



**HAL**  
open science

# Sum-of-products Evaluation Schemes with Fixed-Point arithmetic, and their application to IIR filter implementation

Benoit Lopez, Thibault Hilaire, Laurent-Stéphane Didier

► **To cite this version:**

Benoit Lopez, Thibault Hilaire, Laurent-Stéphane Didier. Sum-of-products Evaluation Schemes with Fixed-Point arithmetic, and their application to IIR filter implementation. Conference on Design and Architectures for Signal and Image Processing (DASIP), Oct 2012, Karlsruhe, Germany. hal-01076043

**HAL Id: hal-01076043**

**<https://hal.sorbonne-universite.fr/hal-01076043>**

Submitted on 20 Oct 2014

**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.

# Sum-of-products Evaluation Schemes with Fixed-Point arithmetic, and their application to IIR filter implementation

Benoit Lopez, Thibault Hilaire and Laurent-Stéphane Didier  
LIP6, University Pierre et Marie Curie (UPMC), Paris, France

**Abstract**—The signal processing and control algorithms are widely based on sum-of-products evaluation. In fixed-point arithmetic, the roundoff errors and coefficient quantization may have an important effect on the application’s performance and characteristics.

As part of a global methodology on optimal fixed-point implementation of filters/controllers, this paper formalizes the various implementation schemes for sum-of-products in fixed-point arithmetic and automates the fixed-point code production. The order of the operations are considered, as their bit-width and the fixed-point representation of the coefficients, variables and partial results. Applied to linear filters, the output roundoff noise error is then evaluated and used as a criteria to find out interesting evaluation scheme. An example illustrates the approach.

## I. INTRODUCTION

The great majority of embedded signal processing or control algorithms is implemented using digital devices such as general purpose micro-controllers, DSP, ASIC or FPGA. The computation on these devices is mainly based on integer arithmetic (rather than floating-point arithmetic) for cost, size and power consumption reasons. The fixed-point arithmetic is used as an approximation of real numbers. Unfortunately the numerical implementation of such algorithms suffers from a deterioration in performance and characteristics, due to the quantization of the embedded coefficients and the roundoff errors occurring in the computations [1], [2], [3].

Classically, the filter implementation problem is tackled in several steps [4]. First, for a given filter or controller, a *good* realization is chosen, such as a direct form, a state-space realization, a cascade decomposition, etc. Criteria based on transfer function sensitivity, or pole sensitivity [1], [5] can be used as indicators of the Finite Word-Length (FWL) effects, in order to compare the realizations. Then, the associated algorithm is realized in fixed-point arithmetic on the dedicated device.

Since the considered algorithms (Linear Time Invariant systems, such as FIR<sup>1</sup> or IIR<sup>2</sup> filters, linear controllers, etc.) are highly regular, we focus in this article on the study of the sum-of-products implementation in fixed-point arithmetic, and the evaluation of the associated roundoff errors. These errors are highly dependent on the order of the additions [6], on each binary-point position and the bit-width of each operator.

So, a complete formalization of sum-of-products in fixed-point arithmetic is proposed here. It is much more detailed and accurate than regular approaches used, in the sense that it provides various possibilities, such as a dedicated fixed-point representation for each coefficient, variable and partial result and it supports the multiple word-length paradigm [7].

We do not focus here on the optimal realization problem (see [8]) nor on hardware consideration (*i.e.* how to map the algorithm on a given architecture), but mainly on arithmetic and algorithmic organization of the computations and its effects on precision. The main objective is to formalize all the possible fixed-point implementations of a given sum-of-products in order to exhibit the *interesting* ones (in a sense to define). This part will be used in the global implementation flow from filter/controller to code [4] and is common to hardware and software implementation. It will be even important for FPGA implementation, where fine-grained computing should be exploited to obtain good tradeoff between implementation cost and precision.

The paper is organized as follows. Section II reminds the FWL effects on IIR filter and the dynamic range evaluation. Then, after reminding the fixed-point arithmetic and its associated operations, section III considers the ordered sum-of-products and the propagation rules of the fixed-point format, with different variations. Finally, an illustrative example is given in section IV before conclusion in section V.

## II. FINITE WORD-LENGTH EFFECTS AND DYNAMIC RANGE EVALUATION

Let us consider a  $n$ -th order IIR filter with  $h$  as transfer function:

$$h(z) = \frac{\mathbf{b}_0 + \mathbf{b}_1 z^{-1} + \dots + \mathbf{b}_n z^{-n}}{1 + \mathbf{a}_1 z^{-1} + \dots + \mathbf{a}_n z^{-n}}, \quad \forall z \in \mathbb{C}. \quad (1)$$

This filter is usually realized with the following algorithm

$$y(k) = \sum_{i=0}^n \mathbf{b}_i u(k-i) - \sum_{i=1}^n \mathbf{a}_i y(k-i) \quad (2)$$

where  $u(k)$  and  $y(k)$  are the input and output at step  $k$ , respectively.

Of course, this direct form (2) is not the only one possible realization. Some other interesting realizations can be considered, such as state-space realizations [1], realizations with the  $\delta$ -operator (defined by  $\delta \triangleq \frac{q-1}{\Delta}$ , where  $q$  is the classical

<sup>1</sup>Finite Impulse Response

<sup>2</sup>Infinite Impulse Response

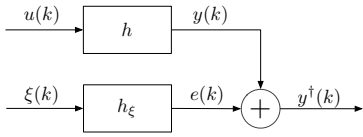


Fig. 1. Equivalent system, with noise extracted

shift operator, and  $\Delta$  a strictly positive constant [9]), the  $\rho$ -Direct Form II transposed ( $\rho$ DFII<sub>t</sub>) [10], the  $\rho$ -modal forms, the wave lattice filter, warped filter, and a lot of other specific realizations (LGC, LCW-structures [11], etc.).

The Specialized Implicit Framework (SIF) proposed in [8] can be used as a unifying framework to describe all these realizations. In that context, this work on sum-of-products evaluation scheme will be extended to that framework, in order to be applied on various realizations.

