



**HAL**  
open science

## Evaluating the hardware cost of posit arithmetic

Luc Forget, Yohann Uguen, Florent de Dinechin

► **To cite this version:**

Luc Forget, Yohann Uguen, Florent de Dinechin. Evaluating the hardware cost of posit arithmetic. 2021. hal-03195756v1

**HAL Id: hal-03195756**

**<https://hal.science/hal-03195756v1>**

Preprint submitted on 12 Apr 2021 (v1), last revised 16 Apr 2021 (v3)

**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.

# Evaluating the hardware cost of posit arithmetic

Luc Forget, Yohann Uguen, Florent de Dinechin  
 Univ Lyon, INSA Lyon, Inria, CITI  
 F-69621 Villeurbanne, France  
 first-name.last-name@insa-lyon.fr

**Abstract**—The posit number system is an elegant encoding of floating-point values, recently proposed as a drop-in replacement for the IEEE-754 floating-point standard. It is more regular and simpler than IEEE-754, discarding many of its less used features, such as subnormals, Not-a-Numbers, signed infinities, or directed rounding modes. Posit arithmetic consistently rounds to nearest even, with overflows managed by saturation. It also offers tapered precision, where numbers near 1 have more significant bits than very large or very small ones (and more significant bits than IEEE-754 for the same word size). For applications whose data is consistently within this accurate region, posits offer more accuracy for the same word size, while accuracy may be degraded for other applications.

A common claim of the posit literature is that posits, being simpler than IEEE-754, are also cheaper to implement. This article refutes this claim. In the same context, posit operators are shown to be intrinsically both larger and slower than IEEE-754 operators of the same size. A first reason is that posits must be decoded into some kind of floating-point before being processed, then they must be reencoded. This article develops a solution where posits are kept decoded in processor registers to avoid paying the decoding/encoding overhead for each operation, using an original in-register rounding approach to avoid double rounding issues. A second reason is that the extra accuracy of posits requires a wider data-path. For instance, the processor registers must be wider than posits themselves if the previous solution is used. Compared to these two overheads, the overhead of managing directed rounding modes, infinities or NaNs in IEEE-754 operators is very small.

Unbiased quantitative comparisons support this analytical study. For this purpose, this work introduces an open-source library of basic parametrized operators (for addition/subtraction, multiplication, and exact accumulation) for posit and IEEE-754. All the operators are developed with the same high design effort, using the same tools and libraries, and are fully compliant to their respective standards. This library is first shown to improve the state of the art of posit operators. Despite this, a comparison on FPGAs shows that posits arithmetic has an overhead between 30% and 60% both in area and delay compared to IEEE-754 for the same data size.

## I. INTRODUCTION

The set of real numbers is infinite and uncountable. Digital computers operate on a finite set of values, the machine codes. Despite the tremendous pace of progress in computer technology, compromises must be made when computing with real numbers in a computer. The most used solution is to define a machine code size  $N$  (a number of bits), then select a finite subset of the real numbers of at most  $2^N$  elements, which can be encoded on the  $2^N$  available binary codes. Computations are then performed within this subset, therefore both the subset and its encoding must be carefully selected so that the computations can be efficiently implemented out

of the machine codes. When the exact result of a computation does not itself belong to the representable subset, a convention (rounding and overflow management) determines what to do.

The mainstream way to represent real numbers in general purpose computers is floating-point. The IEEE-754 standard [1] defines floating-point number sets and encoding schemes for  $N \in \{16, 32, 64, 128\}$ , and a versatile set of rounding and overflow conventions.

A recently introduced alternative to this standard, for comparable applications, is the posit encoding scheme [2], [3]. This format is presented as a drop-in replacement of the IEEE-754 formats, while providing better performance and accuracy for the same  $N$  using a more efficient encoding scheme [2]. These two competing ways of representing numbers are compared in Section II.

Many works have compared the accuracies of posits and IEEE floats [4], [5], [6], [7]. The focus of the present article is the comparison of the hardware implementation cost of basic arithmetic operators for the two encoding schemes.

The method for this comparison has been to develop an open-source<sup>1</sup>, fully parametric library of hardware operators called MARTo that covers both IEEE-754 operators and posits of arbitrary sizes, such that posits can effectively be used as a drop-in replacement for IEEE floats. It is hoped that this tool will allow the community to assess their relative speed and cost, just like software libraries such as SoftFloat and SoftPosit have enabled like-for-like comparison in accuracy.

To make this comparison as unbiased as possible, all the operators are developed with similar design effort using the same core library of fixed-point components. They will be evaluated using the same tools in the same context. The common design goals are 1/ full standard compliance, 2/ combinatorial designs that can be pipelined for higher frequency, and 3/ area-oriented designs.

Classical techniques that improve latency at the expense of area, such as dual-path floating-point addition [8] or hardware speculation [9] are not considered here – as the reader will see, they can benefit posits as well as IEEE. However, the datapath sizes are carefully minimized. An example of such a design effort is the proposed posit adder that has a narrower data-path than the state of the art. This is explained by the fact that the “far” and “near” cases are exclusive in a floating-point addition.

The IEEE-754 library of MARTo is a competent reimplementation of the state of the art and does not need to be

<sup>1</sup><https://gitlab.inria.fr/lforget/marto>

described in detail. Conversely, posit literature is younger, and this paper discusses in detail posit hardware support and the choices made in MARTo.

Posit really is a clever encoding of a selected subset of a larger set of floating-point numbers. To compute on posits, one must first decode them into this set [10], [11], [12]. This set is defined formally in Section II, which also discusses an efficient encoding of this set called the Posit Intermediate Format (PIF).

Section III defines two alternative approaches that can be used in a Posit Arithmetic Unit (PAU): one that does the posit encoding and decoding for each operation, and one where the data kept in posit registers is actually decoded PIF data. This second approach reduces the latency of the operations, but requires larger registers than the first. It also requires more complex rounding hardware in order to be bit-for-bit compatible with the first.

Section IV describes in detail all the hardware blocks that constitute a PAU in each of these alternatives. Section V discusses the implementation of the quire, a posit version of Kulisch's exact floating-point accumulator. The latter is also proposed in MARTo for comparison purpose.

Finally, Section VI evaluates and compares the costs and delay of posit and IEEE-754 operators on FPGA hardware. To ensure an unbiased comparison, we first show that the proposed posit implementation improves the state of the art, while the IEEE-754 operators compare favorably to it. The comparisons then suggest that replacing IEEE with posit is costly, both in terms of hardware and in terms of latency. For instance, for both addition and subtraction, the area-delay product of 32-bit posits is worse by a factor two compared to IEEE 754-compliant 32-bit floats with full hardware subnormal support. This is to contrast with the extra representation accuracy offered by posits: at best 4 bits, or a 17% improvement, on a small domain, all the more as posits actually degrade accuracy out of this domain [6].

## II. BINARY FLOATING-POINT NUMBERS, AND THEIR IEEE AND POSIT ENCODING

### A. Precision and range of sets of floating-point numbers

A set of binary floating-point numbers is characterized by three integers [8]:

- a precision  $P \geq 2$ ,
- a minimal and maximal exponent  $E_{\min}$  and  $E_{\max}$  respectively

The precision is the maximum number of significant bits this set allows.

A real value  $x$  belongs to this set if there exists at least one couple of integers  $(m, e)$  such that

$$x = m \cdot 2^{e-P+1} \quad (1)$$

where  $E_{\min} \leq e \leq E_{\max}$  and  $m < 2^P$ .

The value  $m$  is the integral significand of the representation of  $x$ , and  $e$  its exponent.

