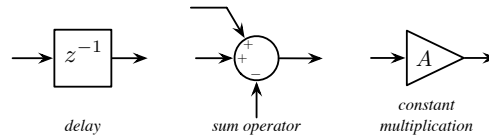


Fig. 1. Blocks of linear data-flows.

- b) deduce the FxP formats and the FxP operations;
- c) run the FxP implementation of the system and compare it with the simulations from step a);
- d) if the result is not convincing, increase the word-length, and return to step b).

This process is unsatisfactory for various reasons. First, a finite-precision (floating-point) simulation is taken as a reference for comparison which can be inaccurate. Secondly, the input sequences to be tested are often not representative. Moreover, these non-exhaustive and time-consuming simulations do not give any **reliable** guarantee on the FxP implementation. The saturation technique is widely used to deal with overflows, however it does not ensure the correct value of the output.

In our work we consider only **linear** algorithms, *i.e.* data-flow systems built only with the sum operators (+ or – operations), constant multiplications and delays (see Fig. 1). The examples of such algorithms are the IIR and FIR (Finite Impulse Response) filters, Linear Time Invariant (LTI) controllers, etc.

In order to overcome the drawbacks of the simulation-based techniques and provide mathematical guarantees on the FxP implementation, we use an analytical data-flow representation called Specialized Implicit Framework (SIF) [1].

We propose a technique on representation of the linear data-flows with SIF formalism. Then, we can easily determine the inputs-outputs relationship as transfer functions, and then establish analytically the impact of the finite precision implementation. The implementation flow is encapsulated into an automatized code generator. A rigorous approach on determining the magnitudes and the Fixed-Point formats has been developed for the generator [2], [3]. Furthermore, without using any simulations we can optimize the word-lengths within the algorithm evaluation to satisfy the output precision constraints.

The paper is organized as follows. Our complete flow from data-flow system to fixed-point code is described in Section II with some necessary prerequisites. Section III focuses on the transformation of Simulink linear data-flow systems to our

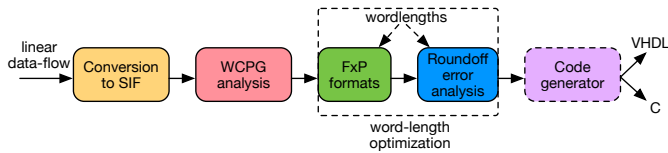


Fig. 2. Data-flow implementation steps.

Specialized Implicit Framework. The fixed-point conversion and the complete roundoff analysis are done in Section IV. Finally, some illustrative examples of fixed-point implementation under precision constraints are given in Section V.

**Notation:** throughout this paper, real numbers are in lowercase, column vectors in lowercase boldface and matrices in uppercase boldface;  $\mathbf{I}_n$  is a  $n \times n$  identity matrix.

## II. METHODOLOGY AND PREREQUISITES

### A. Methodology

In order to provide rigorous error analysis of finite-precision implementation effects, we propose to use an analytical representation of a linear system, SIF, which is described in details in Section II-B. This is a matrix representation, which can be obtained for any data-flow graph. It preserves all the properties of the system but is easier to operate with, in our context.

Our approach consists of following. First, we represent the data-flow with SIF formalism. Then, we deduce the magnitudes of all variables within the system. In comparison to the commonly used approach, we do not use any simulations to compute them. Considering that the inputs are bounded, our approach is based on reliable evaluation of the maximum amplification of the output of a system using the Worst-Case Peak Gain theorem [2]. This gives us a theoretical guarantee on the computed magnitudes and does not depend on a specific choice of inputs. See Section II-D for detailed description.

Given an output error bound and word-length constraints for all variables we choose Fxp formats for the implementation (using **reliably** computed magnitudes) such that the output code ensures that no overflow occurs.

Finally, we deduce the parameters for algorithm evaluation in Fxp such that the computed output satisfies an *a priori* given error bound for each iteration of the system. See Section IV-B.

However, the propagation of the roundoff errors through recursive systems is non-linear and can have a significant impact on the output. We take into account the roundoff errors propagation and provide a rigorous bound on the difference between the implemented and the ideal systems.

Therefore, our approach can be represented as following steps, as shown in Fig. 2:

- a) transform a data-flow to the SIF formalism;
- b) **rigorously** compute magnitudes using the Worst-Case Peak-Gain analysis;
- c) given word-length constraints determine the Fxp formats based on rigorous magnitudes, while taking into account the propagation of the computational errors;
- d) generate the software or hardware code.

Moreover, if the word-lengths are not fixed, the Fixed-Point Formats computed on step c) can be optimized to meet an output error-bound criteria without using any simulations.