#### A. Roundoff Noise Analysis

When implemented in finite precision, equation (2) is modified by the addition of some noise  $\xi(k)$ , and only  $y^\dagger$  (the output contaminated with roundoff error) can be computed:

$$y^\dagger(k) = \sum_{i=0}^n \mathbf{b}_i u(k-i) - \sum_{i=1}^n \mathbf{a}_i y^\dagger(k-i) + \xi(k) \quad (3)$$

This added noise depends on:

- the way the computations are organized and realized (the order of the sums, etc.);
- the fixed-point representation of all the signals used in the computations (inputs, outputs, coefficients)
- and the fixed-point representation of each step of the operations (inner signals)

In fixed-point arithmetic, it can be modeled as independent white noise with given first and second order moments (see section III-F).

**Proposition 1** *It is possible to express the implemented system as the initial system with a noise  $e(k)$  added on the output, as shown in Figure 1.*

Moreover  $e(k)$  is the result of the noise  $\xi(k)$  through the filter  $h_\xi$  defined by:

$$h_\xi(z) = \frac{1}{1 + \sum_{i=1}^n \mathbf{a}_i z^{-i}}, \quad \forall z \in \mathbb{C}. \quad (4)$$

*Proof:*  $e(k)$  is defined by  $e(k) \triangleq y^\dagger(k) - y(k)$ , so

$$e(k) = - \sum_{i=1}^n \mathbf{a}_i (y^\dagger(k-i) - y(k-i)) + \xi(k) \quad (5)$$

$$= \xi(k) - \sum_{i=1}^n \mathbf{a}_i e(k-i) \quad (6)$$

and then  $e(k)$  is the result of noise  $\xi(k)$  through the transfer function  $h_\xi$ . ■

The first ( $\mu$ ) and second ( $\sigma^2$ ) order moments of a noise  $\xi$  are defined and denoted by:

$$\mu_\xi \triangleq E \{ \xi(k) \} \quad (7)$$

$$\sigma_\xi^2 \triangleq E \{ (\xi - \mu_\xi)^2(k) \}, \quad (8)$$

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

One classical criteria used to measure the roundoff noise effect is the output noise power:

**Definition 1 (Output noise power)** *The output noise power is defined as the power of the noise added on the output:*

$$p \triangleq E \{ e^2(k) \}. \quad (9)$$

**Proposition 2** *The output noise power is obtained from the first and second-order moment of  $\xi$  by:*

$$p = \|h_\xi\|_2^2 \sigma_\xi^2 + (h_\xi(1)\mu_\xi)^2 \quad (10)$$

where  $\|h_\xi\|_2$  is the  $\ell_2$ -norm of the filter  $h_\xi$  and  $h_\xi(1)$  its DC-gain:

$$h_\xi(1) = \frac{1}{1 + \sum_{i=1}^n \mathbf{a}_i}. \quad (11)$$

*Proof:* Eq. (9) leads to  $p = \sigma_e^2 + \mu_e^2$ . Moreover, the basic properties of noise transmission through a linear system gives [12], [13]

$$\mu_e = h_\xi(1)\mu_\xi, \quad \sigma_e^2 = \|h_\xi\|_2^2 \sigma_\xi^2. \quad (12)$$

■

#### B. Quantization of the coefficients

In addition to roundoff noise, the coefficients of the implemented filter are different from the ideal filter. Indeed, the coefficients are approximated by their fixed-point best approximation because of the finite precision of the representation of real numbers (see section III-A).

We denote  $h^\dagger$  the implemented transfer function  $h$ , i.e. the transfer function with quantized coefficients:

$$h^\dagger(z) = \frac{\sum_{i=0}^n \mathbf{b}_i^\dagger z^{-i}}{1 + \sum_{i=1}^n \mathbf{a}_i^\dagger z^{-i}}, \quad \forall z \in \mathbb{C}. \quad (13)$$

where the  $(\mathbf{a}_i^\dagger)$  and  $(\mathbf{b}_i^\dagger)$  are the quantized approximations of  $(\mathbf{a}_i)$  and  $(\mathbf{b}_i)$ , respectively. These approximations depend on the fixed-point representation used.

The criterion used to measure the transfer function deviation is the following:

**Definition 2 (Transfer function error)** *The transfer function error  $\Delta h$  due to the quantization is defined as the  $\ell_2$ -norm of the gap between the ideal transfer function and the implemented one, i.e.:*

$$\Delta h \triangleq \|h - h^\dagger\|_2. \quad (14)$$

**Remark 1** This measure depends explicitly on the exact fixed-point representation (bit-width, binary-point position)

used for each coefficient. Since this information is not always available when the FWL effects are studied, a more interesting criterion is preferred: the  $\ell_2$ -sensitivity measure. A lot of work have been done in order to evaluate this sensitivity with respect to the coefficients and evaluate how much the coefficients' quantization changes the transfer function and the poles or zeros. For further explanations, see [14], [1], [5].

### C. Dynamic Range analysis

If  $x(k)$  is a scalar signal, we denote  $\|x\|_\infty$  the  $L_\infty$ -norm of this signal, *i.e.*

$$\|x\|_\infty \triangleq \sup_k |x(k)|. \quad (15)$$

The following lemma reminds the basic properties of a Bounded Input Bounded Output (BIBO) system, and is useful to determine the fixed-point format of the output of the considered filter.

**Lemma 1** *The filter  $h$  is said to be BIBO if the  $\ell_1$ -norm of  $h$  is finite, with*

$$\|h\|_{\ell_1} \triangleq \sum_{k=0}^{\infty} |h(k)| \quad (16)$$

and  $\{h(k)\}_{k \geq 0}$  is the impulse response of  $h$ . Moreover, in that case, the peak value of the output  $y$  can be bounded by [15], [16]

$$\|y\|_\infty \leq \|h\|_{\ell_1} \|u\|_\infty. \quad (17)$$

where  $u$  is the input.

