

# Formatting Bits to Better Implement Signal Processing Algorithms

Benoit Lopez<sup>1</sup>, Thibault Hilaire<sup>1</sup> and Laurent-Stéphane Didier<sup>2</sup>

<sup>1</sup>LIP6, Pierre and Marie Curie University (UPMC Univ Paris 06), Paris, France

<sup>2</sup>IMATH, University of the South, Toulon-Var (USTV), Toulon, France

**Keywords:** Fixed-point Arithmetic, Accurate Sum-of-Products, Bit Formatting, Digital Signal Processing Implementation.

**Abstract:** This article deals with the fixed-point computation of the sum-of-products, necessary for the implementation of several algorithms, including linear filters. Fixed-point arithmetic implies output errors to be controlled. So, a new method is proposed to perform accurate computation of the filter and minimize the word-lengths of the operations. This is done by removing bits from operands that don't impact the final result under a given limit. Then, the final output of linear filter is guaranteed to be a faithful rounding of the real output.

## 1 INTRODUCTION

Usually, embedded digital signal processing algorithms are specified using floating-point arithmetic and next implemented using fixed-point (FxP) arithmetic (Padgett and Anderson, 2009) for cost, size and power consumption reasons. FxP arithmetic is used as an approximation of real numbers based on integers and implicit fixed scaling by a power of 2. Of course, the quantization of coefficients and the rounding errors due to FxP computations lead to a degraded numerical accuracy of the implemented algorithm. Therefore, it is a great interest for the designer of embedded system to determine and control the implementation error while maintaining low computational effort.

In fixed-point arithmetic, a main current problem is to minimize the word-lengths of operands under constraints of precision in order to minimize area and/or power consumption (Constantinides et al., 2004). In this paper, a new method to reduce the number of bits to consider in each sum-of-products (SoP, also called Multiply-And-Accumulate) is proposed. The SoPs are one of the elementary operations of DSP algorithms. The main point of our approach is that if the final fixed-point format is known, then the bits having no impact in the final result can be detected and therefore discarded. Each term of the sum-of-products can be then reformatted into a new fixed-point format having less bits.

Some fixed-point arithmetic definitions and notations are reminded in section ?? . Section 3 formalizes

the proposed approach, which is decomposed into two formatting, for most significant bits and least significant bits, respectively. Section 4 describes the error analysis for Direct Form I filters implemented with the bit formatting technique. Finally, an illustrative example is given with a 4<sup>th</sup> order Butterworth filter, before conclusion in section 6.

## 2 FIXED-POINT ARITHMETIC AND SUM-OF-PRODUCTS

In this article we consider *signed* FxP arithmetic in two's complement representation. Let  $x$  be such a FxP number with  $w$  bits as word-length:

$$x = -2^m x_m + \sum_{i=\ell}^{m-1} 2^i x_i \quad (1)$$

where  $x_i \in \mathbb{B} \triangleq \{0, 1\}$  is the  $i^{\text{th}}$  bit of  $x$ ,  $m$  and  $\ell$  are the position of the *most* significant bit (MSB) and *least* significant bits (LSB), respectively (Fig. 1). It can be noted that  $m > \ell$  and

$$w = m - \ell + 1. \quad (2)$$

In a digital system,  $x$  is represented by an integer  $X$ , composed by the  $w$  bits  $\{x_i\}_{\ell \leq i \leq m}$ . In other words,  $X = x \cdot 2^{-\ell}$ , or equivalently

$$X = -2^{m-\ell} x_m + \sum_{i=0}^{m-\ell-1} 2^i x_{i+\ell}. \quad (3)$$

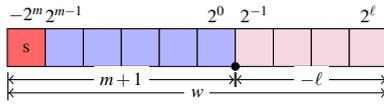


Figure 1: Fixed-point representation.  $m$  and  $\ell$  are the position of the MSB and LSB respectively (in this figure,  $m = 5$  and  $\ell = -4$ ).

Through this paper the notation  $(m, \ell)$  is used to denote the Fixed-Point Format (FPF) of such a fixed-point number and  $m$ ,  $\ell$  and  $w$  will be suffixed by the variable or constant they refer to.

**Remark 1.** In FxP arithmetic, there is no restriction on the the position of the MSB and LSB. The FPF is often chosen with  $m \geq 0$  and  $\ell \leq 0$ . FPF with  $\ell > 0$  are also possible (the quantization step is greater than 1) or  $m < 0$  (the largest represented number is lower than  $\frac{1}{2}$ ).

## 2.1 Conversion from Real to Fixed-point

Many SoP-based DSP algorithms involve real coefficients that have to be converted into FxP arithmetic. Let consider a real constant  $c \in \mathbb{R}^*$ . The position of the most significant bit  $m$  of its  $w$ -bit wide FxP representation in binary two's complement is:

$$m = \begin{cases} \lceil \log_2 |c| \rceil & \text{if } c < 0 \\ \lceil \log_2 |c| \rceil + 1 & \text{if } c > 0 \end{cases} \quad (4)$$

where  $\lfloor \cdot \rfloor$  and  $\lceil \cdot \rceil$  are the *round to the integer towards minus infinity* and *round towards plus infinity* operators, respectively.

For some very special cases, eq. (4) should be adapted (Hilaire and Lopez, 2013). The position of the least significant bit  $\ell$  is deduced from eq. (2) and the  $w$ -bit integer  $C$  representing  $c$  is computed:

$$C = \lfloor c \cdot 2^\ell \rfloor \quad (5)$$

where  $\lfloor \cdot \rfloor$  is the *round to the nearest integer* operator.