### B. Specialized Implicit Framework

The Specialized Implicit Framework has been first proposed in [1] as a unifying tool in order to describe and encompass all the possible realizations of a given transfer function (like the direct forms, state-spaces, cascade or parallel decompositions, lattice filters, etc.). It allows the study and comparison of their finite precision effects. SIF is an extension of the state-space realization, modified in order to allow chained Sum-of-Products (SoP) operations.

In fact, all linear data-flows (those with delays, multiplications by constants and additions) can be represented with the SIF. This macroscopic description is more suited for the analysis than a graph relationship as it gives direct analytical formula for the finite precision error analysis [1]. Both the Single Input Single Output (SISO) and the Multiple-Inputs Multiple Outputs (MIMO) can be described with SIF.

Denote  $\mathbf{u}(k)$  and  $\mathbf{y}(k)$  the vector of  $q$  inputs and the vector of  $p$  outputs respectively. The  $n$  variables that are stored from one step to the other are in the state vector  $\mathbf{x}(k)$ , while the  $l$  intermediate results are collected in the vector  $\mathbf{t}(k)$ . Then, the SIF is the following set of equations:

$$\begin{pmatrix} \mathbf{J} & \mathbf{0} & \mathbf{0} \\ -\mathbf{K} & \mathbf{I}_n & \mathbf{0} \\ -\mathbf{L} & \mathbf{0} & \mathbf{I}_p \end{pmatrix} \begin{pmatrix} \mathbf{t}(k+1) \\ \mathbf{x}(k+1) \\ \mathbf{y}(k) \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \mathbf{M} & \mathbf{N} \\ \mathbf{0} & \mathbf{P} & \mathbf{Q} \\ \mathbf{0} & \mathbf{R} & \mathbf{S} \end{pmatrix} \begin{pmatrix} \mathbf{t}(k) \\ \mathbf{x}(k) \\ \mathbf{u}(k) \end{pmatrix} \quad (1)$$

The vector  $\mathbf{t}(k)$  is not used for computations at step  $k$ , which characterizes the concept of intermediate variables.

The matrix  $\mathbf{J}$  is lower-triangular with 1 on its diagonal, so the first value of  $\mathbf{t}(k+1)$  is first computed, then the second one is computed using the first and so on (thus, the computation of  $\mathbf{J}^{-1}$  is not necessary). The *implicit* term  $\mathbf{J}\mathbf{t}(k+1)$  naturally serves for describing the specific order of the computation. Therefore, the general algorithm to compute the SIF is:

$$\begin{aligned} \mathbf{J}\mathbf{t}(k+1) &\leftarrow \mathbf{M}\mathbf{x}(k) + \mathbf{N}\mathbf{u}(k) \\ \mathbf{x}(k+1) &\leftarrow \mathbf{K}\mathbf{t}(k+1) + \mathbf{P}\mathbf{x}(k) + \mathbf{Q}\mathbf{u}(k) \\ \mathbf{y}(k) &\leftarrow \mathbf{L}\mathbf{t}(k+1) + \mathbf{R}\mathbf{x}(k) + \mathbf{S}\mathbf{u}(k) \end{aligned} \quad (2)$$

The matrices  $\mathbf{J}$ ,  $\mathbf{K}$  and  $\mathbf{L}$  allow us to describe the sequence of computations. For example,  $\mathbf{y} = \mathbf{M}_2\mathbf{M}_1\mathbf{x}$  can be computed as  $\mathbf{y} = (\mathbf{M}_2\mathbf{M}_1)\mathbf{x}$  or as  $\mathbf{y} = \mathbf{M}_2(\mathbf{M}_1\mathbf{x})$ . The latter expression will be described as

$$\begin{pmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{M}_2 & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} \mathbf{M}_1 \\ \mathbf{0} \end{pmatrix} \mathbf{x}, \quad (3)$$

with  $\mathbf{t}$  holding the intermediate value  $\mathbf{M}_1\mathbf{x}$ .

The same approach will be used to transform any linear data-flow into a SIF in Section III.

Throughout the paper we consider that the computations associated to the equations (1) are ordered from top to bottom<sup>2</sup>, associated in a one to one manner to the Algorithm (2).

<sup>2</sup>Of course, there is no need to consider the null coefficients. The associated computations may be removed.

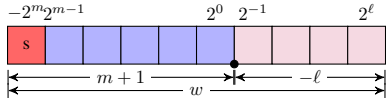


Fig. 3. Representation of a fixed-point number in format  $(m, \ell)$ .