**Remark 2** The  $\ell_1$ -norm can be very conservative, specially if the input  $u$  is known to have some particularities. For example, if  $u$  has a known frequency distribution, it can be interesting to use a weighted  $\ell_1$ -norm, in order to avoid over-estimation of  $\|y\|_\infty$ .

## III. FIXED-POINT SUM-OF-PRODUCTS

The sum-of-products (SoP) considered in this paper is a dot-product between a constant vector  $\mathbf{c} = (c_i)_{1 \leq i \leq n}$  and a variable vector  $\mathbf{x} = (x_i)_{1 \leq i \leq n}$ :

$$r = \mathbf{c} \cdot \mathbf{x} = \sum_{i=1}^n c_i x_i. \quad (18)$$

In this section, a lot of possible schemes for evaluate  $r$  will be shown, so the problem can be stated as follow :

**Problem 1** *Given a sum-of-products as in eq. (18) where the constants  $c_i$ 's, and their wordlength, the fixed-point representation of the variable  $x_i$ 's and the fixed-point representation of the result  $r$  are known, determine all the possible evaluation schemes for (18) and parametrize them.*

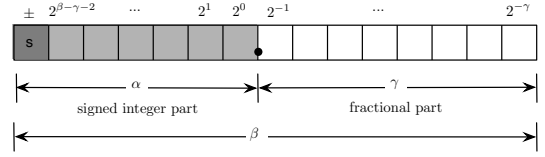


Fig. 2. FPR scheme

### A. Fixed-point arithmetic and operations

Compared to floating-point arithmetic, fixed-point arithmetic is a number representation widely used in embedded system because of its ratio precision/lightness.

Fixed-Point Representation (FPR) can be formalized by the tuple  $(\beta, \alpha, \gamma)$ , where  $\beta$  is the bit-width of a number  $x$  in fixed-point arithmetic,  $\alpha$  the number of bits of the integer part and  $\gamma$  the number of bits of its fractional part (see Figure 2). Number representation used in this paper is signed number with 2's complement, so the sign bit is included in  $\alpha$ .

Let  $x \in \mathbb{R}^*$  be a real number we want to convert in a signed-fixed-point number. The bit-width of  $x$  ( $\beta_x$ ) is a chosen value, generally a multiple of 8 bits. Then, the length of its integer part is determined as the exact number of bits necessary to write its integer part, so:

$$\alpha_x \triangleq \lfloor \log_2 |x| \rfloor + 2 \quad (19)$$

with  $\lfloor x \rfloor$  the operation that rounds  $x$  to the nearest integer lower than or equal to  $x$ .

Note that if an unsigned representation were used, we should have  $\alpha_x = \lfloor \log_2 |x| \rfloor + 1$ .

Finally the fractional part is evaluated as  $\gamma_x = \beta_x - \alpha_x$ , and the FPR of  $x$  is defined as  $FPR_x = (\beta_x, \alpha_x, \gamma_x)$ .

Moreover, in fixed-point arithmetic,  $x$  is represented by the integer  $N_x$  determined by the first  $\beta_x$  significant bits of its binary representation:

$$N_x \triangleq \lfloor x \cdot 2^{\gamma_x} \rfloor \quad (20)$$

where  $\lfloor x \rfloor$  is the operation that rounds  $x$  to the nearest integer. So,  $x$  is approximated by  $x^\dagger = N_x \cdot 2^{-\gamma_x}$ .

After those recalls, the only two types of operations needed for the sum-of-products, *i.e.* additions and multiplications, can be described.

For both of multiplication and addition, FPR's of operands and bit-width of the result are given data, and complete FPR of the result (*i.e.* its  $\alpha$  and  $\gamma$  values) need to be determined.

1) *Multiplication:* Let  $z = x \times y$  be a multiplication, with  $(\beta_x, \alpha_x, \gamma_x)$  and  $(\beta_y, \alpha_y, \gamma_y)$  the respective FPR of  $x$  and  $y$ .  $FPR_z$  is given by:

$$FPR_z = (\beta_x + \beta_y, \alpha_x + \alpha_y, \gamma_x + \gamma_y) \quad (21)$$

and the operation is realized by  $N_z \leftarrow N_x \times N_y$ .

In fact, this is the optimal case, because generally the bit-width of the hardware operation is fixed, and possibly different from  $\beta_x + \beta_y$ .

**Proposition 3 (Fixed-Point Multiplication)** Let  $\beta_{op}$  be the fixed bit-width of the result. Then, in the general case, the result FPR of the operation  $z = x \times y$ , i.e.  $FPR_z$ , is deduced from the operands' FPR, from  $\beta_{op}$ ,  $FPR_x$  and  $FPR_y$ , by:

$$FPR_z = (\beta_{op}, \alpha_x + \alpha_y, \beta_{op} - (\alpha_x + \alpha_y)) \quad (22)$$

and the operation is realized by

$$N_z \leftarrow (N_x \gg s_x) \times (N_y \gg s_y) \quad (23)$$

where  $\gg$  is the right-shift operation (with troncation or round-to-the-nearest),  $s_x$  and  $s_y$  are signed right bit-shifts to apply on  $N_x$  and  $N_y$  such that  $s_x + s_y = \beta_x + \beta_y - \beta_{op}$  (if  $\beta_{op} = \beta_x + \beta_y$ , then  $s_x = s_y = 0$  is choosen).

See Figures 5(a) and 5(b) for an illustrative example.

**Remark 3** As said in Proposition 3,  $\beta_{op}$  is a given data, so determine  $FPR_z$  means determine  $\alpha_z$  or  $\gamma_z$  (the other one is given by completion with  $\beta_{op}$ ).