## 2.2 Sum-of-Products

In digital signal processing, the computation of filter or controller algorithms requires the evaluation of one or several SoP. Their type and number depend on the algorithm chosen (Hanselmann, 1987; Istepanian and Whidborne, 2001; Gevers and Li, 1993). For instance, the direct forms require only 1 SoP, whereas the  $n$ -th order state-space require  $n + 1$  SoPs.

The products considered in such a SoP are products of real constants and real variables. But, in the

context of fixed-point design, only fixed-point variables and fixed-point constants are considered. In this article, we consider SoPs whose constants have already been converted in FxP format.

More formally, we consider SoPs

$$s = \sum_{i=1}^n c_i \cdot v_i, \quad (6)$$

where  $\{c_i\}_{1 \leq i \leq n}$  are given non-null FxP constants and  $\{v_i\}_{1 \leq i \leq n}$  FxP variables only known to be in known intervals  $[\underline{v}_i; \bar{v}_i]$ . We focus on the best way (*i.e.* employing the minimum word-lengths) to obtain a rounding of the exact sum  $s$  at a given format.

**Remark 2.** It is also possible to consider the  $\{c_i\}$  to be real constants instead of FxP constants, so as to analyze the impact of their quantization that is not considered here.

However, this impact is well studied with sensitivity measures such as the transfer function sensitivity (Tavşanoğlu and Thiele, 1984; Gevers and Li, 1993; Hinamoto et al., 2006), the pole/zero sensitivity (Gevers and Li, 1993; Li, 1998) or IIR stability (Lu and Hinamoto, 2003).

## 3 BITS FORMATTING

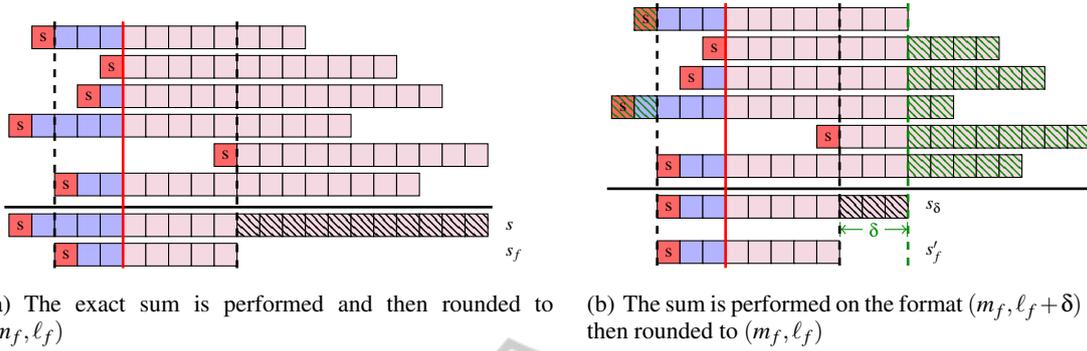
The main point of the proposed approach is that if the final fixed-point format of a sum  $s = \sum p_i$ , denoted  $FPF_f = (m_f, \ell_f)$  is known, then it is probably possible to discard some useless bits.

More formally, this paper is focused on bits of  $p_i$ s with positions lower than  $\ell_f$  (section 3.2) and greater than  $m_f$  (section 3.3) in order to determine their impact on the result. We determine the useless bits and remove them from  $p_i$ s before the sum is computed. The  $p_i$ s are rounded into an intermediate format  $(m_i, \ell_f - \delta)$ , where  $\delta$  is the number of non-useless bits with position lower than  $\ell_f$ . Then, the sum of these modified  $p_i$ s is computed and rounded into the final format  $FPF_f$  in order to obtain the final result (see Figure 2(b)).

### 3.1 Definitions and notations

In this section the fixed-point rounding modes used in this article are defined.

**Definition 1** (Fixed-point rounding modes). Let  $x$  be a real value. The notations  $\circ_d(x)$ ,  $\nabla_d(x)$  and  $\Delta_d(x)$  express the rounding to the nearest, the rounding down (*i.e.* truncation) and the rounding up of  $x$  according to the  $d^{\text{th}}$  bit, respectively. These operators



(a) The exact sum is performed and then rounded to  $(m_f, \ell_f)$  (b) The sum is performed on the format  $(m_f, \ell_f + \delta)$  and then rounded to  $(m_f, \ell_f)$

Figure 2: Two different ways to perform the FxP accumulation.

are defined by:

$$\circ_d(x) \triangleq 2^d \cdot \left\lfloor \frac{x}{2^d} \right\rfloor, \quad (7)$$

$$\nabla_d(x) \triangleq 2^d \cdot \left\lfloor \frac{x}{2^d} \right\rfloor, \quad (8)$$

$$\Delta_d(x) \triangleq 2^d \cdot \left\lceil \frac{x}{2^d} \right\rceil. \quad (9)$$

The operator  $\star_d(x)$  is the faithful rounding of  $x$  at the  $d^{\text{th}}$  bit, i.e.

$$\star_d(x) \in \{\nabla_d(x), \Delta_d(x)\}. \quad (10)$$

The round-to-the-nearest operation always returns the nearest representable point of the real exact value, while the faithful rounding operation produces either the nearest or next-nearest point.