To manipulate floating-point values on a digital computer, an associated binary encoding must be specified. Two encoding schemes are described in the following sections.

### B. The IEEE-754 binary encoding

An IEEE-754 binary encoding scheme is defined by two positive integers:

- $W_e$ , the exponent field width and
- $W_f$ , the fraction field width.

A value is represented by a bit vector of size  $1 + W_e + W_f$ .

The first bit of this vector is the sign bit  $s$ , then follow the exponent field  $e$  of size  $W_e$  and the fraction field  $f$  of size  $W_f$ .

The way a given code should be decoded depends on the value of the exponent field.

1) *General case: decoding normal values:* The general case is when bits of  $e$  are not all ones or all zeros. In this case, the encoding represents the value

$$x = (-1)^s \times 1.f \times 2^{e-E_{\max}} \quad (2)$$

with  $e$  interpreted as a positive integer and  $E_{\max}$  the maximal exponent of the encoding scheme

$$E_{\max} = 2^{W_e-1} - 1 \quad (3)$$

The '1' before the fractional point is used to avoid duplicated representation of the same value. Indeed, without such a constraint on the significand, a given number might have many representations (e.g.  $0.0101 \times 2^0 = 0.1010 \times 2^{-1}$ ). As this 1 is always set, it is not necessary to store it in the encoding.

The general case allow the representation of every floating-point number of precision  $W_f + 1$  with minimal exponent  $E_{\min} = 2^{W_e} + 2$  and maximal exponent  $E_{\max}$ , excepted those of the form  $m \times 2^{E_{\min}-W_f}$  with  $m < 2^{W_f}$ . This last set of values is referred as the subnormals.

2) *The subnormal case:* Subnormal value representations have all the exponent bits set to zero.

In this case, the implicit leading 1 is replaced by a leading zero, and the value represented is interpreted as

$$x = (-1)^s \times 0.f \times 2^{E_{\min}} \quad (4)$$

Zero is encoded as a subnormal number with all fraction bits set to zero. Two encodings exist for zero differing, on the sign bit value.

3) *Non-numerical values:* In addition to numerical values, the IEEE-754 binary encoding reserves machine codes for special non-numerical values.

Such values are identified with all the exponent field bits set to one. The represented value can be either infinity if all fraction bits are zeros, or Not a Number (NaN) if at least one fraction bit is one.

Not a Number represents the output of an illegal arithmetic operation, such as a division by zero.

The infinity value is an overflow marker: it is the result of an operation whose exact result is greater than the greatest representable numerical value of the encoding.

### C. Posits

The posit number system [2] also encodes floating-point values. A posit format is also defined by two positive integers:

- the word size  $N$ ,
- the exponent size,  $ES$

A value  $v$  is represented by a bit vector of  $N$  bits, starting with the sign bit  $s$ . Next comes the Regime, a variable length field encoding the coarse grain exponent using a thermometer encoding. The regime end is marked by an alternating bit or the end of the word. Then follows the exponent scale field,  $e_s$ , of at most  $ES$  bits. The remaining bits constitute the fraction part. Figures 1 and 2 present the decomposition of posit codes. They will be used as examples for the description of the encoding.

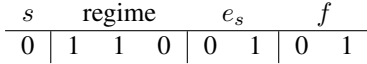


Fig. 1: Posit decomposition example ( $N = 8$ ,  $ES = 2$ ).

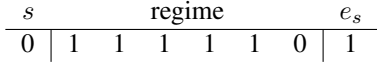


Fig. 2: Posit example with a long regime ( $N = 8$ ,  $ES = 2$ ).

1) *The posit encoding*: The posit numeric values are encoded in 2's complement. When the sign bit is zero, the encoded value is positive.

The exponent  $e$  of the represented value is split in two parts. The most significant part is computed out of the regime field. This field is constituted by a sequence of  $l$  identical bits  $b$  followed by an opposite bit or the end of the word. The encoded range  $e_h$  is  $-l$  if the bits of this sequence are equal to  $s$ , otherwise  $l - 1$ . In the encoding of Figure 1, the sequence consists in two ones:  $l = 2$ , therefore  $e_h = 1$ . The  $ES$  following bits are xored with  $s$  to obtain the lower exponents bits  $e_l$ : the exponent  $E$  is the concatenation of  $k$  and  $e_s$ . In other words,

$$e = e_h \cdot 2^{ES} + e_l \quad (5)$$

In our example,  $e = 101$  as  $e_s = 01$ .

The remaining bits encode the fractional part  $f$  of the significand. An implicit leading bit  $i$  is obtained by negating  $s$ , here  $i = 1$ . Finally, the value of the posit can be defined as:

$$2^e \times (i \cdot f - 2s). \quad (6)$$

The value represented by the example is

$$2^{101_2} \times (1.01_2 - 2 \times 0) = 2^5 \times 1.25 = 40.$$

Note that the regime can extend to the point where there is no room for  $f$  or  $e_s$ . In this case, the shifted out bits are assumed to be zeros. This is what happens for the value encoded in Figure 2. Here, the length of the regime  $l$  is 5 and it is constituted by ones so the regime is  $e_h = 4_{10} = 100_2$ . The exponent shift value is  $e_l = 10_2$  and there is no fraction bit. The represented value is then

$$2^{10010_2} \times (1.0 - 2 \times 0) = 1 \times 2^{10} = 1024$$

TABLE I: Parameters of standard posit formats.

N	Standard posit		smallest FP superset		quire
	$ES$	$E_{\max}$	$W_e$	$W_f$	$W_q$
8	0	6	4	5	32
16	1	28	6	12	128
32	2	120	8	27	512
64	3	496	10	58	2048

Posit formats admit two special values, 0 and Not a Real (NaR). For encoding 0, all the posit fields are null, including the implicit bit. NaR is the equivalent of IEEE-754 NaN (Not a Number). Its encoding only has the sign bit set. There is no special encoding to catch overflows: posit arithmetic saturates instead.

2) *Posit smallest floating-point superset*: Due to the run length encoding of the range, posit values with low magnitude exponents have a greater precision than high magnitude exponent values. As the hardware bit-widths cannot change dynamically to adapt to varying precision inputs, posit values are handled as fixed precision floating-point values. Parameters for the smallest fixed precision floating-point superset are derived as follows.

The maximum precision  $W_F$  of posit values is obtained for the minimum length of the regime (2), therefore

$$W_f = N - (3 + ES) \quad (7)$$

On the other hand, maximal exponent is obtained when the regime running length is  $N - 1$ . In this case, all the  $e_s$  and  $f$  bits are pushed out by the regime. Hence, the maximum exponent value is

$$E_{\max} = (N - 2)2^{ES} \quad (8)$$

As the opposite exponent can also be reached, the number of bits needed to store the exponent is

$$W_e = 1 + \lceil \log_2((N - 2)2^{ES}) \rceil \quad (9)$$

$$= 1 + ES + \lceil \log_2(N - 2) \rceil \quad (10)$$

The  $ES$  parameter allows trading between the range of the format and its maximal precision. The posit draft standard [3] defines four formats with an encoding size  $N$  of 8, 16, 32 and 64 respectively, such that

$$ES = \log_2(N) - 3 \text{ for standard posits.} \quad (11)$$

These formats are used for evaluation in this paper, although the MARTo library is fully parameterized in  $N$  and  $ES$ .

Table I gives, for each of the standard posit formats, the exponent and fraction sizes of the smallest floating-point superset.