2) *Addition*: The other operation required, for which we have a constraint. Indeed, for adding two numbers in fixed-point arithmetic, they need to share the same FPR. Let  $z = x + y$  be an addition, with  $FPR_x$  and  $FPR_y$  as defined in previous subsection.  $FPR_z = (\beta_z, \alpha_z, \gamma_z)$  is determined in the optimal case by:

$$\begin{cases} \alpha_z = \max(\alpha_x, \alpha_y) + 1 \\ \gamma_z = \max(\gamma_x, \gamma_y) \\ \beta_z = \alpha_z + \gamma_z \end{cases} \quad (24)$$

Once again, generally the bit-width of the operator, i.e.  $\beta_{op}$  may be different than the optimal one. The following proposition gives the fixed-point representation for the general case:

**Proposition 4 (Fixed-Point Addition)** The result FPR  $FPR_z = (\beta_z, \alpha_z, \gamma_z)$  is inferred from the input's FPR in the general case, from  $\beta_{op}$ ,  $FPR_x$  and  $FPR_y$ , by:

$$\begin{cases} \alpha_z = \max(\alpha_x, \alpha_y) + 1 \\ \gamma_z = \beta_{op} - \max(\alpha_x, \alpha_y) - 1 \\ \beta_z = \beta_{op} \end{cases} \quad (25)$$

Moreover,  $z = x + y$  is realized by:

$$N_z \leftarrow (N_x \gg s_x) + (N_y \gg s_y) \quad (26)$$

with  $s_x = \gamma_x - \gamma_z$  and  $s_y = \gamma_y - \gamma_z$ .

## B. Ordered Sum-of-products

A consequence of proposition 4 is that the addition in fixed-point arithmetic is commutative but not necessarily associative. Indeed, before an addition, one of the two operands will be quantized. If  $\oplus$  represents the fixed-point addition,  $(a \oplus b) \oplus c$  could be different from  $a \oplus (b \oplus c)$ , because the fixed-point representation of  $(a \oplus b)$  and  $(b \oplus c)$  can be different, leading to different quantizations of  $a$ ,  $b$  and  $c$  in the two cases.

Consequently, a given dot-product can be realized in fixed-point arithmetic with several sum-of-products with particular

order, depending on the parenthesizings of the sum. We denote them *ordered-SoP* (oSoP).

Moreover, for a  $n^{\text{th}}$ -order SoP, there are  $\prod_{i=1}^{n-1} (2i - 1)$  different oSoPs to consider (the proof can be obtained by induction) [17].

For instance, for  $n = 6$ , there are 945 possibilities. Here are two of them, with  $z_i = c_i \cdot x_i$  (see Figures 3(a) and 3(b)):

$$r = ((z_0 + z_3) + z_2) + ((z_4 + z_5) + z_1) \quad (27)$$

and

$$r = (((z_3 + z_5) + z_2) + z_4) + (z_0 + z_1) \quad (28)$$

Obviously, they can lead to different result, when implemented in fixed-point arithmetic.

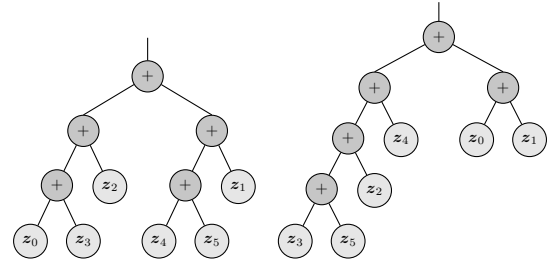


Fig. 3. Two 6<sup>th</sup>-order oSoPs, among the 945 possibilities

Thus, it is of interest to consider all of them and find the most interesting ones, with respect to roundoff noise analysis (and for our further work with respect to architecture adequacy).

It is quite easy to generate them all, except that it is not reasonable from  $n \geq 10$ .

## C. FPR propagation in an oSoP

The FPR propagation is the way to determine the FPR among the various operations of the oSoP.

In our approach, oSoP are binary trees, where internal nodes are adders and multipliers and leaves are the couples  $(c_i, x_i)_{1 \leq i \leq n}$ . In fact, if we consider a multiplier as a cell composed by the operator and its operands, then all leaves of the trees are these multiplier cells and internal nodes are adders.

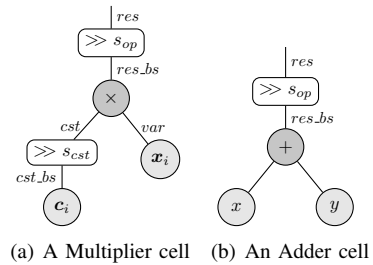


Fig. 4. Graphical representations of adder and multiplier cells

The two operators used in our oSoP and their associated rules of propagation can be now described .

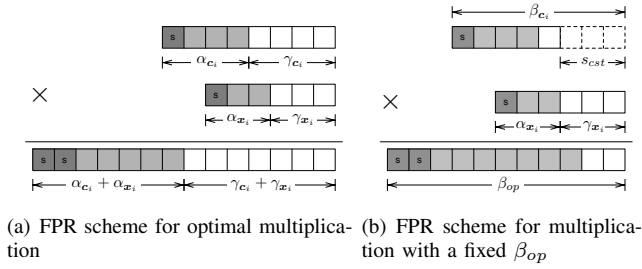


Fig. 5. Optimal and general multiplication

1) *Multiplier cell*: As shown in Figure 4(a), the multiplier cell is composed by a multiplier ( $\beta_{op}$  as resulting bit-width), followed by a potential signed right bit-shift of  $s_{op}$  bits, used to align the result. The result of the multiplication, before the shift, is denoted  $res.bs$  (*result before shift*), whereas the shifted result is denoted  $res$ .

The first operand of the multiplier is the variable  $x_i$  (with  $FPR_{var}$  as FPR) and the second operand is the constant  $c_i$  potentially followed by a signed right-shift of  $s_{cst}$  bits.  $cst.bs$  and  $cst$  denote the constant before shift and the shifted constant, respectively.