**Notations.** Some notations need to be explicitly defined before explaining the proposed approach:

- $p_i \triangleq c_i \times v_i$  denotes the result of the fixed-point product of  $c_i$  and  $v_i$ . According to the fixed-point multiplication rule (Lopez et al., 2012), the fixed-point format of  $p_i$  is defined as  $FPF_{p_i} \triangleq (m_i, \ell_i) = (m_{c_i} + m_{v_i} + 1, \ell_{c_i} + \ell_{v_i})$ .
- $p_{i,j}$  is the  $j^{\text{th}}$  bit of  $p_i$ , for  $1 \leq i \leq n$  and  $\ell_i \leq j \leq m_i$ .
- $(M, L)$  is the FPF of the exact sum  $s = \sum_{i=1}^n p_i$ , where

$$M \triangleq \max_i(m_{p_i}) + \lceil \log_2(n) \rceil \quad (11)$$

and

$$L \triangleq \min_i(\ell_{p_i}). \quad (12)$$

$\lceil \log_2(n) \rceil$  corresponds to the number of carry bits to consider for the sum of  $n$  terms.

Moreover, three different sums are also considered, where  $\diamond_d$  is a given common rounding mode (round-to-nearest or truncate):  $\diamond_d \in \{\circ_d, \nabla_d\}$ :

- $s_f \triangleq \diamond_{\ell_f}(s)$  is the rounding of the exact sum  $s$  into the final format  $FPF_f$ .
- $s_\delta \triangleq \sum_{i=1}^n \diamond_{\ell_f - \delta}(p_i)$  is the sum of the products  $p_i$ s rounded into format  $(m_i, \ell_f - \delta)$  where  $\delta$  is a given positive constant to be discussed later.
- $s'_f \triangleq \diamond_{\ell_f}(s_\delta)$  is the rounding of the sum  $s_\delta$  into the final format  $FPF_f$ .

Figures 2(a) and 2(b) illustrate these different approximation of the sum  $s$ .

**Modular fixed-point sum.** As reminded in equation (3), a fixed-point number  $x$  is coded in computer with a  $w$ -bit signed integer  $X$ . As a consequence, all the operations are done modulo  $2^w$  on this integer. Proposition 1 specifies the modular fixed-point sum as an extension of the modular sum on integers.

**Proposition 1** (Modular fixed-point sum). *The sum modulo  $2^d$  of two fixed-point numbers  $x$  and  $y$  sharing the same FPF  $(m, \ell)$ , is noted  $x \oplus^d y$  and is computed as:*

$$x \oplus^d y \triangleq \left( \left( (X + Y + 2^{d-\ell}) \bmod 2^{d-\ell+1} \right) - 2^{d-\ell} \right) 2^\ell. \quad (13)$$

Moreover, the modular fixed-point sum of  $n$  fixed-point numbers  $x_i$  is noted and computed as:

$$\bigoplus_{1 \leq i \leq n}^d x_i \triangleq x_1 \oplus^d x_2 \oplus^d \dots \oplus^d x_n. \quad (14)$$

*Proof:* The fixed-point sum modulo  $d$ ,  $x \oplus^d y$ , corresponds to the  $d - \ell$ -bits sum of the integers  $X = x \cdot 2^{-\ell}$  and  $Y = y \cdot 2^{-\ell}$ .

Therefore, adding two signed fixed-point numbers requires to convert them in positive integers, add them modulo  $2^{d-\ell+1}$ , and convert the result back into signed fixed-point number. ■

**Example 1.** Adding 12.5 and 3.75 in FPF (4, -3) (two's complement with 8 bits) is given by  $12.5 \oplus 3.75$  and leads to -15.75 according to eq. (13) because of the overflow. 12.5 is coded by 01100.100<sub>B</sub>, 3.75 is coded by 00011.110<sub>B</sub>, so the modular sum leads to 10000.010<sub>B</sub> into format (4, -3), that is interpreted as -15.75.

### 3.2 LSBs Formatting

Let consider the final FxP format  $(m_f, \ell_f)$  of a SoP. It appears that not all the least significant bits are useful in order to correctly round the result of a SoP to  $(m_f, \ell_f)$ . The value  $\delta$  is the position such that all bits with a position lower than  $\ell_f - \delta$  are insignificant and can be removed from  $p_i$ s (Fig. 2(b)). Therefore, only  $p_i$ s such that  $\ell_i < \ell_f - \delta$  are rounded, the other remain unchanged. The sum  $s_\delta$  of rounded  $p_i$ s is computed on format  $(m_f, \ell_f - \delta)$  and finally rounded onto the final format  $(m_f, \ell_f)$ .

The following proposition formalizes the choice of  $\delta$ .

**Proposition 2.** For both rounding mode ( $\diamond_{l_f} = \circ_{l_f}$  round-to-nearest or  $\diamond_{l_f} = \nabla_{l_f}$  truncation), the integer  $\delta$  that provides  $s'_f = \star_{l_f}(s_f)$  is given by:

$$\delta = \lceil \log_2(n_f) \rceil \quad (15)$$

with  $n_f = \text{Card}(I_f)$  and  $I_f \triangleq \{i \mid \ell_i < \ell_f\}$ .

*Proof:* This proof is done for truncation round-to-nearest mode. The same reasoning can be established for round-to-nearest mode.