A posit-compliant environment must also provide a *quire*. This large fixed point vector allows for the exact accumulation of posit products. It is based on the floating-point Kulisch accumulator [13]. Non-zero product exponents range from

$$P_{\min} = -2 \times E_{\max} \quad (12)$$

to

$$P_{\max} = 2 \times E_{\max} \quad (13)$$

The quire is a very large fixed-point number spreading over all the product exponent ranges. A positive number of carry guard bits  $C$  can be added to allow the sum of up to  $2^C$  maximal magnitude products before an overflow occurs.

The weight of the associated quire most significant bit is then

$$Q_{\text{msb}} = P_{\text{max}} + C \quad (14)$$

while the weight of its least significant bit is

$$Q_{\text{lsb}} = P_{\text{min}} \quad (15)$$

The width of such a quire is

$$W_q = Q_{\text{msb}} - Q_{\text{lsb}} + 1 \quad (16)$$

From (8) and (11), one can see that for standard posit formats, product exponents range from  $-\frac{N^2-2N}{4}$  to  $\frac{N^2-2N}{4}$ . Hence, without carry guard bits the quire width would be  $W_q = \frac{N^2}{2} - N + 1$ . The standard motivates that the quire should easily be transferred to and from memory. To do so, it should have a size which is a multiple of 8. With the sign bit and the addition of  $N - 2$  carry guard bits, this goal is attained. Hence the width of the quire is

$$W_q = \frac{N^2}{2} \text{ for standard posits.} \quad (17)$$

and the number of carry guard bits is

$$C = N - 2 \text{ for standard posits.} \quad (18)$$

3) *Posit Intermediate Format: a hardware-friendly posit recoding*: With this smallest posit FP superset, it becomes possible to define a new intermediate encoding for numbers in this set that is more adapted to hardware arithmetic operations, in the sense that all its fields have a fixed width. In this work, this encoding scheme is denoted Posit Intermediate Format or PIF.

Figure 3 highlights the role of this format in an end-to-end posit arithmetic operator. Posits are first decoded to PIF. As PIF can represent all posit values, this operation is exact. Then, the arithmetic operation is performed on PIF data. Finally, the result is encoded back to posit. Since rounding (to an exponent-dependent position) must be performed in this encoding step, the output format of the PIF operation must be an Unrounded PIF, which is a PIF extended with all the additional information needed for standard-compliant rounding. The UPIF format is detailed below.

We believe this 3-step approach is inevitable for stand-alone posit operators (except for very small formats where a simple tabulation may be used). It is followed (more or less explicitly) by leading hardware posit implementations [10], [11], [12].

The PIF should then be designed with two objectives in mind:

- Posit to PIF conversion should be as simple as possible,
- Arithmetic operations should be efficiently computed on this representation.

Because of the second objective, PIF is a direct floating-point representation that uses the parameters of Table I. Concerning the first objective, the proposed PIF encodes both the

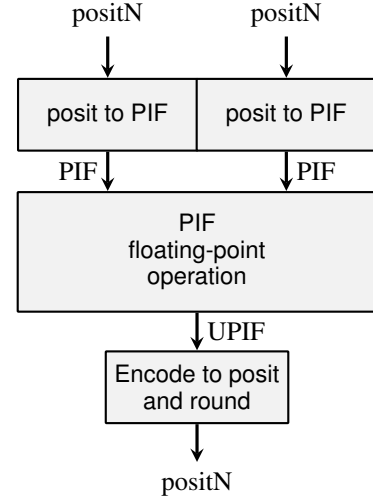


Fig. 3: Architecture of a posit operator in a PAU that uses posit registers and posit-to-posit operators.

exponent and significand in two's complement (where IEEE-754 uses a biased encoding for exponent and a sign-magnitude representation for the significand). This avoids unnecessary two's complement to sign-magnitude conversions (which may cost a carry propagation). It also has the side effect of slightly simplifying PIF addition of values with opposite signs.

Two's complement encoding for a normalized significand consists of a sign bit  $s$  and a fraction  $f$  on  $W_f$  bits. Significand value is then

$$v = -2s + \bar{s}.f \quad (19)$$

However, (19) does not allow the encoding of zero. The proposed PIF encoding scheme introduces an extra bit  $i$  which is one if and only if the represented value is strictly positive. Significand value is then

$$v = -2s + i.f \quad (20)$$

Zero is the only posit value whose PIF representation has both  $s$  and  $i$  set to zero. This fact allows efficient zero detection in arithmetic operators. For non zero values, exactly one of  $s$  or  $i$  is set.

PIF also has an extra  $isNaR$  bit, set to one if and only if the represented value is NaR. An alternative option could have been to use a non-posit value such as  $s$  and  $i$  set to one. This would however have implied more logic in the computation of  $s$  and  $i$ , and more logic to decode the value inside the operator.

To summarize, the PIF encoding scheme is composed of the following fields:

- a  $isNaR$  flag,
- the sign bit  $s$ ,
- the exponent  $e$  stored in two's complement on  $W_e$  bits,
- the weight one bit  $i$ ,
- the fraction bits  $f$  on  $W_f$  bits.

The encoded value is

$$v = \begin{cases} NaR & \text{if } isNaR \text{ is one} \\ (-2 + i + 0.f) \times 2^e & \text{otherwise} \end{cases} \quad (21)$$

TABLE II: Width of PIF and UPIF for standard posit.

N	$W_{es}$	$W_{PIF}$	$W_{UPIF}$
8	0	12	14
16	1	21	23
32	2	38	40
64	3	71	73

4) *Unrounded PIF encoding of the result of basic operations*: In the general case, the result of an operation on two PIF values is not exactly representable as a PIF, and must be rounded. As PIF is a floating-point format, we may use textbook techniques [14], [8] for this. For the basic operations (addition, multiplication, division and square root) the exact result can always be computed on at most  $2W_{PIF}$  bits, then for the purpose of rounding all the extra bits can be condensed into only two bits:

- an extra fraction bit at the LSB, called the *round* bit;
- a *sticky* bit, set if and only if the exact value is strictly greater than what is represented by the fraction  $f$  extended with the *round* bit (but still smaller than the next representable value). In other words, a *sticky* bit of zero means that the value represented by the extended fraction is exact.

We define the UPIF (Unrounded PIF) format as a PIF with these two extra bits.

The floating-point literature often uses a third additional bit (called the guard bit), useful in the case when a 1-bit normalization of the significand may be needed. In the big picture of Figures 3 and 4, the PIF operator is in charge of this normalization, so no guard bit is needed.

In this paper we only demonstrate the use of the UPIF format on addition/subtraction and multiplication, but it is equally suitable for division and square root. Digit recurrence algorithms [15] compute a remainder along with the quotient or square root, out of which the round and sticky bits can be computed (actually, for division the sticky bit is necessary non-zero if the round bit is set [8, Section 4.7.2], and therefore need not be computed if only round to nearest is needed). Multiplication-based algorithms [14], [8] also can output their result in UPIF format – for instance by computing the remainder.

Table II gives the width for PIF and UPIF associated with standard posit formats.

5) *Overflow management*: Posit arithmetic does not offer an overflow detection mechanism to the user. When the exact result of an operation is bigger than the biggest representable value, this biggest representable value is returned.

This saturation could in principle be handled in a generic way in the “Encode to posit and round” block of Figure 3, or in the “UPIF inplace round” block of Figure 4. However, as each operation leads to different overflow situations, it is more efficient to manage saturation in each PIF operator. Indeed, the UPIF specification exposed previously assumes that saturation has been performed by the operator, otherwise more bits would be needed. Another advantage is that some saturation situations may be detected in parallel with computation, thus reducing latency.

### III. ALTERNATIVES FOR A HARDWARE POSIT UNIT

This section describes the two options considered in this article for a processor supporting posit arithmetic. The detailed study of the second option is, to the best of our knowledge, novel. The purpose of this section is to define the hardware components that will be described in detail in Section IV.