The following parameters are supposed to be known:  $\beta_{op}$  the bit-width of the multiplier,  $FPR_{var} = (\beta_{x_i}, \alpha_{x_i}, \gamma_{x_i})$  the fixed-point representation of the variable  $x_i$  and  $FPR_{cst.bs} = (\beta_{c_i}, \alpha_{c_i}, \gamma_{c_i})$  the fixed-point representation of the constant  $c_i$  (obtained with equation (19)).

Then the propagation is made by determining  $FPR_{res.bs}$  and  $s_{op}$  from these three values and proposition 3:

- The right shift is only applied here on the constant, so

$$s_{cst} = \beta_{c_i} + \beta_{x_i} - \beta_{op}, \quad (29)$$

- the FPR of the constant operand is then

$$FPR_{cst} = (\beta_{op} - \beta_{x_i}, \alpha_{c_i}, \gamma_{c_i} - s_{cst}), \quad (30)$$

- and finally the FPR of the result (before a potential shift) is

$$FPR_{res.bs} = (\beta_{op}, \alpha_{c_i} + \alpha_{x_i}, \beta_{op} - \alpha_{c_i} - \alpha_{x_i}). \quad (31)$$

Figures 5(a) and 5(b) illustrate the operation.

2) *Adder cell*: The adder cell, as shown in Figure 4(b) is composed by an adder (with  $\beta_{op}$  as bit-width), followed by a potential signed right bit-shift of  $s_{op}$  bits. The same notation  $res.bs$  (*result before shift*) and  $res$  is used for consistency. The two operands are denoted  $x$  and  $y$ .

In our approach, propagation algorithm is a recursive algorithm, *i.e.* it is applied on the final adder of the oSoP (the tree root) which calls the algorithm for its two operands, and so on.

This algorithm, applied on an adder cell, determines the FPR of the result, before the potential shift. It is composed of three steps:

- Apply the propagation algorithm to its two operand  $x$  and  $y$  in order to obtain their respective  $FPR_{res.bs}$ . Note that

if  $x$  or  $y$  are a multiplier cell, then equations (29) to (31) are used.

- A common FPR is determined for the addition, by applying proposition 4:

$$FPR_{res.bs} = (\beta_{op}, \max(\alpha_x, \alpha_y) + 1, \beta_{op} - \max(\alpha_x, \alpha_y) - 1) \quad (32)$$

Moreover, since the FPR of the final result of the sum-of-products is supposed to be known, we can limit the integer part of the partial results to the integer part of the final result. So the  $FPR_{res.bs}$  is given by  $FPR_{res.bs} = (\beta_{op}, \alpha_{op}, \beta_{op} - \alpha_{op})$  with

$$\alpha_{op} = \min(\max(\alpha_x, \alpha_y) + 1, \alpha_r) \quad (33)$$

where  $\alpha_r$  is the integer part of the final result  $r$  (even if an intermediate overflow occurs here, it will be compensated in a next addition, thanks to Jackson's rule<sup>3</sup>[18]).

- Then, we know what must be the FPR of the two operands. They are the  $FPR_{res}$  of the operands  $x$  and  $y$ . Since their  $FPR_{res.bs}$  is known, the right shift  $s$  to be applied can be evaluated with

$$s_x = \gamma_x - \gamma_{res.bs}, \quad s_y = \gamma_y - \gamma_{res.bs}. \quad (34)$$

It corresponds to the  $s_{op}$  of the two cells  $x$  and  $y$ .

Figures 6(a), 6(b) and 6(c) illustrate respectively the steps a), b) and c) of the propagation algorithm.

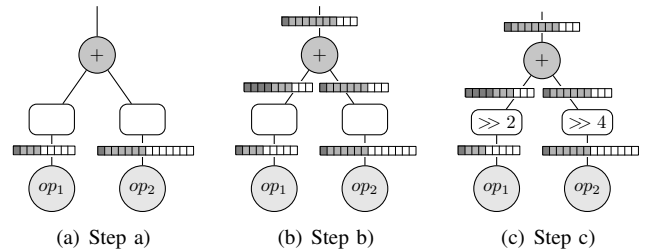


Fig. 6. Propagation algorithm for the addition (with  $\beta_{op} = 10$ )

In addition, some further rules can be added for best roundoff or propagation, with respect to material architectures or extra optimization. We are going to discuss these options in the next parts.

#### D. Roundoff Before and/or After Multiplication

Inside the multiplication cell, it has been seen that there could be a shift on the constant when the multiplier bit-width is smaller than the optimal bit-width. Moreover, an adder could yield to a shift on its operands, like a multiplier.

So a multiplier cell may have two different shifts, one on the constant and the other on the result (see Figure 7(a)). Of course, the first one is a *virtual* shift, *i.e.* there is no real shift to perform, the constant is only replaced by its shifted value.

<sup>3</sup>Jackson's Rule says that in any number of consecutive additions and/or subtractions, some intermediate results and operands may overflow. As long as the final result representation can handle the final result without overflow, then the result is valid. This is due to the fact that fixed-point arithmetic is a modulo arithmetic.

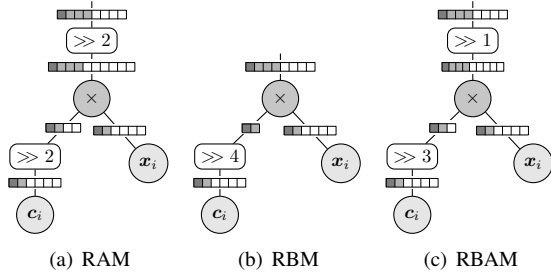


Fig. 7. Roundoff Before and/or After Multiplication

Depending on the architecture, there is not always the possibility to proceed to a shift after the multiplication (for example, most DSP doesn't have a barrel shifter in their Multiplier-and-Accumulator (MAC) operator, and do a shift after the multiplication implies a lot of extra cycles). In that case, a Roundoff After Multiplication (RAM) scheme is not possible, so we have to group the two shifts onto the constant shift. These scheme is called Roundoff Before Multiplication (RBM): the bit-shift  $s_{op}$  is then deferred onto the bit-shift  $s_{cst}$  (see Figure 7(b)).