It can be easily established that equation (1) is equivalent in infinite precision to the state-space filter  $(\mathbf{A}_Z, \mathbf{B}_Z, \mathbf{C}_Z, \mathbf{D}_Z)$ :

$$\begin{cases} \mathbf{x}(k+1) &= \mathbf{A}_Z \mathbf{x}(k) + \mathbf{B}_Z \mathbf{u}(k) \\ \mathbf{y}(k) &= \mathbf{C}_Z \mathbf{x}(k) + \mathbf{D}_Z \mathbf{u}(k) \end{cases} \quad (4)$$

with:

$$\begin{aligned} \mathbf{A}_Z &= \mathbf{K}\mathbf{J}^{-1}\mathbf{M} + \mathbf{P}, & \mathbf{B}_Z &= \mathbf{K}\mathbf{J}^{-1}\mathbf{N} + \mathbf{Q}, \\ \mathbf{C}_Z &= \mathbf{L}\mathbf{J}^{-1}\mathbf{M} + \mathbf{R}, & \mathbf{D}_Z &= \mathbf{L}\mathbf{J}^{-1}\mathbf{N} + \mathbf{S}. \end{aligned} \quad (5)$$

However, equation (4) corresponds to a different set of coefficients than the one in (1). Therefore, while in infinite precision (1) and (4) are equivalent, in finite precision they have different numerical properties. Thus, the SIF not only describes the input/output relationship but also fully captures properties of computational algorithm.

The main interest of such a framework is that it allows expressing various algorithms of the same filter in a unifying form, in order to study and compare the effect of the finite precision arithmetic on them, and then to choose the best FxP algorithm to implement this filter. All the classical measures used to evaluate the impact of the quantization of the coefficients (sensitivity-based measure like in [4], [5], [6]) and the impact of the roundoff errors ([5], [7], [8], [9], [10]) have been extended to the SIF [11], [12].

### C. Fixed-Point Arithmetic

In this paper, only the signed fixed-point arithmetic with 2's complement representation is used. Let  $z$  be a  $w$ -bit fixed-point number:

$$z = -2^m z_m + \sum_{i=\ell}^{m-1} 2^i z_i \quad (6)$$

where  $m$  and  $\ell$  are the positions of the most significant bit (MSB) and the least significant bit (LSB) of  $z$ , respectively. Let the couple  $(m, \ell)$  denote the Fixed-Point format of  $z$  (see Fig. 3). The word-length  $w$  can be obtained with:

$$w = m - \ell + 1, \quad (7)$$

In order to convert a non-zero real number  $r$  into a fixed-point number  $\tilde{r}$ , the first step is to evaluate the binade in which  $r$  resides, and more precisely the position of the MSB of  $\tilde{r}$ . The main formula to compute  $m_r$  is given by:

$$m_r = \begin{cases} \lceil \log_2(r) \rceil & \text{if } r > 0 \\ \lfloor \log_2(-r) \rfloor + 1 & \text{if } r < 0 \end{cases} \quad (8)$$

For some very special cases (where  $r$  rounded to  $\tilde{r} \in [-2^{m_r-1}, 2^{m_r-1} - 2^{\ell_r}]$  while  $r$  is outside this interval), (8) needs to be adjusted (see [13]). Then the fractional part  $\ell_r$  is obtained from (7) and we have:

$$\tilde{r} = \lfloor r \cdot 2^{-\ell_r} \rfloor \cdot 2^{\ell_r}, \quad (9)$$

where  $\lfloor \cdot \rfloor$  is the round-to-nearest-integer operator. In the machine  $r$  is represented as a  $w$ -bit signed integer  $R = \lfloor r \cdot 2^{-\ell_r} \rfloor$ .

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

The following theorem is used to provide the maximum possible value for the outputs of the algorithm, and thus define their FxP formats (using (8)) with a guarantee that no overflow will occur. It is also used to determine the impact of the internal roundoff error on the output (see Section IV-C).

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

$$\mathcal{H} \begin{cases} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k) &= \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{cases} \quad (10)$$

with  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times q}$ ,  $\mathbf{C} \in \mathbb{R}^{p \times n}$  and  $\mathbf{D} \in \mathbb{R}^{p \times q}$ .

**Theorem 1 (Worst-Case Peak Gain theorem)** Let  $\mathcal{H}$  be a state-space system. If an input  $\{\mathbf{u}(k)\}_{k \geq 0}$  is known to be bounded by  $\bar{\mathbf{u}}$  ( $\forall k \geq 0$ ,  $|\mathbf{u}_i(k)| \leq \bar{u}_i$ ,  $1 \leq i \leq q$ ), then the output  $\{\mathbf{y}(k)\}_{k \geq 0}$  will be bounded<sup>3</sup> iff the spectral radius  $\rho(\mathbf{A})$  is strictly less than 1.