Computation of  $s$  involves all bits  $p_{i,j}$  whereas  $s_\delta$  requires only bits  $p_{i,j}$  for  $j \geq \ell_f - \delta$ , so:

$$s \geq s_\delta \quad (16)$$

The trivial case  $s = s_\delta$  implies  $s_f = s'_f$ , so thereafter only the case  $s - s_\delta > 0$  is considered and the difference  $s - s_\delta$  is evaluated precisely as the sum of bits  $p_{i,j}$  for  $1 \leq i \leq n$  and  $L \leq j \leq \ell_f - \delta - 1$ :

$$s - s_\delta = \sum_{i=1}^n \sum_{j=L}^{\ell_f - \delta - 1} 2^j p_{i,j} \quad (17)$$

Since  $\sum_{j=L}^{\ell_f - \delta - 1} 2^j p_{i,j} < 2^{\ell_f - \delta}$  for  $1 \leq i \leq n$ ,  $s - s_\delta$  can be bounded as follows:

$$s - s_\delta < n_f \cdot 2^{\ell_f - \delta} \quad (18)$$

Now, the difference  $s_f - s'_f$  corresponds to the rounding of the difference  $s - s_\delta$  according to the  $\ell_f$ <sup>th</sup> bit:

$$s_f - s'_f = \nabla_{\ell_f}(s - s_\delta) \quad (19)$$

It also can be viewed as the carry bits greater than  $2^{\ell_f}$  implied by the difference  $s - s_\delta$ .

Using equation (18), the difference  $s_f - s'_f$  can also be bounded:

$$\lfloor (s - s_\delta) \cdot 2^{-\ell_f} \rfloor \leq (s - s_\delta) \cdot 2^{-\ell_f} < n_f \cdot 2^{-\delta} \quad (20)$$

$$s_f - s'_f < n_f \cdot 2^{\ell_f - \delta} \quad (21)$$

Since  $\delta$  needs to be determined in order to verify  $|s_f - s'_f| \leq 2^{\ell_f}$ , the following inequality comes from equation (21):

$$n_f \cdot 2^{\ell_f - \delta} \leq 2^{\ell_f} \quad (22)$$

The smallest integer solving inequality (22) is  $\delta = \lceil \log_2(n_f) \rceil$ . ■

**Remark 3.** With Proposition 2, it may happen that one  $p_i$  (or more) has a MSB lesser than  $\ell_f - \delta$ , and so all bits of this  $p_i$  will be removed by applying this technique. Therefore, a good idea will be to redetermine  $\delta$  with  $n_f$  minus the number of removed  $p_i$ . So equation (15) can be replaced by Algorithm 1.

---

**Algorithm 1:** Evaluation of the integer  $\delta$ .

---

**Input:** Operands  $p_i$ s in format  $(m_i, \ell_i)$   
The final format  $FPF_f = (m_f, \ell_f)$

**Output:**  $\delta \in \mathbb{N}$

```

1  $n_\delta \leftarrow n$ ;
2 repeat
3    $n' \leftarrow n_\delta$ ;
4    $\delta \leftarrow \lceil \log_2(n') \rceil$ ;
5    $n_\delta \leftarrow \text{Card}(\{i \mid 1 \leq i \leq n \text{ and } m_i < \ell_f - \delta\})$ ;
6 until  $n' = n_\delta$ ;
7 return  $\delta$ 
    
```

---

**Formatting method.** The first step of LSB formatting is a direct application of proposition 2: it removes useless bits. After this step  $\ell_i \geq \ell_f - \delta, \forall 1 \leq i \leq n$ .

The second step involves having  $\ell_i = \ell_f - \delta, \forall 1 \leq i \leq n$ . To do this, either  $FPF_{p_i}$  can be changed from  $(m_i, \ell_i)$  to  $(m_i, \ell_f - \delta)$  for  $p_i$ s such that  $\ell_i > \ell_f - \delta$  (consisting to add  $\ell_i - \ell_f + \delta$  zeros to the right of these  $p_i$ s), or multipliers can be rewritten to perform operation into a given word-length. Let  $M_i$  be the multiplier computing  $p_i = c_i \times v_i$ . Then,  $w_{M_i}$ , the word-length of the result of  $M_i$  is given by:

$$w_{M_i} = m_{M_i} + \ell_{M_i} + 1 \quad (23)$$

with

$$m_{M_i} = m_{c_i} + m_{v_i} + 1 \quad (24)$$

$$\ell_{M_i} = \ell_f - \delta \quad (25)$$

where  $m_{c_i}$  and  $m_{v_i}$  are MSBs of  $c_i$  and  $v_i$  respectively.

**Remark 4.** For a better accuracy,  $m_{M_i}$  can be evaluated using formulas from section 2.1, it avoids double sign bit in general case (given by the +1 in eq. (24)). Moreover, if  $\ell_{c_i} + \ell_{v_i} > \ell_{M_i}$ , then  $\ell_{c_i} + \ell_{v_i} - \ell_{M_i}$  zeros are added to the right of  $p_i$ s to ensure  $\ell_i = \ell_f - \delta$  for these  $p_i$ s.

**Error evaluation.** Adding two numbers in FxP arithmetic requires to align them onto the same LSB using right-shifts. A rounding error may occur, which introduces a numerical error. After the second step of LSB formatting, where rounding errors may be introduced, all  $p_i$ s have the same LSB, i.e.  $\ell_f - \delta$ . There is no need of right-shift to align operands of additions and therefore no additional rounding errors are introduced by the global sum of  $p_i$ s.

The total number of right-shifts involved in the first step of our method can be bounded as follows.

**Proposition 3.** With this LSB formatting technique and for a  $n^{\text{th}}$ -order SoP, the number of right-shifts is bounded by  $n + 1$ , at most one right-shift by multiplier (denoted  $d_i$  for multiplier  $M_i$ ) and exactly one final right-shift (denoted  $d_f$ ). Their values are:

$$d_i \triangleq \ell_f - \delta - \ell_{c_i} - \ell_{v_i} \quad \forall i \in I \quad (26)$$

and

$$d_f \triangleq \delta \quad (27)$$

where  $I = \{i \mid 1 \leq i \leq n \text{ and } \ell_f - \delta > \ell_{c_i} + \ell_{v_i}\}$ .