RBM doesn't produce any roundoff noise after the multiplication, but there is a loss of precision into the coefficients used, leading to a greater transfer function error. This effect mainly depends on the sensitivity of the application (here a filter or controller) with respect to its coefficients.

Another option is to let the roundoff after multiplication. Indeed, as it said below, and if material architecture permits it, RAM saves constant informations for multiplication. One can note that we can't defer  $s_{cst}$  onto  $s_{op}$  in this case because  $s_{cst}$  is due to a bit adjustment compared to operator width.

A third and last option is a roundoff before and after multiplication (RBAM), where we choose to defer a part of the bit-shift  $s_{op}$  onto the bit-shift  $s_{cst}$  and let the other part (see Figure 7(c)). With this option we can control the loss of information and the noise appearing in oSoP after multipliers, as a tradeoff.

### E. Shifts removal

As seen in the previous part, RBM scheme removes the shifts directly after multiplication in our oSoP (*i.e.* between multipliers and adders), but what if we want no shift at all in the oSoP? A good solution to this problem is to consider a single common FPR for all additions. Indeed, once we have calculated final fixed-point representation  $FPR_r = (\beta_r, \alpha_r, \gamma_r)$ , we can choose a single FPR with  $\alpha_r$  as integer part for all additions in our oSoP.

Then, in the propagation algorithm, the integer part of each adder is not determined by equation (33), but by  $\alpha_r$ . Once again, some overflows may occurs in intermediate additions, but the final result will still be valid thanks to Jackson's rule.

### F. Roundoff noise model

In fixed-point arithmetic, any positive shift can be considered as equivalent to the add of noise, called *roundoff*

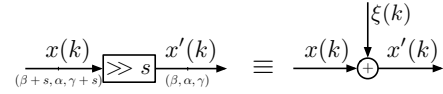


Fig. 8. Right shift can be considered as equivalent to add roundoff noise.

noise (see Figure 8). These noises are modeled as white noise uniformly distributed and statistically independent [2].

The following proposition reminds the noise produced during shift:

**Proposition 5** Let  $x(k)$  be a signal with  $(\beta + s, \alpha, \gamma + s)$  as FPR. Right shifting  $x(k)$  of  $s$  bits is similar to add to  $x(k)$  the independent white noise  $\xi(k)$ .

If  $s > 0$ , the moments of  $\xi(k)$  are given by:

$$\mu_\xi = 2^{-\gamma-1}(1 - 2^{-s}) \quad (35)$$

$$\sigma_\xi^2 = \frac{2^{-2\gamma}}{12}(1 - 2^{-2s}) \quad (36)$$

otherwise  $\xi(k)$  is null.

*Proof:* See [2], [19]. ■

**Remark 4** We consider here right-shift with troncation, but proposition 5 can also be generalized for round-to-the-nearest quantization.

So, when FPR have been propagated through the whole oSoP, all shift values are known and so we can evaluated final noise. Indeed, according to the roundoff noise model, every noises are independents, so we can sum noises through the oSoP as well as we propagate FPR.

Finally, a sum-of-products implemented in fixed-point arithmetic (with a given oSoP and implementation scheme) is equivalent to the original infinite precision sum-of-products, corrupted by the addition of a roundoff noise, that is the sum of all the roundoff noises occurring during the fixed-point evaluation.

So the moments of the final noise  $\xi(k)$  added by the implementation of the IIR in equation (3) can be known and used with proposition 2 to evaluate the roundoff noise power of a given implementation.

## IV. ILLUSTRATIVE EXAMPLE

A 3-th order Butterworth filter is used as illustrative example. Its coefficients ( $a_i$ ) and ( $b_i$ ) are given by the Matlab command `butter(3, 0.166)`:

$$\begin{aligned} b_0 &= 0.0112074993 \\ a_1 &= -1.9670683622 & b_1 &= 0.0336224979 \\ a_2 &= 1.4046495057 & b_2 &= 0.0336224979 \\ a_3 &= -0.3479211490 & b_3 &= 0.0112074993 \end{aligned} \quad (37)$$

Its  $\ell_1$ -norm is  $\|h\|_{\ell_1} = 1.21820$ .

We will consider 8 and 16-bit implementation, with different evaluation schemes, and compare them according to their respective roundoff noise error  $p$  and transfer function error  $\Delta h$ . We will also consider the theoretical latency of the

computations (assuming infinite resources), *i.e.* the height of the oSoP, given by the number of successive additions required.

The oSoPs will be displayed as a tree, with all the FPR visible (coefficients, variables and partial results), and as fixed-point algorithm (see Figures 9(a) and 9(b), and equations (38) to (39)).

We suppose that the input belongs to  $[-10, 10]$  ( $\|u\|_\infty \leq 10$ ), so with lemma 1 the output peak value is bounded by  $\|y\|_\infty \leq 12.19$ . Then  $u$  and  $y$  will have the same fixed-point representation, namely  $FPR_x = FPR_y = (8, 5, 3)$  for 8-bit implementation and  $FPR_x = FPR_y = (16, 5, 11)$  for 16-bit.

1) *Basic evaluation (oSoP #1)*: As said in introduction, classic evaluation (like those proposed by the usual tools from Mathworks and FPGA vendors) considers the same fixed-point representation for all the coefficients. In that case, the common format for all the coefficients is  $(8, 3, 5)$  and  $(16, 3, 13)$  for 8 and 16-bit implementation, respectively. One can directly note that 5 bits for the fractional part are not enough to represent the  $(b_i)$  coefficients that will be rounded to zero (underflow). So 8-bit implementation is not possible in that context (at least 11 bits are necessary to represent each coefficient without underflow).

For 16-bit implementation, the additions/accumulations will be performed on 32 bits, with  $FPR_{op} = (32, 8, 24)$ . The associated fixed-point algorithm is given by equation (38). Note that any evaluation scheme can be used (any parenthesis of the sum) can be used, since the additions are all done on the same FPR and are thus associative.