In that case, the output is (element-by-element) bounded by

$$\forall k, \quad |\mathbf{y}(k)| \leq \|\langle \mathcal{H} \rangle \bar{\mathbf{u}}\|, \quad (11)$$

where  $\langle \mathcal{H} \rangle$  is the Worst-Case Peak Gain (WCPG) matrix of the system [15], [16]. It can be computed as

$$\langle \mathcal{H} \rangle \triangleq |\mathbf{D}| + \sum_{k=0}^{\infty} |\mathbf{C}\mathbf{A}^k \mathbf{B}|. \quad (12)$$

Moreover, for any  $\varepsilon > 0$ , there always exists a finite input sequence  $\{\mathbf{u}(k)\}_{0 \leq k \leq N}$  such that the inequality (11) is almost attained for one output  $i$ , *i.e.*

$$|\mathbf{y}_i(N) - (\langle \mathcal{H} \rangle \bar{\mathbf{u}})_i| < \varepsilon. \quad (13)$$

There exist numerous other approaches to determine the output interval, such as involving Affine Arithmetic [17], which may lead to overestimations, or simulations [18], which in turn are inefficient and not rigorous. It can be shown that the WCPG approach gives the smallest interval containing all possible values of  $\mathbf{y}(k)$ . Moreover, in [2] an algorithm to evaluate the WCPG at arbitrary precision was proposed.

## III. SIMULINK-TO-SIF CONVERSION

### A. Simulink block diagram

Simulink, by Mathworks, is a widely used environment for model-based design of dynamic systems. It uses mainly a graphical block diagram to describe the model.

Internally, the block diagram is stored in an xml file (more specifically in `.slx` format). In this format the `<System>` tag contains the model description, and the `<Block>` and `<Line>` tags inside hold blocks of elements and their interconnections.

Consider a very simple data-flow diagram given on Fig. 4. It consists of standard gain, delay and sum blocks. All the

<sup>3</sup>This property is known as the Bounded Input Bounded Output (BIBO) stability [14].

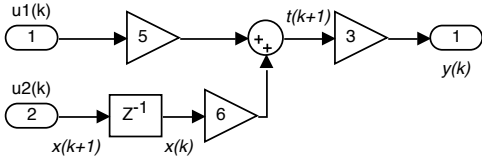


Fig. 4. A very simple Simulink block diagram.

information on the block parameters can be easily retrieved by parsing the xml file.

### B. The conversion algorithm

The SIF representation is a set of equations which compute the output, the states and the intermediate variables from the value of inputs, current states and the intermediate variables, and where each equation is a Sum-of-Products (SoPs).

In order to systematically transform any Simulink linear data-flow model to the SIF representation, we assume following rules. Since each delay element represents a computation that is saved from one time instance to another, in SIF we note it as a state variable. The gain elements are naturally interpreted in SIF as coefficients of the variables in corresponding SoPs. The outputs of gain and sum operators are naturally considered as temporary variables in SIF (unless they feed a delay or system output) since they are not stored from one system iteration to another. After noting all the variables of the system, their corresponding SoPs can be easily formed. Therefore, a system of equations corresponding to the data-flow can be obtained.

By design, most of the time there is a chain of sum blocks in a data-flow. For a finite-precision implementation the order of computations has a significant impact on the precision of the result. To apply our SoP computation technique presented in IV-B, merging all directly cascading sum blocks is necessary. Moreover, in case of a subsystem present in the diagram, a care was taken to flatten the design before getting blocks equations. For instance, the simple example of a data-flow diagram given in Fig. 4 corresponds to the following system of equations:

$$\begin{cases} t(k+1) &= 6 \cdot x(k) + 5 \cdot u_1(k) \\ x(k+1) &= 1 \cdot u_2(k) \\ y &= 3 \cdot t(k+1) \end{cases}$$

From this system of equations, it is then straightforward to identify matrices  $\mathbf{J}$ ,  $\mathbf{K}$ ,  $\mathbf{L}$ ,  $\mathbf{M}$ ,  $\mathbf{N}$ ,  $\mathbf{P}$ ,  $\mathbf{Q}$ ,  $\mathbf{R}$ ,  $\mathbf{S}$  of the SIF. An important detail is to preserve the order in which the variables appear in SoPs in order to match exactly the data-flow diagram. In our toy example, we have  $\mathbf{J} = \begin{pmatrix} 1 \\ \end{pmatrix}$ ,  $\mathbf{L} = \begin{pmatrix} 3 \\ \end{pmatrix}$ ,  $\mathbf{M} = \begin{pmatrix} 6 \\ \end{pmatrix}$ ,  $\mathbf{N} = \begin{pmatrix} 0 & 5 \\ \end{pmatrix}$ ,  $\mathbf{P} = \begin{pmatrix} 0 \\ \end{pmatrix}$ ,  $\mathbf{Q} = \begin{pmatrix} 0 & 1 \\ \end{pmatrix}$ , and the other matrices are null.