*Proof:*  $d_i$  is the right-shift in multiplier  $M_i$  if a right-shift is necessary, i.e. if  $\ell_f - \delta > \ell_{c_i} + \ell_{v_i}$  so the number of bits to remove to ensure  $\ell_i = \ell_f - \delta$  is  $\ell_f - \delta - \ell_{c_i} - \ell_{v_i}$ . All the additions are computed on  $w_f + \delta$  bits, and the result is  $w_f$  bits long, so the final right-shift value is  $\delta$ . ■

**Remark 5.** In Proposition 3 only non-zero right-shifts are considered. If  $d_i$  is defined as  $\max(\ell_f - \delta - \ell_{c_i} - \ell_{v_i}, 0)$  rather than just  $\ell_f - \delta - \ell_{c_i} - \ell_{v_i}$ , all multipliers have a right-shift, possibly null, and so the exact number of right-shifts is  $n + 1$ .

Finally, it is possible to bound the error introduced by our method. As seen in (Hilaire and Lopez, 2013), the right shifting of  $d$  bits of a variable  $x$  (with  $(m, \ell)$  as FPF) is equivalent to add an interval error  $[e] = [e; \bar{e}]$  with

	Truncation	Round to the nearest
$[e, \bar{e}]$	$[-2^{\ell+d} + 2^\ell; 0]$	$[-2^{\ell+d-1} + 2^\ell; 2^{\ell+d-1}]$

(28)

So the global interval error for the LSB technique can be evaluated with the following properties.

**Proposition 4.** The global interval error using LSB formatting technique is  $[e] = [e; \bar{e}]$  with:

Truncation:

$$e = \sum_{i \in I} (-2^{\ell_i+d_i} + 2^{\ell_i}) - 2^{\ell_f} + 2^{\ell_f-\delta} \quad (29)$$

$$\bar{e} = 0 \quad (30)$$

Round to nearest:

$$e = \sum_{i \in I} (-2^{\ell_i+d_i-1} + 2^{\ell_i}) - 2^{\ell_f-1} + 2^{\ell_f-\delta} \quad (31)$$

$$\bar{e} = \sum_{i \in I} (2^{\ell_i+d_i-1}) + 2^{\ell_f-1} \quad (32)$$

with  $I = \{i \mid 1 \leq i \leq n \text{ and } \ell_f - \delta > \ell_{c_i} + \ell_{v_i}\}$  and  $\ell_i = \ell_{c_i} + \ell_{v_i}$ , where  $\ell_{c_i}$  and  $\ell_{v_i}$  are positions of LSBs of  $c_i$  and  $v_i$  respectively.

*Proof:* By using (28) on a multiplier  $i$ ,  $e$  equals to  $-2^{\ell_i+d_i-1} + 2^{\ell_i}$  where  $d_i$  is the right-shift value given by Proposition 3 and  $p_i$  is the initial LSB of the multiplier result, i.e. the optimal LSB which is the sum of LSBs of product operands  $c_i$  and  $v_i$ . For the final right-shift, the initial LSB equals to  $\ell_f - \delta$  and the final result is  $\delta$  bits right-shifted. ■

**Remark 6.** The precise bounds of the global interval error shown in Proposition 4 can be bounded by a power of 2. Indeed, for truncation rounding mode the global interval error is included in  $]-2^{\ell_f+1}; 0]$ , whereas for round-to-nearest rounding mode it is included in  $]-2^{\ell_f}; 2^{\ell_f}[$ .

In (Lopez et al., 2012), all  $p_i$ s have different LSBs, and therefore the global error depends on the order of the additions. Consequently, all the different evaluation schemes (ES), i.e. all the different possible orders of the additions, are generated and the choice is made mainly for ES with a minimal error. Here, all ES have the same global error value (Proposition 4), so error can not be a criteria to choose the best ES representing the sum. The criteria chosen in the section 5 is the infinite parallelism criteria, i.e. the most parallelizable ES.

### 3.3 MSBs Formatting

The MSBs of  $p_i$ s having a greater positions than the final MSB,  $m_f$  can be removed using a new formalization of the Jackson's Rule (Jackson, 1970). This Rule states that in consecutive additions and/or subtractions in two's complement arithmetic, 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.

**Example 2.** Let us consider a sum  $S$  of three 8-bit integers with two's complement arithmetic, for example  $104 + 82 - 94$ . The result  $S = 92$  is in the range of 8-bit signed numbers, but the intermediate sum  $104 + 82$

produces an overflow and equals to  $-70$  into this format (instead of 186 that cannot be represented). The final sum  $-70 - 94$  also produces an overflow and equals to 92 into the final format, that is the correct result.

With this paper's notations, it means that bits with a greater position than  $m_f$  can be removed from concerned  $p_i$ s.