#### A. Posit as a register encoding

Figure 3 shows that building posit-to-posit operators consists mainly in three steps:

- 1) decoding the posit representation to PIF,
- 2) performing the computation on PIF (building upon the floating-point literature, but without the overhead of the IEEE encoding, decoding, and special value support), and
- 3) rounding the result back to the nearest posit value.

Here, the two conversions between posit and PIF are essentially a posit overhead compared to minimal (non-IEEE) floating-point operators. As these blocks involve leading zero counting (LZC) and shifting, they could in principle be compared to hardware subnormal support in an IEEE-754 multiplier (subnormal support in addition is comparatively lighter [8]). However, it is well known that in an IEEE-754 multiplier only one subnormal input needs to be considered (the product of two subnormals is tinier than the tiniest representable value), whereas both posit inputs must be converted. Besides, subnormals are part of the floating-point encoding itself, which allows to hide their latency overhead e.g. by speculation. To the best of our knowledge this is not possible with posits. Finally, the LZCs and shifters in IEEE-754 only operate on significands, not on full words.

In any case, attacks by posit supporters pointing out the cost of subnormal support in IEEE-754 are misplaced: from a hardware perspective, the overhead of subnormals is due to the variable position of their rounding bit with respect to the leading 1. From this point of view, all posits are as bad as IEEE-754 subnormals.

#### B. Posit as a memory-only encoding

An alternative architecture to limit the latency impact of this conversion consists in using posit as a storage-only encoding. This architecture is depicted on Figure 4. The decoding and encoding blocks are placed on the memory path, while values in CPU registers are PIF-encoded.

Since multiple operations may occur before the result is written back to memory, rounding and encoding can no longer be fused. To ensure standard compliance, posit rounding must occur after each computation, in such a way that the PIF result exactly represents result of the operation defined by the posit standard. This is the role of the “UPIF inplace round” box.

In other words, in this architecture, the two conversion boxes, “posit to PIF” and “PIF to posit” must be exact, while the rounding and saturation logic must still be performed after each operation. This will still involve large shifters, however this logic is potentially cheaper than classical rounding: since the PIF is kept normalized, what must be shifted is masks

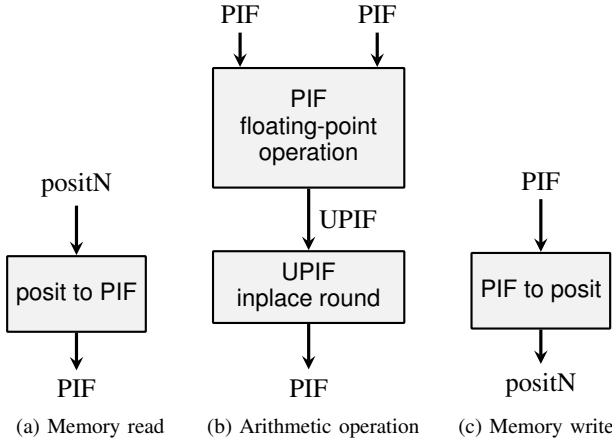


Fig. 4: Architecture of a PAU using posits as a memory-only encoding, with PIF registers and PIF-to-PIF operators.

and rounding bits, and these bit vectors that are simpler (more structured) than arbitrary significands. Details will be given in Section IV-C.

#### IV. POSIT HARDWARE DETAILED ARCHITECTURE

This section details the implementation of posit operators components provided by the MARTo library. The provided schematics aim at representing the usage of high level primitive inside operator function. As such, it should be easy for the interested reader to follow and check the library code. In case of discrepancy, the code is the reference.

##### A. Posit to PIF decoder

The proposed posit decoder is depicted in Figure 5.

The “LZOC + Shift” block (LZOC stands for “leading zero/one counter”) counts the range bits while discarding them, resulting in a normalized fraction.

The PIF exponent most significant bits  $e_h$  are computed out of the range count. If the leading bit is equal to  $s$ , then  $e_h = -l(= \bar{l} + 1)$ ; else  $e_h = l - 1$ . An optimization is to skip the first range bit when counting, effectively computing  $l' = l - 1$ . Indeed, if the first range bit is equal to  $s$ ,  $e_h = \bar{l}' + 1 + 1 = \bar{l}'$ , or  $e_h = l'$  otherwise. This high bit decoding method differs from the literature and avoids an addition when computing  $-l$ . The exponent least significant bits  $e_l$  are obtained by taking the  $ES$  first bits of the aligned fraction and xoring them with  $s$ .

The PIF exponent  $e$  is the concatenation of  $e_h$  and  $e_l$  (or is equal to  $e_h$  if the corresponding format has  $ES = 0$ ).

The PIF fraction  $f$  is directly read from the  $W_f$  least significant bits of the aligned fraction.

An OR reduction over the  $N - 1$  rightmost bits of the posit input is used to detect both zero and NaR values, in conjunction with  $s$ .

The weight 0 significand bit  $i$  is computed out of  $s$  and the detection of zero value.

The most expensive parts of this architecture are the “OR reduce” over  $N - 1$  bits to detect NaR numbers and the

combined leading zero/one counter and shifter (“LZOC + Shift”) that consumes the regime bits while aligning the significand.

It is interesting to compare this conversion to the decoding of special cases from IEEE-754 floats (which similarly must be performed on the inputs). There, one OR reduction on the exponent bits is needed to detect subnormals, another one on the significand bits is needed to detect zeros, and two similar AND reductions are needed to detect respectively NaN and infinities. The two OR reductions operate in total on the same width as the posit OR reduction, so the cost is the same. Then the shifter consists of  $\log_2 N$  multiplexer steps. In our combined LZOC + shifter implementation, the multiplexer at step  $i$  is driven by an AND reductions on  $2^i$  bits<sup>2</sup>. The combined sizes of these AND reduction is  $N$ , again matching the IEEE-754 ones. In summary, the posit decoding has the overhead of a shifter for  $N$ -bit data.

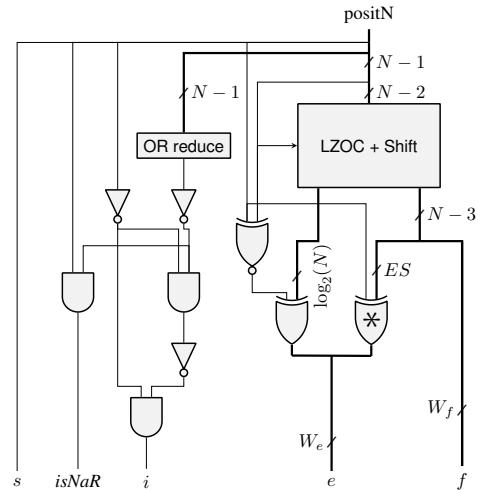


Fig. 5: Architecture of a posit to PIF decoder.

The decoder is slightly simplified with  $ES = 0$  posit formats, as it saves the XOR gates labeled \* on Figure 5.

##### B. UPIF to posit and PIF to posit

The complete UPIF to posit encoder architecture is shown in Figure 6.

Due to the variable-length encoding of posits, the fraction bit after which the rounding occurs depends on the exponent value. The UPIF to posit encoding process will handle both the rounding and shifting to the right position of the PIF fraction. In order to do this, the rounding bit is appended to the fraction, and the range is prepended.

Extra fraction bits are simultaneously shifted out and OR-reduced to a unique bit, which is OR-ed with the PIF sticky bit to get the final sticky bit.

A round-up bit is computed out of the sticky and round bits, which is added to the final encoding. The final result is

<sup>2</sup>Asymptotically faster implementations of LZC exist [8], but the one chosen here is better on the FPGAs used for our numerical experiments, thanks to very fast AND/OR reductions through the fast-carry logic.