However, the order of labeling the temporary variables during the conversion can compromise the lower triangularity of the matrix  $\mathbf{J}$ . Since  $\mathbf{J}$  represents the order of computation of the intermediate variables within the data-flow and dependencies between them, it can be interpreted as an adjacency matrix for a directed acyclic graph. Therefore, a topological sort algorithm can be applied to determine the correct order of computations. Here the depth-first search algorithm was used

to ensure the lower triangular form of  $\mathbf{J}$  along with necessary reorganizations of matrices  $\mathbf{N}$ ,  $\mathbf{K}$  and  $\mathbf{L}$ .

## IV. FIXED-POINT IMPLEMENTATION AND ROUND-OFF ERROR ANALYSIS

### A. Fixed-Point formats

The first step of the fixed-point implementation is the determination of the FxP format of every internal variable within algorithm computation. Applying the WCPG theorem we can obtain the magnitudes for the outputs of the system. However, we can obtain the intervals for the intermediate and state variables as well [13]. To achieve this we concatenate the output vector  $\mathbf{y}$  with vectors  $\mathbf{t}$  and  $\mathbf{x}$ , and make necessary concatenations in the right side of (10). Thus, instead of considering the filter  $\mathcal{H}$  we consider the filter  $\mathcal{H}_u = (\mathbf{A}_Z, \mathbf{B}_Z, \mathbf{N}_1, \mathbf{N}_2)$  with

$$\mathbf{N}_1 \triangleq \begin{pmatrix} \mathbf{J}^{-1} \mathbf{M} \\ \mathbf{A}_Z \\ \mathbf{C}_Z \end{pmatrix}, \mathbf{N}_2 \triangleq \begin{pmatrix} \mathbf{J}^{-1} \mathbf{N} \\ \mathbf{B}_Z \\ \mathbf{D}_Z \end{pmatrix}. \quad (14)$$

Then, applying the WCPG theorem upon the filter  $\mathcal{H}_u$  we obtain the magnitudes  $\bar{\mathbf{t}}$ ,  $\bar{\mathbf{x}}$  and  $\bar{\mathbf{y}}$  by

$$\begin{pmatrix} \bar{\mathbf{t}} \\ \bar{\mathbf{x}} \\ \bar{\mathbf{y}} \end{pmatrix} = \langle \langle \mathcal{H}_u \rangle \rangle \bar{\mathbf{u}}. \quad (15)$$

The MSB positions  $m_t$ ,  $m_x$  and  $m_y$  are obtained with (8).

The evaluation of the magnitude can be improved by taking into account the propagation of roundoff errors that may, in some rare cases, yield an overflow. An algorithm for determining the FxP formats with a non-overflow guarantee was presented in [3].

### B. Sum-of-Products

As shown in Section II-B, the evaluation of linear data flow involves the evaluation of some sum-of-products, *i.e.* sums of variables (inputs, states and intermediate variables) multiplied by constant coefficients (those of the gain blocks). Let us consider here one sum-of-product (SoP):

$$s = \sum_{i=1}^N \mathbf{c}_i \cdot \mathbf{v}_i = \sum_{i=1}^N \mathbf{p}_i, \quad (16)$$

where  $\mathbf{p}_i$  denotes the product  $\mathbf{c}_i \times \mathbf{v}_i$ ,  $\mathbf{c}_i$  is a constant and  $\mathbf{v}_i$  a variable, for  $1 \leq i \leq N$ .

In our context, the MSB positions of the variables are known (deduced using the reliably computed magnitudes) and the MSB of the constants are computed using (8). In the SIF the sum-of-product results are stored either in the intermediate, state or output variables, so their MSB are also known. Thus, it is easy to deduce the FxP format of the accumulation in order to guarantee a roundoff error less than a given bound.

Denote  $(m_s, \ell_s)$  the required FxP format for the result. Then from [12], in order to guarantee a faithful rounding (*i.e.* the SoP evaluation error bounded by  $2^{\ell_s}$ ) it is sufficient to perform the multiplication and accumulation with the format  $(m_s, \ell_s - g)$  with  $g = \lceil \log_2 N \rceil$ . The  $g$  extra guard bits are discarded after the accumulation is performed, as shown in Fig. 5.