**Proposition 5** (Fixed-Point Jackson's Rule). *Let  $s$  be a sum of  $n$  fixed-point number  $p_i$ s, in format  $(M, L)$ . If  $s$  is known to have a final MSB equals to  $m_f$  with  $m_f < M$ , then:*

$$s = \bigoplus_{1 \leq i \leq n} \left( \sum_{j=L}^{m_f+1} 2^j p_{i,j} \right) \quad (33)$$

*Proof:*  $s = \sum_{i=1}^n p_i$ , so, from (1):

$$s = \sum_{i=1}^n \left( -2^M p_{i,M} + \sum_{j=L}^{M-1} 2^j p_{i,j} \right) \quad (34)$$

All bits of  $s$  greater than  $2^{m_f}$  (from  $p_{i,j}$  with  $j \geq m_f + 1$  and from the carry bits produced by  $p_{i,j}$  with  $j < m_f + 1$ ) are repetitions of the sign bits, since  $-2^{m_f} \leq s < 2^{m_f}$  (by definition of the final FPF). Therefore,  $s$  can be only computed with  $p_{i,j}$  with  $j < m_f + 1$  in the format  $FPF_f$ . ■

Thus, in our method, the MSB formatting is an application of the proposition to a sum  $s_\delta$  previously *LSB-formatted*  $p_i$ s, i.e. with  $L = \ell_f - \delta$ . Therefore,  $s_\delta$  can be computed only using bits  $p_{i,j}$  with  $1 \leq i \leq n$ ,  $\ell_f - \delta \leq j \leq m_f$  without considering intermediate overflows.

## 4 OUTPUT ERROR ANALYSIS

Let us consider a  $n$ -th order IIR<sup>1</sup> filter having  $H$  as a transfer function:

$$H(z) = \frac{b_0 + b_1 z^{-1} + \dots + b_n z^{-n}}{1 + a_1 z^{-1} + \dots + a_n z^{-n}}, \quad \forall z \in \mathbb{C}. \quad (35)$$

This filter is usually realized with the following algorithm

$$y(k) = \sum_{i=0}^n b_i u(k-i) - \sum_{i=1}^n a_i y(k-i) \quad (36)$$

where  $u(k)$  is the input at step  $k$  and  $y(k)$  the output at step  $k$ .

So the evaluation of the filter relies on the evaluation of a SoP. As seen in previous sections, the fixed-point evaluation of eq. (36) implies the add of an error

<sup>1</sup>Infinite Impulse Response

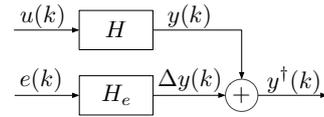


Figure 3: Equivalent system, with output error extracted.

$e(k)$  at time  $k$ , and only  $y^\dagger$  (the output contaminated with roundoff error) can be computed:

$$y^\dagger(k) = \sum_{i=0}^n b_i u(k-i) - \sum_{i=1}^n a_i y^\dagger(k-i) + e(k). \quad (37)$$

In (Lopez et al., 2012), it has been shown that the implemented system eq. (37) can be seen as the initial system (36) with an error added on the output, as shown in Figure 3: by subtracting equations (37) and (36), it comes

$$y^\dagger(k) - y(k) = e(k) - \sum_{i=1}^n a_i (y^\dagger(k-i) - y(k-i)) \quad (38)$$

So the output error  $\Delta y(k) \triangleq y^\dagger(k) - y(k)$  can be seen as the result of the error  $e(k)$  through the filter  $H_e$  defined by

$$H_e(z) = \frac{1}{1 + a_1 z^{-1} + \dots + a_n z^{-n}}, \quad \forall z \in \mathbb{C}. \quad (39)$$

Since the error  $e(k)$  done in the evaluation of the SoP is known to be in a given interval  $[e; \bar{e}]$  (see Proposition 4), then the following proposition (Hilaire and Lopez, 2013) gives the output error interval:

**Proposition 6** (Output error interval).  *$\Delta y(k)$  is the output of the error  $e(k)$  through the filter  $H_e$ . If the error  $e(k)$  is in  $[e; \bar{e}]$ , then  $\Delta y(k)$  is in  $[\underline{\Delta y}; \overline{\Delta y}]$  with:*

$$\underline{\Delta y} = \frac{\bar{e} + e}{2} |H_e|_{DC} - \frac{\bar{e} - e}{2} \|H_e\|_{\ell^\infty} \quad (40)$$

$$\overline{\Delta y} = \frac{\bar{e} + e}{2} |H_e|_{DC} + \frac{\bar{e} - e}{2} \|H_e\|_{\ell^\infty} \quad (41)$$

and  $|H_e|_{DC}$  is the DC-gain (low-frequency gain) of  $H_e$  and  $\|H_e\|_{\ell^\infty}$  its worst-case peak gain:

$$\|H_e\|_{\ell^\infty} \triangleq \frac{\sup_{k \geq 0} |y(k)|}{\sup_{k \geq 0} |u(k)|} \quad \forall u \text{ and } y \text{ input and output of } H_e. \quad (42)$$

They can be computed by:

$$|H_e|_{DC} = H_e(1), \quad \|H_e\|_{\ell^\infty} = \sum_{k \geq 0} |h_e(k)| \quad (43)$$

where  $h_e(k)$  is impulse response of the filter  $H_e$ .

*Proof:* Since  $H_e$  is linear,  $\Delta y(k)$  can be seen as the sum of a constant term  $\frac{\bar{e}+e}{2}$  through the filter  $H_e$  and a variable term bounded by  $\frac{\bar{e}-e}{2}$ . The constant term is amplified by the low-frequency gain  $|H_e|_{DC}$ , whereas the bound of the variable term is amplified by  $\|H_e\|_{\ell^\infty}$  (eq. (42)). ■

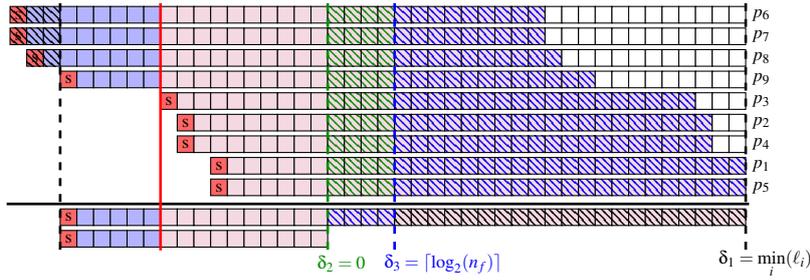


Figure 4: Bits representation of the sum of the example.