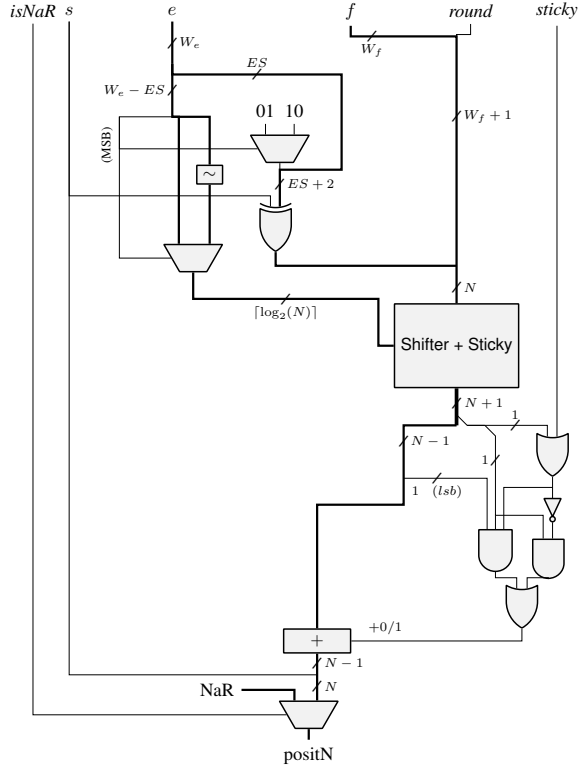


Fig. 6: Architecture of a UPIF to posit encoder. The PIF to posit encoder is similar, with the round and sticky logic (including the final adder) removed.

either the computed encoding or NaR representation if the PIF  $isNaR$  bit is set.

When working with  $ES = 0$  posit format, a special case detection box has to be added to detect and forbid the rounding up that would cause the value to round to NaR or 0.

PIF to posit conversion (as in Figure 4) is simpler, as this process is exact: in this case there is no need for the rounding logic. Its complexity is delegated to the UPIF inplace rounding architecture which we detail now.

### C. UPIF inplace round

When working with posit as a memory format (Figure 4), the rounding and encoding step are distinct. Indeed, the UPIF result must be rounded as if it had been converted to posit, then converted back to PIF.

One option would indeed be to shift the significand, round it, then shift it back. But then there would be no latency advantage to this architecture. A cheaper alternative is to keep the significand fixed, and perform the rounding and mantissa resizing using bitwise binary operations with bit masks that are shifted to the proper place.

In details, the last range bit,  $e_s$  bits and complete fraction are concatenated, giving a posit stem of width  $W_s = 1 + ES + W_f$ . A cut bound  $c(e)$  is determined from the exponent  $e$ , which is the number of bits that “would fall behind posit representation limit” due to the range width.

The masks used are

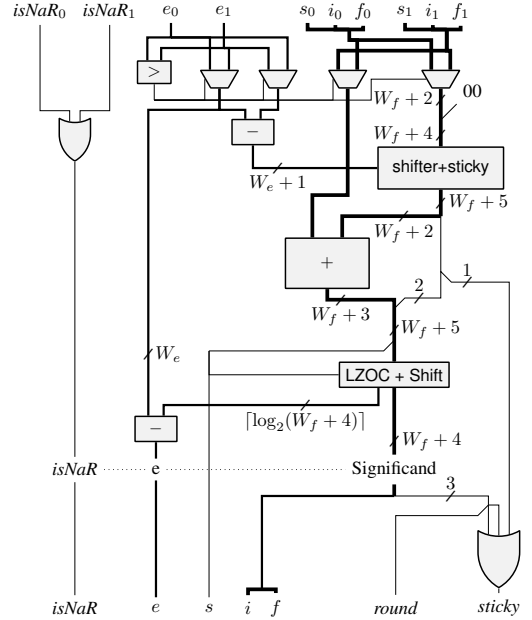


Fig. 7: Architecture of a PIF adder. Exponent comparison block denoted with “>” also takes the operand  $i$  and  $s$  bits to detect zero values, but wires have been omitted here for clarity.

- a round mask with only the last untruncated bit (of weight  $c(e)$ ) set to one,
- a guard mask corresponding to the previous round mask shifted one bit to the right,
- a sticky mask with all bits below guard bit set to one,
- a keep mask with all bit above the round bit (included) set to one.

These masks could be constructed by dedicated shifters, but our current implementation retrieves them from Look-Up Tables (LUTs) inputting the exponent. This avoids the need of any shift, and is extremely efficient on the LUT-based FPGAs on which we conduct our quantitative evaluation. When  $ES = 0$ , the same precaution as for the fused round/encode operator should be taken.

### D. PIF floating-point operations

The architectures of the PIF adder/subtractor (Figure 7) and multiplier (Figure 8) first compute the exact result (top part of the figures) using the transposition to the PIF format of classical floating-point algorithms.

Although the adder is a single-path architecture [8], its datapath can be minimized thanks to the classical observation that large shifts in the two shifters are mutually exclusive. Indeed, the normalizing LZOC+Shift of Figure 7 will only perform a large shift in a cancellation situation, but such a situation may only occur when the absolute exponent difference is smaller than 1, which means that the first shift was a very small one. Conversely, when the first shifter performs a large shift, the rightmost part of the significand can be immediately compressed into a sticky bit, since we know that it will not be shifted back by the second LZOC+Shift. All this allows us



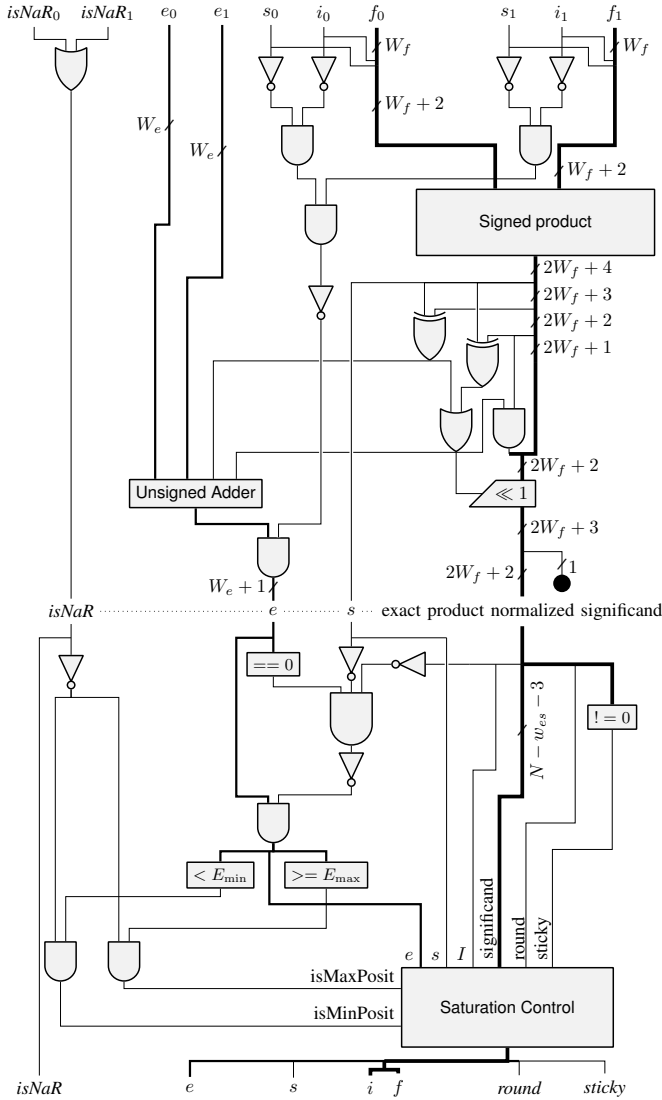


Fig. 8: Architecture of a PIF multiplier.

to keep most intermediate signals on  $W_f + 2$  to  $W_f + 6$  bits, where previous work [10], [11] seem to use datapaths that are twice as large.