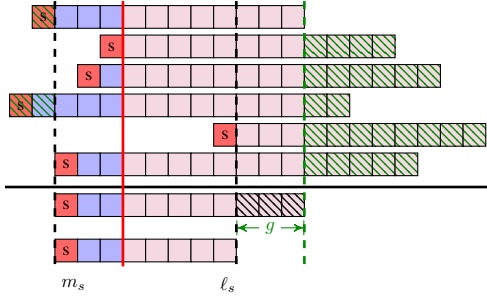


Fig. 5. Sum-of-Products performed with FxP format  $(m_s, l_s - g)$ .

### C. Roundoff error analysis

The evaluation of each SoP in fixed-point arithmetic may provide an additional error. When implemented, the Algorithm (2) is changed to

$$\begin{aligned} \mathbf{J}t^*(k+1) &\leftarrow \mathbf{M}\mathbf{x}^*(k) + \mathbf{N}\mathbf{u}(k) + \varepsilon_t(k) \\ \mathbf{x}^*(k+1) &\leftarrow \mathbf{K}t^*(k+1) + \mathbf{P}\mathbf{x}^*(k) + \mathbf{Q}\mathbf{u}(k) + \varepsilon_x(k) \\ \mathbf{y}^*(k) &\leftarrow \mathbf{L}t^*(k+1) + \mathbf{R}\mathbf{x}^*(k) + \mathbf{S}\mathbf{u}(k) + \varepsilon_y(k) \end{aligned} \quad (17)$$

where  $\varepsilon_t(k)$ ,  $\varepsilon_x(k)$  and  $\varepsilon_y(k)$  are the vectors of roundoff errors due to the sum-of-products evaluation.

Denote  $\varepsilon(k)$  the column vector that aggregates those error vectors; and  $\ell$  and  $w$  the vectors aggregating the LSB positions and word-lengths of the intermediate variables, states and outputs (in that order), respectively.

From the Section IV-B it follows that we can compute SoPs with faithful rounding. Using equations (7) and (15) we obtain that the errors  $\varepsilon(k)$  are element-by-element bounded by  $\bar{\varepsilon}$ :

$$\bar{\varepsilon} \triangleq 2^\ell = 2 \lceil \lceil \langle \mathcal{H}_u \rangle \bar{\mathbf{u}} \rceil_2 \times 2^{-w} \rceil, \quad (18)$$

where  $\lceil x \rceil_2$  is the operator that gives the power-of-2 immediately greater than  $x$ ,  $\times$  is the entrywise product and assuming the power of 2 of a vector is done element-by-element.

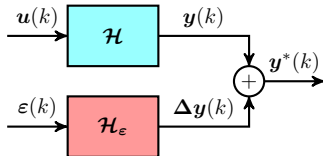


Fig. 6. Equivalent system, with errors separated.

In order to capture the effects of the FxP implementation we must take into account the propagation of the roundoff errors through the data-flow.

Due to the linearity of the system  $\mathcal{H}$ , we can express the implemented system as the initial system with an error  $\Delta y(k)$  added to the outputs, as shown in Fig. 6. The vector  $\Delta y(k)$  shows the propagation of the error  $\varepsilon(k)$  through the special error-system  $\mathcal{H}_\varepsilon$ . It is obtained by performing a difference between the implemented system (17) and the exact one. Its state-space representation is  $(\mathbf{A}_Z, \mathbf{M}_1, \mathbf{C}_Z, \mathbf{M}_2)$ , with:

$$\mathbf{M}_1 \triangleq (\mathbf{K}\mathbf{J}^{-1} \quad \mathbf{I}_n \quad \mathbf{0}), \quad \mathbf{M}_2 \triangleq (\mathbf{L}\mathbf{J}^{-1} \quad \mathbf{0} \quad \mathbf{I}_p). \quad (19)$$

Using such a decomposition and the bound (18) on the roundoff errors, we can apply the WCPG theorem upon the system  $\mathcal{H}_\varepsilon$  to obtain the maximum amplification of the errors

to the output. Then, for all time instances  $k$  the difference between the exact output  $y(k)$  and the implemented output  $y^*(k)$  is bounded by  $\overline{\Delta y}$ :

$$\overline{\Delta y} \triangleq 2 \langle \mathcal{H}_\varepsilon \rangle (\lceil \langle \mathcal{H}_u \rangle \bar{\mathbf{u}} \rceil_2 \times 2^{-w}). \quad (20)$$

Therefore, we have a reliable analytical bound on the data-flow implementation error, which takes into account the (possible) non-linear propagation of roundoff errors. As we can see, our approach does not involve any simulations.