## 5 RESULTS AND COMPARISONS

A 4-th order Butterworth filter is used as an illustrative example. The chosen realization to compute this filter is the Direct Form I:

$$y(k) = \sum_{i=0}^4 b_i u(k-i) - \sum_{i=1}^4 a_i y(k-i). \quad (44)$$

Its coefficients  $a_i$ s and  $b_i$ s are given by the Matlab command `butter(4, 0.136)`:

$$\begin{array}{ll} b_0 = 0.001328017792779 & \\ a_1 = -2.871116228316502 & b_1 = 0.005312071171115 \\ a_2 = 3.208250066295749 & b_2 = 0.007968106756673 \\ a_3 = -1.634594881084453 & b_3 = 0.005312071171115 \\ a_4 = 0.318709327789667 & b_4 = 0.001328017792779 \end{array}$$

For implementations, the output variables  $y(i)$ s, input variables  $u(i)$ s, constants  $a_i$ s and  $b_i$ s are 16-bit words.

Moreover, variables  $u(i)$ s are considered in this example to be in the interval  $[-13; 13]$  (the corresponding FPF is  $(4, -11)$ ), and variables  $y(i)$ s (including result output  $y(k)$ ) are known to be in the interval  $[-17.123541221107534; 17.123541221107534]$  (corresponding to the final FPF  $(m_f, \ell_f) = (5, -10)$ ).

From these informations, the operands to be summed,  $p_i$ s, can be obtained with their respective FPF:

$$p_i \triangleq \begin{cases} b_{i-1}u(k-(i-1)) & \text{if } 1 \leq i \leq 5 \\ a_{i-5}y(k-(i-5)) & \text{if } 6 \leq i \leq 9 \end{cases}$$

$$\begin{array}{ll} FPF_{p_1} = (-4, -35) & \\ FPF_{p_2} = (-2, -33) & FPF_{p_6} = (8, -23) \\ FPF_{p_3} = (-1, -32) & FPF_{p_7} = (8, -23) \\ FPF_{p_4} = (-2, -33) & FPF_{p_8} = (7, -24) \\ FPF_{p_5} = (-4, -35) & FPF_{p_9} = (5, -26) \end{array}$$

Four implementations are compared: a double precision implementation and three fixed-point implementations using bit formatting approach using different values of  $\delta$ . In the first FxP implementation (denoted  $Fix_1$ ), all bits are considered. This means

that  $\delta_1$  is chosen such that  $\ell_f - \delta_1 = \min_i(\ell_i)$ . In other word, we have  $\delta_1 = \ell_f - \min_i(\ell_i) = 25$ . The  $Fix_1$  implementation corresponds to the computation of the sum  $s$  with no LSB formatting.

In the second FxP implementation  $Fix_2$ , no additional bits are considered. The intermediate format is the final format,  $(m_f, \ell_f)$ , which corresponds to  $\delta_2 = 0$ . A large LSB reduction is performed.

The third FxP implementation  $Fix_3$  is the faithful implementation, with  $\delta_3$  determined from Proposition 2, i.e.  $\delta_3 = \lceil \log_2(n_f) \rceil$  with  $n_f = n = 9$ , so  $\delta_3 = 4$ . Only 4 guards bits are used in the LSB formatting.

Figure 4 illustrates this example. For  $Fix_1$ , since the intermediate format is  $(5, -35)$ , additional bits equal to 0 are considered to align all  $p_i$ s onto this format. For  $Fix_2$  and  $Fix_3$ , the intermediate formats,  $(5, -10)$  and  $(5, -14)$  respectively, permit to remove bits, blue and green hatched bits and blue hatched bits respectively. Finally, the intermediate sums are rounded to the final format  $(5, -10)$ , except  $s_{\delta_2}$  which is already in the final format.

The global interval error  $[\underline{e}; \bar{e}]$  is computed, for the implementation  $Fix_3$ , using Proposition 4 :

$$\underline{e} = -1.4645302 \times 10^{-3}, \quad \bar{e} = 0. \quad (45)$$

From equation (43), DC-gain and worst-case peak gain of  $H_e$  are obtained :

$$|H_e|_{DC} = 49.5647, \quad \|H_e\|_{\ell^\infty} = 66.8474. \quad (46)$$

Finally, the output error interval (Proposition 6)  $[\underline{\Delta y}; \bar{\Delta y}]$  is computed from equations (45) and (46) :

$$\underline{\Delta y} = -8.52445240 \times 10^{-2}, \quad \bar{\Delta y} = 1.26555189 \times 10^{-2}. \quad (47)$$

As illustration (but not proof), of the theoretical result (46), a simulation in FxP and floating point arithmetic has been done with a white noise input  $u(k)$  in  $[-13; 13]$ . The error between the double floating result and each of the FxP implementations is shown in Figure 5.