The posit multiplier shown in Figure 8 is straightforward. It performs the addition of the exponent, the signed product of the significands, shifts the output and corrects the exponent if necessary. This architecture aims at a minimal area and time cost. If energy efficiency is the goal, alternative architecture have been proposed that exploit the relation between exponent magnitude and precision to disable the computation of unneeded product lsb [16].

The bottom part of Figures 7 and 8 normalize the exact result computed by the top parts to a PIF. For both operators, the exact significand must be realigned, correcting the exponent accordingly.

## V. HARDWARE SUPPORT FOR EXACT ACCUMULATION

The idea of an exact accumulator currently has a lot of momentum. Several existing machine learning accelerators

TABLE III: Quire bit-width parameters for standard posit.

Posit		Quire sizes					
N	ES	C	$W_q$	$W_O$	$W_R$	$W_L$	$W_Z$
8	0	6	32	12	13	16	6
16	1	14	128	42	57	64	28
32	2	30	512	150	241	256	120
64	3	62	2048	558	993	1024	496

[17], [18] already use variations of the exact accumulator to compute on IEEE-754 16-bit floating-point. Other application-specific uses have been suggested [19], [20]. For larger sizes, this could be a useful instance of “dark silicon” [21]. This trend was also anticipated with the reduction operators in the IEEE 754-2008 standard, although without the requirement of exactness.

### A. Quire specification and parameters

The quire is not yet completely specified in the posit draft standard [3]. Currently, the standard specifies a binary interchange format, which consists in a very large fixed point number of size  $W_q$  defined by (17). In the sequel, we discuss the cost of hardware support for a quire with the range of the interchange format. Note however that the draft standard defines fused operations as *those expressible as sums and differences of the exact product of two posits; no other fused operations are allowed*. Hardware quire support is a way to implement such fused operations, but this formulation does not prevent cheaper specific implementations of useful fused operations such as *Fused Multiply-Add (FMA)* [8], complex multiplication [22], or even full convolutions for neural networks.

The parameter  $W_q$  defines the storage requirement, and thus a lower bound of the area cost. It also entails a large delay: Figure 9 shows a high-level functional description of a quire accumulation, and shows that there is a  $W_q$ -bit addition on the critical path from the quire to itself. A technique that enables 1-cycle accumulation (the architecture must be able to add one new summand to the quire at each cycle) at high frequencies is to use for the quire a high-radix carry-save internal format [19], [23]. It is briefly reviewed below. With this technique, the high latency is still there, but delayed to the conversion from the quire to a posit.

The posit draft standard [3] specifies NaR as a special quire value. Testing this special value at each new quire operation would require to check the quire equality with this special value. Instead, this work proposes to add a flag bit that signals that the value held in the quire is NaR. This bit is set when NaR is added to the quire and stays set until the end of the computation. This extra bit can replace one of the quire carry bits. A slightly more expensive alternative would be to encode and decode NaR value when transferring quire to/from memory.

In the posit context, it is natural to use two’s complement for managing signs in the quire. Note that some implementations of Kulisch’s exact accumulator seem to use a sign-magnitude representation for the accumulator [13], matching the sign-magnitude representation of IEEE floating-point, but even

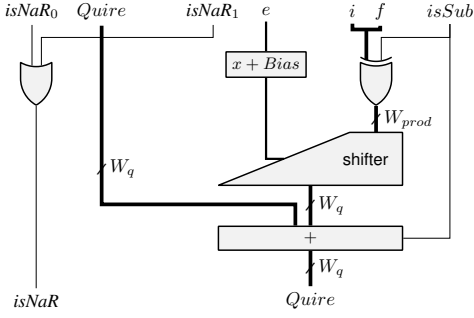


Fig. 9: Architecture of a posit quire addition/subtraction.

there a two's complement representation of the accumulator is more efficient [17], [23].

### B. Addition of products to the quire

The posit quire is able to perform exact sums and sums of products. Therefore, the input format of the quire is defined as the output of the exact multiplier from Figure 8 (top).

To add a simple posit to the quire, it is first converted to PIF, then the PIF value is trivially cast to the same exact multiplier format (the details are skipped for brevity).

The simplest implementation of the quire addition/subtraction is depicted in Figure 9. An exact posit product fraction is shifted to the correct place to the quire format according to its exponent. A large adder then performs the addition with the previous quire value. The two's complement subtraction is performed at the cost of a XOR on the input and a carry-in to the adder, as in the posit adder/subtractor. For this the shifter must be a sign-extending one.

The simple architecture of Figure 9 can be used directly for small sizes (up to posit16). For larger sizes, the long carry propagation delay of the addition in this architecture will restrict the maximum frequency achievable. To address this, a cost-effective solution [19], [23] is to segment the quire into smaller words (typically standard 32-bit or 64-bit words). Carry propagation is then limited to a segment, and the carries between segments are stored in registers and propagated to the next segment during the next cycle. Another point of view is that the quire is kept in a high-radix carry-save redundant format (radix is  $2^{32}$  or  $2^{64}$ ). If such a format is used, its conversion to a non-redundant format will incur additional overhead to complete carry propagation. Some hardware can be dedicated to this, but a cheap alternative is simply to dedicate a few cycles to the completion of the carry propagation, during which the summand input to the quire is kept at zero. The number of carry-propagation cycles is  $W_q/32$  for 32-bit segments. These extra cycles are amortized if the quire is used for summing large numbers of values.

Several variants of unsegmented and segmented quires will be evaluated in Section VI.

### C. Conversion from quire to posit

The conversion of the quire value to a posit is divided in two steps. The quire is first converted to a UPIF value (architecture

depicted in Figure 11) before the latter is encoded to a posit (Section IV-B).

There are four distinct cases to take into account when converting the quire to the UPIF:

- If the quire holds a NaR value, the result is NaR;
- If the quire value is larger in magnitude than the maximum-magnitude posit (overflow), the latter should be returned (saturation);
- If the quire value belongs to the representable posit range, it should be converted;
- If the quire value is smaller in magnitude than the minimum-magnitude non-zero posit (underflow), the latter should be returned (saturation);

Figure 10 illustrates the different interesting zones of a quire. The values of the parameters appearing in this figure are determined as follows.

Detection of overflow consists in comparing all the bit in the overflow zone with the sign bit. If at least one differs, the posit overflows. The width of the overflow zone  $W_O$  is computed as follow:

$$\begin{aligned} W_O &= Q_{\text{msb}} - E_{\text{max}} \\ &= E_{\text{max}} + C \end{aligned} \quad (22)$$

For quire values inside the posit range, a normalization should be performed, which uses a LZOC + shifter of ideal input width  $W_R$ , with

$$W_R = E_{\text{max}} - E_{\text{min}} + 1 \quad (23)$$

Actual LZOC + shifter input width  $W_L$  is the immediate next power of two for implementation efficiency reason.