### D. Word-length optimization

We can note that the output error bound  $\overline{\Delta y}$  depends only on the word-lengths vector  $w$ . For hardware implementation, where it is possible to have different word-lengths for each variable and operator (*multiple word-length paradigm*), it should be possible to find the word-lengths that minimize a criteria, like area, power consumption, etc., while satisfying a constraint on the output error.

Although this is out of the scope of this article, such optimization problem can be solved when the cost function is not too complicated: for linear cost function (like the sum of the word-lengths), it can be solved with some Mixed-Integer Non Linear Programming solver, like Bonmin<sup>4</sup>.

When considering homogeneous word-lengths, the minimum word-length  $w_{\min}$  that guarantees that the output error is less than a given  $\varepsilon$  is

$$w_{\min} = \left\lceil \log_2 \left( \frac{\langle \mathcal{H}_\varepsilon \rangle (\lceil \langle \mathcal{H}_u \rangle \bar{\mathbf{u}} \rceil_2)}{\varepsilon} \right) \right\rceil + 1. \quad (21)$$

It can be used to compare realizations.

## V. ILLUSTRATIVE EXAMPLES

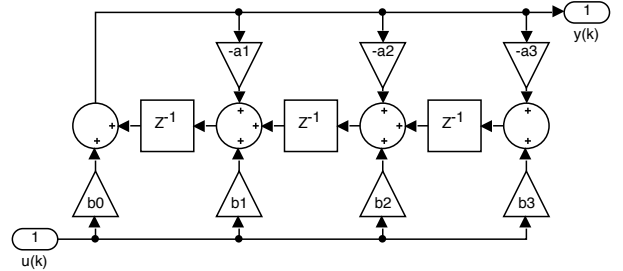


Fig. 7. Direct Form II transposed data-flow.

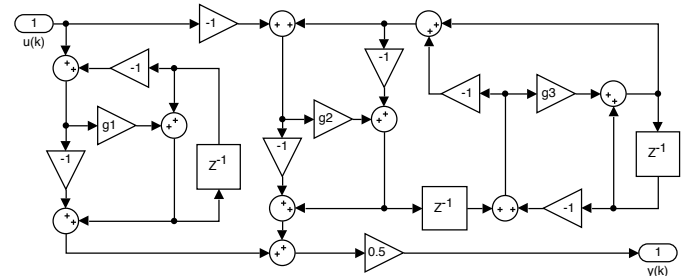


Fig. 8. Lattice Wave Digital Filter data-flow.

<sup>4</sup>Basic Open-source Nonlinear Mixed INteger programming, <https://projects.coin-or.org/Bonmin>

To illustrate our flow from a linear data-flow diagram to reliable fixed-point implementation, we use two realizations of the same 3rd order Butterworth filter with cutoff frequency 0.2. The transfer function was designed with Matlab. The first realization is a Direct Form II transposed, which requires 7 coefficients to be stored. The second one is a Lattice Wave Digital Filter [19], which requires only 3 coefficients. The corresponding data-flows are given in Fig. 7 and Fig. 8.

The two data-flow graphs were translated into our internal framework using the Simulink-to-SIF converter described in the Section III. Then, for a given  $\bar{u}$ , equal here to 5.25, the MSBs of the intermediate variables, states and output were deduced using (15).

For a given homogeneous word-length  $w$  (i.e. the same  $w$  for all variables), the error for each sum-of-products was deduced along with the final output error, which takes into account the propagation of the computational errors.

Finally, instead of exhibiting the output error with respect to the variable word-lengths, we prefer to use a more realistic indicator, like the total number of bits used for the computations as shown in Fig. 9. We can observe that Lattice Wave realization requires less coefficients to be stored and, consequently, slightly smaller total amount of bits than the Direct Form II transposed to attain the same output error. In our example, the Direct Form II transposed realization requires 17 bits per variable to ensure an output error less than  $\varepsilon = 10^{-2}$ , whereas 16 bits per variable are enough for the Lattice Wave one.

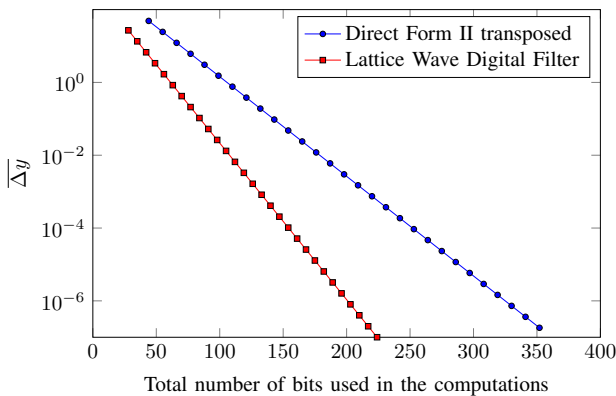


Fig. 9. Output error when the word-length  $w$  goes from 4 to 32.