The number of additional bits considered in  $Fix_3$  is small compared with  $Fix_1$  which considers all bits,

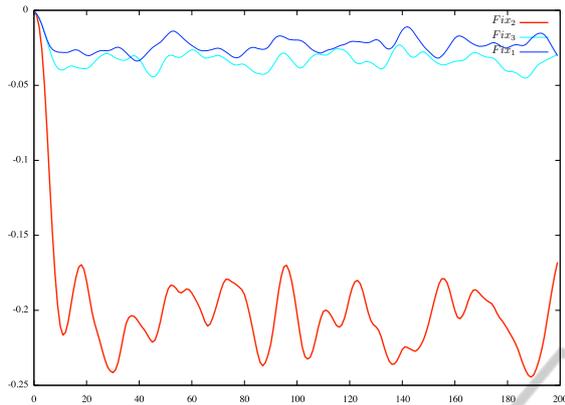


Figure 5: Computed error between double implementation and the three different fixed-point implementations.

but it is good enough to have errors measures far better than  $Fix_2$  and really close to  $Fix_1$ . The plotted error for implementation  $Fix_3$  is in the bound predicted by the theory.

The fixed-point implementation  $Fix_3$  is given by algorithm 2. In this implementation, sum modulo  $2^{m_f+1}$  (i.e.  $2^6$ ) is performed, but since algorithm considers integer computations, the sum is performed modulo  $2^{m_f+1-\ell_f+\delta}$  (i.e.  $2^{20}$ ).

**Algorithm 2:** Fixed-point algorithm.

<b>Input:</b>	$R2 \leftarrow (-20887 * Y4) \gg 12;$
$U0$ to $U4$ : 16-bit input (4, -11)	$R3 \leftarrow R0 \oplus R2;$
$Y1$ to $Y4$ : 16-bit input (5, -10)	$R0 \leftarrow R1 \oplus R3;$
<b>Output:</b> $Y$ : 16-bit output (5, -10)	$R1 \leftarrow (22280 * U4) \gg 21;$
<b>Data:</b> $Rx$ : 20-bit registers	$R2 \leftarrow (26781 * Y3) \gg 10;$
$\oplus$ : the 20-bit sum	$R3 \leftarrow R1 \oplus R2;$
$R0 \leftarrow (23520 * Y1) \gg 9;$	$R1 \leftarrow (16710 * U2) \gg 18;$
$R1 \leftarrow (-26282 * Y2) \gg 9;$	$R2 \leftarrow R3 \oplus R1;$
$R2 \leftarrow R0 \oplus R1;$	$R1 \leftarrow (22280 * U1) \gg 19;$
$R0 \leftarrow (22280 * U0) \gg 21;$	$R3 \leftarrow R2 \oplus R1;$
$R1 \leftarrow R0 \oplus R2;$	$R1 \leftarrow R0 \oplus R3;$
$R0 \leftarrow (22280 * U3) \gg 19;$	// Output computation
	$Y \leftarrow R1 \gg 4;$

## 6 CONCLUSIONS

Throughout this paper, a new method of formatting bits has been described, in order to design fixed-point sum-of-products and then linear filters. This method allows to remove some bits and keep only the bits that impact the final result. The computed result is a faithful rounding of the final result considering all the bits. The example has shown the utility of applying this method to a linear filter expressed in a very common form, and the gain in term of number of bits is significant.

Future work will consist of a word-length optimization step that will consider the bit formatting

method, and a code generation for algorithm-to-code mapping.

## ACKNOWLEDGEMENTS

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

The authors would like to thank Florent de Dinechin for the instructive discussions about fixed-point implementation.

## REFERENCES

- Constantinides, G., Cheung, P., and Luk, W. (2004). *Synthesis and Optimization of DSP Algorithms*. Kluwer Academic Publishers.
- Gevers, M. and Li, G. (1993). *Parametrizations in Control, Estimation and Filtering Problems*. Springer-Verlag.
- Hanselmann, H. (1987). Implementation of digital controllers - a survey. *Automatica*, 23(1):7–32.
- Hilaire, T. and Lopez, B. (2013). Reliable implementation of linear filters with fixed-point arithmetic. In *Proc. IEEE Workshop on Signal Processing Systems (SiPS)*.
- Hinamoto, T., Omoifo, O., and Lu, W.-S. (2006). L2-sensitivity minimization for mimo linear discrete-time systems subject to l2-scaling constraints. In *Proc. IS-CCSP 2006*.
- Istepanian, R. and Whidborne, J., editors (2001). *Digital Controller implementation and fragility*. Springer.
- Jackson, L. (1970). Roundoff-noise analysis for fixed-point digital filters realized in cascade or parallel form. *Audio and Electroacoustics, IEEE Transactions on*, 18(2):107–122.
- Li, G. (1998). On the structure of digital controllers with finite word length consideration. *IEEE Trans. on Autom. Control*, 43(5):689–693.
- Lopez, B., Hilaire, T., and Didier, L.-S. (2012). Sum-of-products Evaluation Schemes with Fixed-Point arithmetic, and their application to IIR filter implementation. In *Conference on Design and Architectures for Signal and Image Processing (DASIP)*.
- Lu, W.-S. and Hinamoto, T. (2003). Optimal design of iir digital filters with robust stability using conic-quadratic-programming updates. In *IEEE Trans. Signal Processing*, volume 51, pages 1581–1592.
- Padgett, W. T. and Anderson, D. V. (2009). Fixed-point signal processing. *Synthesis Lectures on Signal Processing*, 4(1):1–133.
- Tavşanoğlu, V. and Thiele, L. (1984). Optimal design of state-space digital filters by simultaneous minimization of sensibility and roundoff noise. In *IEEE Trans. on Acoustics, Speech and Signal Processing*, volume CAS-31.