Finally, to detect the difference between an underflow value and a real zero, a wide or is performed on the underflow zone of width  $W_Z$ , with

$$\begin{aligned} W_Z &= E_{\text{min}} - Q_{\text{lsb}} \\ &= E_{\text{max}} \end{aligned} \quad (24)$$

With (18) and (11) we get

$$W_O = \frac{N^2}{8} + \frac{3N}{4} - 2 \text{ for standard posits.} \quad (25)$$

$$W_Z = \frac{N^2}{8} - \frac{N}{4} \text{ for standard posits.} \quad (26)$$

$$W_R = \frac{N^2}{4} - \frac{N}{2} + 1 \text{ for standard posits.} \quad (27)$$

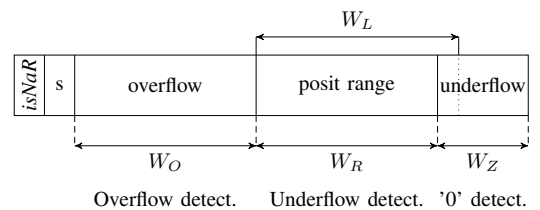


Fig. 10: The bits of a standard quire.

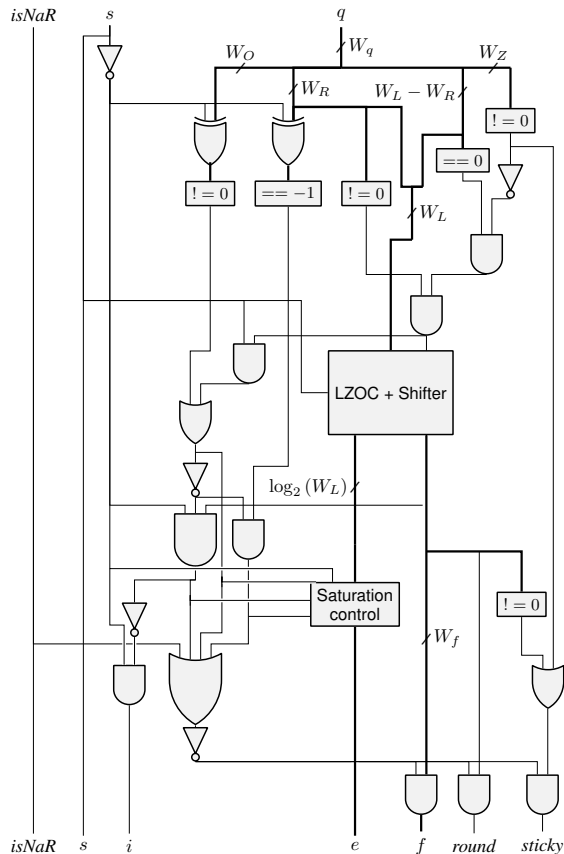


Fig. 11: Architecture of the conversion from qire to UPIF.

$$W_L = \frac{N^2}{4} \text{ for standard posits.} \quad (28)$$

Numerical values are reported in Table III.

## VI. EVALUATION

This section introduces the MARTo library, compares it with the state of the art, then uses it to compare posit and floating-point operators.

Comparisons of MARTo with other works use the target and toolchain that is closest to those used in the work being compared to. All other results in this section are post place and route, obtained using Vivado HLS and Vivado 2020.1 targeting a Virtex-7 FPGA (xc7vx330tffg1157-1).

### A. The MARTo library of posit and floating-point operators

The architectures described in the previous sections have been implemented using HLS-compliant templated C++ code. The library defines parameterized types such as `PositEncoding<N, ES>` and `IEEENumber<WE, WF>`, and functions to perform arithmetic operations on those types. The operators are built on top of the HINT library [24] in order to be compatible with multiple HLS tools, and to benefit from optimised primitives such as the fused LZOC+Shift of Figure 5 or the fused Shifter+sticky of Figures 6 and 7.

### B. Correctness of the operators

In order to verify that the proposed architectures are correct, the following functional tests have been run:

- Exhaustive test of addition and multiplication of standard posit8 and standard posit16 against softposit, for both end-to-end strategy and in place rounded strategies,
- Exhaustive test of inplace rounding for standard posit16,
- Some corner case tests of qire addition/subtraction and conversion back to posit for posit16.
- Exhaustive test for addition/product of IEEE16 against SoftFloat, for the five IEEE-754 rounding modes.

The previous tests were purely software runs of the C++ code. As HLS synthesis tools can not yet be completely trusted, we also checked that the tool used for most experiments, Vivado HLS, worked well on the proposed operators.

First, we checked that the VHDL file produced by the Vivado HLS compiler for a posit adder was functionally correct by performing an exhaustive test using a VHDL simulator. Here the test bench inputs a pair of posits and the expected output posit, and raises an error if the computed result differs from the expected one. Scripts and sources to reproduce this experiment are accessible from the MARTo repository.

Then the standard posit16 multiplier was synthesized, placed and routed for the Zynq FPGA of a Zybo board using the Vivado HLS toolchain, and the resulting FPGA circuit was exhaustively checked against softPosit executed on the ARM core of the Zynq.

Both tests were successful, and confirmed that we could trust the tools used here for evaluation.

### C. Evaluation of pipelined designs

Vivado HLS synthesis proceeds in two steps. The first step consists in converting the C++ code to a pipelined design described in a Hardware Description Language (HDL), the second step is the synthesis and implementation of this HDL. Based on a target clock frequency, the first step estimates how many pipeline levels are required. However, this estimation is not accounting for some optimization the synthesis tool is able to make. The estimation is therefore quite often a large overestimation of what is really required. In order to get meaningful results, we used an automated exploration to find the smaller pipeline depth that still allows the design to run at the target clock period. The exploration starts by taking the initial estimation of the HLS compiler as the higher bound of the required depth, and use a dichotomy and pipeline depth constraints on the HLS tool to get this minimal pipeline depth level.

The script performing this dichotomy is open source, and also accessible from the MARTo repository for reproducibility.

### D. Comparison with the posit state-of-the-art

As we eventually observe that posits are larger and slower than IEEE floats, it is important to be convincing that we are not using a substandard posit implementation. For this purpose, Table IV gathers the results of best-effort comparisons with the current state of the art in posit hardware. It shows that MARTo is a definite improvement of this state of the art.

TABLE IV: Comparison with state-of-the-art hardware posit implementations [11], [10], [12], [25]

(a) Comparison with [11] for standard posit addition and product

	Op	Format	LUT	DSP	Delay (ns)
[11]	+	<16, 1>	391	0	32.4
		<32, 2>	981	0	40.0
	×	<16, 1>	218	1	24.0
		<32, 2>	572	4	33.0
MArTo	+	<16, 1>	<b>299</b>	0	<b>24.2</b>
		<32, 2>	<b>704</b>	0	<b>33.9</b>
	×	<16, 1>	<b>213</b>	1	<b>19.4</b>
		<32, 2>	<b>483</b>	4	<b>28.9</b>

As no sources is provided, we report as-is the figures from [11], obtained for a Zynq-7000 (xc7z020clg484-1) with Vivado 2017.4. To limit the possible effect of tool improvement on the synthesis, MArTo synthesis results have been obtained for the same part with Vivado HLS/Vivado 2018.3, the oldest version available for download at time of experimentation.

(b) Comparison with [10] on standard posit addition and product

	Op	Format	ALM	DSP	Cycles	FMax (MHz)
[10]	+	<16, 1>	~500	0	~49	~550
		<32, 2>	~1000	0	~51	~520
	×	<16, 1>	~330	1	~35	~600
		<32, 2>	~600	<b>1</b>	~38	~550
MArTo	+	<16, 1>	<b>274</b>	0	<b>11</b>	<b>564</b>
		<32, 2>	<b>696</b>	0	<b>17</b>	<b>562</b>
	×	<16, 1>	<b>280</b>	1	<b>15</b>	<b>600</b>
		<32, 2>	<b>452</b>	2	<b>21</b>	445

Synthesis reported in [10] target Stratix V FPGA. Results are read from a graphic plot, hence the approximate values. As there is no version of the Intel HLS toolchain that supports both Stratix V and the C++ 11 standard used in MArTo, the C++ to HDL compilation is done using Vivado HLS. The obtained HDL is then synthesised and routed for Stratix V using Quartus. Despite being baroque, this toolchain seems to give good results, except for the <32, 2> product where it lacks the knowledge of the target’s DSP possible configurations. Indeed, the product is computed using a 36x36 configuration of the DSP block, where a 27x27 configuration would be faster.