## VI. CONCLUSION

We proposed a reliable method to implement linear data-flows with fixed-point arithmetic, while ensuring that no overflow occurs and controlling the output error. It is based on an analytical description that describes all the computations of a data-flow, and allows performing analytical roundoff error analysis, thanks to the Worst-Case Peak Gain theorem.

The output error depends exclusively on the word-lengths used to store the intermediate variables and states, and to perform the sum-of-products. However, these variables have different impact on the overall output error. Therefore, as a perspective, we should focus on solving the word-length

optimization problem, in order to find the word-lengths that minimize a surface or power consumption criteria, while guaranteeing an output error constraint. Finally, the corresponding code (C or VHDL) will be generated with third-party tools, like FloPoCo<sup>5</sup> [20] or Stratus<sup>6</sup> [21].

## REFERENCES

- [1] 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.
- [2] A. Volkova, T. Hilaire, and C. Lauter, "Reliable evaluation of the worst-case peak gain matrix in multiple precision," in *Computer Arithmetic (ARITH), 2015 IEEE 22nd Symposium on*, June 2015, pp. 96–103.
- [3] —, "Determining fixed-point formats for a digital filter implementation using the worst-case peak gain measure," in *2015 49th Asilomar Conference on Signals, Systems and Computers*, Nov 2015, pp. 737–741.
- [4] G. Li, "On the structure of digital controllers with finite word length consideration," *IEEE Trans. on Autom. Control*, vol. 43, May 1998.
- [5] M. Gevers and G. Li, *Parameterizations in Control, Estimation and Filtering Problems*. Springer-Verlag, 1993.
- [6] R. Istepanian and J. Whidborne, Eds., *Digital Controller implementation and fragility*. Springer, 2001.
- [7] B. Widrow and I. Kollár, *Quantization Noise: Roundoff Error in Digital Computation, Signal Processing, Control, and Communications*. Cambridge, UK: Cambridge University Press, 2008.
- [8] 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.
- [9] T. Hinamoto, H. Ohnishi, and W. Lu, "Roundoff noise minimization of state-space digital filters using separate and joint error feedback/coordinate transformation optimization," in *IEEE Transactions on Circuits and Systems, Fundamental Theory and Applications*, vol. 50, 2003.
- [10] D. Menard and O. Sentieys, "A methodology for evaluating the precision of fixed-point systems," in *ICASSP, 2002*, pp. 3152–3155.
- [11] T. Hilaire and P. Chevrel, "Sensitivity-based pole and input-output errors of linear filters as indicators of the implementation deterioration in fixed-point context," *EURASIP Journal on Advances in Signal Processing*, vol. special issue on Quantization of VLSI Digital Signal Processing Systems, January 2011.
- [12] B. Lopez, T. Hilaire, and L.-S. Didier, "Formatting bits to better implement signal processing algorithms," in *4th international Conference on Pervasive and Embedded Computing and Communication Systems (PECCS)*, Lisbon, Portugal, Jan. 2014.
- [13] T. Hilaire and B. Lopez, "Reliable implementation of linear filters with fixed-point arithmetic," in *Proc. IEEE Workshop on Signal Processing Systems (SiPS)*, 2013.
- [14] T. Kailath, *Linear Systems*. Prentice-Hall, 1980.
- [15] V. Balakrishnan and S. Boyd, "On computing the worst-case peak gain of linear systems," *Systems & Control Letters*, vol. 19, 1992.
- [16] S. P. Boyd and J. Doyle, "Comparison of peak and rms gains for discrete-time systems," *Syst. Control Lett.*, vol. 9, no. 1, pp. 1–6, June 1987.
- [17] 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, November 2007.
- [18] 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, November 1998.
- [19] A. Fettweis, "Wave digital filters: Theory and practice," *Proc. of the IEEE*, vol. 74, no. 2, 1986.
- [20] F. de Dinechin and B. Pasca, "Designing custom arithmetic data paths with FloPoCo," *IEEE Design & Test of Computers*, vol. 28, no. 4, pp. 18–27, Jul. 2011.
- [21] S. Belloeil, D. Dupuis, C. Masson, J.-P. Chaput, and H. Mehrez, "Stratus: A procedural circuit description language based upon Python," in *International Conference on Microelectronics*, 2007, pp. 275–278.

<sup>5</sup><http://flopoco.gforge.inria.fr/>

<sup>6</sup><https://soc-extras.lip6.fr/en/coriolis/stratus/>