2) *Optimal oSoPs*: We also consider various optimal schemes, according to the different possible options presented in section III.

- For oSoP #2 (see Figure 9(a) and equation (38)), we consider all the possible oSoPs, with Roundoff After Multiplication and allowing bit-shifts in the computations. Among the oSoPs with the smallest roundoff error  $p$ , we consider the oSoPs with minimal height (in that case, there is only one).
- oSoP #3 is similar to oSoP #2 but with Roundoff Before Multiplication. The associated algorithm is given by eq. (40).
- We finally also look at the oSoPs involving accumulations only, with Roundoff After Multiplication. Four oSoPs satisfy the lowest roundoff error  $p$ , and they are equivalent (only permutations of multiplications with same FPR). One of them is given by (41), and denoted as oSoP #4.

All the results are summed up in table I.

We have considered a 7-term sum-of-products, with various schemes. In 16 bits, there is not much differences between the four realizations, since the final quantization (to bring back the 32-bit result in the 16-bit variable  $y(n)$ ) is the predominant source of degradation (the roundoff noise error  $p$  is not displayed here with enough precision to notice the differences). However, the differences are much more important in 8 bits. First, the basic approach is not possible, due to underflows of the coefficients. Then, one can remark that oSoP #3 produces

	8-bit		16-bit		height
	$p$	$\Delta_h$	$p$	$\Delta_h$	
oSoP #1	-	-	$1.4688e^{-4}$	$4.60875e^{-4}$	6
oSoP #2	$3.77955e^{-2}$	$3.22968e^{-4}$	$1.4688e^{-4}$	$1.98361e^{-4}$	4
oSoP #3	$1.21495e^{-2}$	$1.53966e^{-2}$	$1.4688e^{-4}$	$1.39950e^{-4}$	4
oSoP #4	$3.78499e^{-2}$	$3.22968e^{-4}$	$1.4688e^{-4}$	$1.98361e^{-4}$	6

TABLE I

ROUND-OFF NOISE POWER AND TRANSFER FUNCTION ERROR FOR THE VARIOUS IMPLEMENTATIONS

less noise than the other oSoPs, but the transfer function is much more changed. This is due to the RBM, since some of the bit-shifts are moved into the coefficients. The two others oSoP are finally quite similar, and will be chosen according to the hardware target used for the implementation.

## V. CONCLUSION

In this paper, new evaluation schemes for sum-of-products in fixed-point arithmetic have been proposed. Those schemes have been formalized according to fixed-point representation of constants and variables, and several options. Not only it automates the fixed-point code generation, but also it provides simultaneously a roundoff noise evaluation. Moreover, this result has been applied to IIR filter, for those we have described FWL effects analysis. Of course, this work presents a first step for the optimal fixed-point implementation. In hardware, some other important steps need to be performed, such as word-length optimization and hardware mapping.

It will be also interesting, in further work, to apply these results to the Specialized Implicit Framework [8], in order to generate and analyze fine fixed-point code for various filter structures. Combined to techniques such as residue feedback error [20] or  $\rho$ -operator, it will provides an interesting tool for optimized fixed-point implementation of filters.

Finally, the optimal implementation problem [4], based on the tradeoff precision vs implementation cost will be tackled.

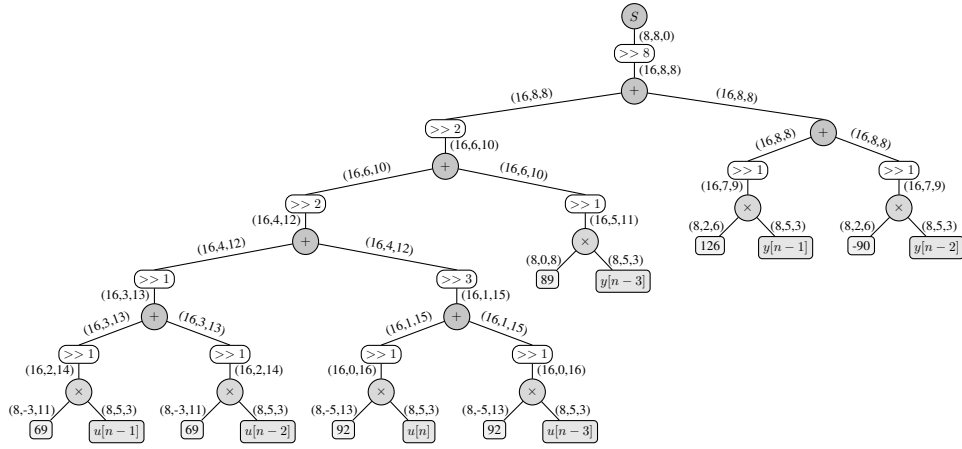
## ACKNOWLEDGMENT

This work has been sponsored by french ANR agency under grant No ANR-11-INSE-008.

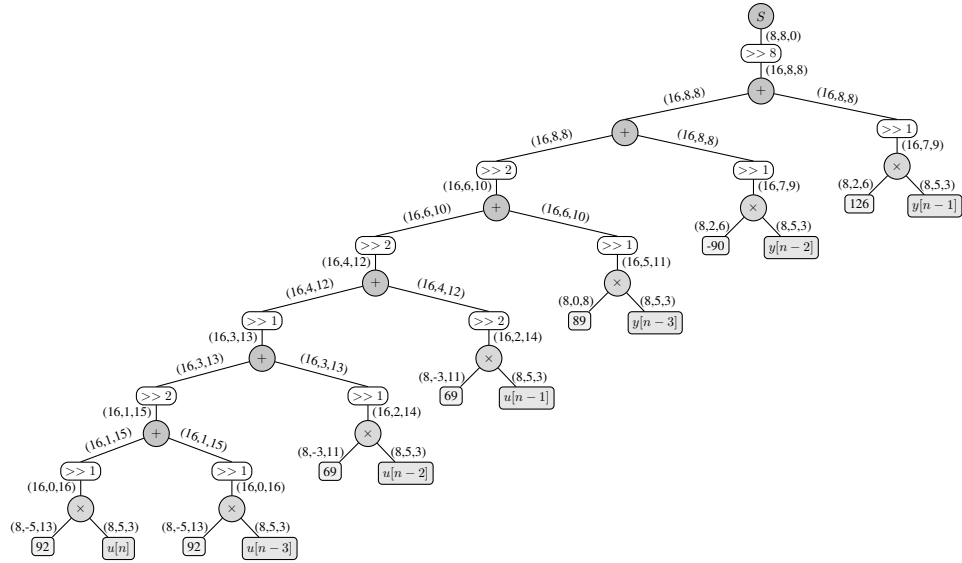
## REFERENCES

- [1] M. Gevers and G. Li, *Parametrizations in Control, Estimation and Filtering Problems*. Springer-Verlag, 1993.
- [2] B. Widrow and I. Kollár, *Quantization Noise: Roundoff Error in Digital Computation, Signal Processing, Control, and Communications*. Cambridge, UK: Cambridge University Press, 2008.
- [3] H. Hanselmann, "Implementation of digital controllers - a survey," *Automatica*, vol. 23, no. 1, pp. 7–32, January 1987.
- [4] T. Hilaire, "Towards tools and methodology for the fixed-point implementation of linear filters," in *Digital Signal Proc. Workshop and IEEE Signal Proc. Education Workshop (DSP/SPE)*, Jan. 2011, pp. 488–493.
- [5] 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.
- [6] P. Langlois, M. Martel, and L. Thévenoux, "Accuracy versus time: a case study with summation algorithms," in *Proceedings of the 4th International Workshop on Parallel and Symbolic Computation*, ser. PASCO '10. New York, NY, USA: ACM, 2010, pp. 121–130.
- [7] G. Constantinides, P. Cheung, and W. Luk, "The multiple wordlength paradigm," in *Proc. IEEE Symposium on Field-Programmable Custom Computing Machines*, 2001.





(a) 8-bit implementation of oSoP #2



(b) 8-bit implementation of oSoP #4

$$N_y(n) \leftarrow (((((((91 * N_{u(n)} + (275 * N_{u(n-1)})) + (275 * N_{u(n-2)})) + (91 * N_{u(n-3)})) + (16114 * N_{y(n-1)})) + (-11507 * N_{y(n-2)})) + (2850 * N_{y(n-3)})) \quad (38)$$

$$N_y(n) \leftarrow (((126 * N_{y(n-1)}) \ggg 1 + (-90 * N_{y(n-2)}) \ggg 1) + (((69 * N_{u(n-1)}) \ggg 1 + (69 * N_{u(n-2)}) \ggg 1) \ggg 1) \ggg 2) + ((92 * N_{u(n)} \ggg 1 + (92 * N_{u(n-3)}) \ggg 1) \ggg 3) \ggg 2 + (89 * N_{y(n-3)}) \ggg 1 \ggg 2) \quad (39)$$

$$N_y(n) \leftarrow (((16114 * N_{y(n-1)} + (-11507 * N_{y(n-2)})) + (((8814 * N_{u(n-1)} + (8814 * N_{u(n-2)})) \ggg 1) + ((11752 * N_{u(n)} + (11752 * N_{u(n-3)})) \ggg 3) \ggg 2 + (11400 * N_{y(n-3)})) \ggg 2) \quad (40)$$

$$N_y(n) \leftarrow (((((((92 * N_{u(n)} \ggg 1 + (92 * N_{u(n-3)}) \ggg 1) \ggg 2 + (69 * N_{u(n-2)}) \ggg 1) \ggg 1 + (69 * N_{u(n-1)}) \ggg 2) \ggg 2) + (89 * N_{y(n-3)}) \ggg 1) \ggg 2 + (-90 * N_{y(n-2)}) \ggg 1) + (126 * N_{y(n-1)}) \ggg 1) \quad (41)$$

- [8] 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.
- [9] R. Middleton and G. Goodwin, *Digital Control and Estimation, a unified approach*. Prentice-Hall International Editions, 1990.
- [10] 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*, vol. 51, no. 4, pp. 769–778, April 2004.
- [11] G. Li, J. Chu, and J. Wu, "A matrix factorization-based structure for digital filters," *Signal Processing, IEEE Transactions on*, vol. 55, no. 10, pp. 5108–5112, October 2007.
- [12] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*. Mc Graw Hill, 1991.
- [13] T. Hilaire, D. Ménard, and O. Sentieys, "Bit accurate roundoff noise analysis of fixed-point linear controllers," in *Computer-Aided Control Systems, 2008. CACSD 2008. IEEE International Conference on*, September 2008, pp. 607–612.
- [14] V. Tavşanoğlu and L. Thiele, "Optimal design of state-space digital filters by simultaneous minimization of sensibility and roundoff noise," in *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. CAS-31, October 1984.
- [15] 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.
- [16] V.-S. Chellaboina, W. M. Haddad, D. S. Bernstein, and D. A. Wilson, "Induced convolution norms of linear dynamical systems," in *Proc. Amer. Contr. Conf.*, June 1999, pp. 3805–3810.
- [17] "The on-line encyclopedia of integer sequences, Sequence A001147." [Online]. Available: <http://oeis.org/A001147>
- [18] L. B. Jackson, *Digital Filters and Signal Processing*, 3rd ed. Norwell, MA, USA: Kluwer Academic Publishers, 1996.
- [19] G. Constantinides, P. Cheung, and W. Luk, "Truncation Noise in Fixed-Point SFGs," *IEE Electronics Letters*, vol. 35, no. 23, Nov. 1999.
- [20] D. Williamson and S. Sridharan, "Residue feedback in digital filters using fractional feedback coefficients," in *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. 33, April 1985, pp. 477–483.