(c) Comparison with [12] on posit<32,6> addition and product

	Op	LUTs	DSPs	Cycles	Delay (ns)
[12]	+	946	0	<b>5</b>	4.1
	×	854	<b>1</b>	<b>6</b>	4.4
MArTo	+	<b>792</b>	0	<b>5</b>	<b>3.9</b>
	×	<b>435</b>	2	<b>6</b>	<b>4.1</b>

MArTo synthesis have been performed using Vivado HLS/Vivado 2020.1 using part xc7vx330t-ffg1157-3. Experimental settings of [12] use the same part, but tool version is not reported.

(d) Comparison with [25] for standard posit addition and product

	Op	Format	LUT	DSP	Delay (ns)
[25]	+	<16, 1>	383	0	27.25
		<32, 1>	939	0	35.8
	×	<16, 1>	<b>201</b>	1	20.9
		<32, 1>	571	4	29.2
MArTo	+	<16, 1>	<b>300</b>	0	<b>25.5</b>
		<32, 1>	<b>672</b>	0	<b>34.5</b>
	×	<16, 1>	205	1	<b>19.2</b>
		<32, 1>	<b>472</b>	4	<b>28.8</b>

MArTo synthesis have been performed using Vivado HLS/Vivado 2020.1 using part xc7z020clg484-1.

There is less pressure to show that the MArTo implementation of IEEE floats is efficient. A comparison with Xilinx implementation of IEEE floats is provided in Table VI. There, the line labeled Xilinx Float corresponds to IP used by Vivado HLS when using the `float` and `double` data types in the HLS C++ (hence the lack of 16-bit results). This hard IP is the industry standard when using Vivado, and can be considered a state-of-the-art implementation of floating-point for Xilinx FPGAs. It supports some of the IEEE features, such as infinity and NaN encoding. However, it is not IEEE-compliant: although the memory format is that of IEEE floats, subnormals are flushed to zero to save resources. A comparison between Xilinx and MArTo on the IEEE format can therefore be used to highlight the cost of subnormal handling.

Table VI shows that the Xilinx Float adders use DSP blocks to implement some of the shifters (a shift being a multiplication by a power of two). This reduces their logic and register cost. Considering this, the hardware costs of Xilinx adders and IEEE adders (obtained with MArTo) are really comparable. This illustrates that the overhead of subnormal handling in floating-point adders is small. Conversely, there is in Table VI a very large difference in the resources used in multipliers. This demonstrates the cost of hardware subnormal handling in floating-point multipliers.

The Xilinx IP also seem to have a fixed pipelined, and do not benefit from a relaxed clock constraint to reduce the latency, hence their important latency compared to other solutions.

Since Xilinx floats lack subnormal support, in the following we must base on MArTo only our posit versus IEEE comparisons.

One may wonder if this comparison doesn’t also highlight some inefficiency of hardware generated using HLS tools, but recent works [26], [27] suggest that such overhead is becoming negligible.

### E. Comparison with floating-point operators

Table V compares combinatorial implementations of posits and floats of the same size on addition and multiplication. In this table, the “posit→posit” lines present results for the classical posit operators of Figure 3. The “PIF→PIF” lines presents results for the posit-compatible PIF operators that use the architecture of Figure 4b, including the inplace round component.

A first observation is that posit arithmetic is indisputably both larger and slower than IEEE-754 arithmetic. This contradicts the comparison in [11], which seems to use a very poor IEEE implementation.

As expected, the PIF-to-PIF operators are lighter and faster than the posit-to-posit ones. They still pay the price in area of a wider significand datapath (see Table I) compared to IEEE operators: for the adders, PIF-to-PIF consume more LUTs than IEEE operators; for multipliers, they consume more DSP blocks (there is a step effect that is due to the discrete nature of DSP blocks). Again we observe in the IEEE multipliers the logic cost of subnormal support, but we also observe a comparable cost in that the PIF multiplier, essentially due to the inplace round logic. Still, the PIF to PIF operators achieve

TABLE V: Synthesis results of posit and IEEE-754 combinatorial adders and multipliers.

(a) Combinatorial adder						
	N	LUT	(ratio)	delay	(ratio)	
posit→posit	16	312	1.33	11.1 ns	1.27	
	32	647	1.49	15.8 ns	1.33	
	64	1550	1.59	21.6 ns	1.35	
PIF→PIF	16	237	1.01	9.7 ns	1.10	
	32	562	1.29	12.9 ns	1.08	
	64	1244	1.27	14.7 ns	0.92	
IEEE→IEEE	16	234	1	8.8 ns	1	
	32	434	1	11.9 ns	1	
	64	976	1	16.0 ns	1	

(b) Combinatorial multiplier						
	N	LUT	(ratio)	DSP	delay	(ratio)
posit→posit	16	182	1.03	1	11.3 ns	1.39
	32	466	1.37	4	15.8 ns	1.62
	64	1213	1.58	16	21.1 ns	1.48
PIF→PIF	16	120	0.68	1	7.8 ns	0.96
	32	291	0.86	4	11.5 ns	1.17
	64	695	0.90	16	15.3 ns	1.08
IEEE→IEEE	16	176	1	1	8.1 ns	1
	32	340	1	2	9.8 ns	1
	64	768	1	9	14.3 ns	1

(c) Posit - PIF converting operators			
	N	LUT	delay
Posit to PIF	16	61	2.59 ns
	32	106	4.74 ns
	64	278	5.52 ns
PIF to posit	16	41	2.12 ns
	32	98	2.50 ns
	64	301	2.83 ns

delays that are closer to those of IEEE operators than to those of posit operators, which was their main motivation.

Note that the area cost of PIF/posit conversions (altogether about half the size of a complete IEEE adder) must still be paid in a posit arithmetic unit that uses the PIF-to-PIF approach. Only its delay (altogether about half the delay of a complete IEEE adder) is avoided. However, there is another advantage in a PIF-to-PIF PAU: these conversions are naturally shared between different operations. Such sharing of the conversion hardware is also possible in a posit-to-posit PAU, but then it will restrict instruction-level parallelism.

Table VI compares pipelined versions of the same operators, targeting a frequency of 333 MHz (3ns cycle time), and producing one output per clock cycle. The method described in Section VI-C is used to get the smallest latency operator that meets these constraints. There is no PIF to PIF line in this table: for this setup, the PIF to PIF approach fails to provide any latency improvement (the arithmetic operators require the same number of cycles, and sometimes require one more cycle). We therefore choose not to report these results, which we consider synthesis artifacts: they are inconsistent

TABLE VI: Synthesis results of posit and IEEE-754 pipelined adders and multipliers.

(a) Pipelined adder						
	N	LUT	Reg.	DSP	cycles	delay
Posit	16	320	128	0	4	2.69 ns
	32	719	460	0	7	2.83 ns
	64	1635	1207	0	10	2.93 ns
IEEE	16	193	137	0	4	2.90 ns
	32	435	337	0	6	2.88 ns
	64	1001	880	0	10	2.99 ns
Xilinx Float	32	167	355	2	10	2.43 ns
	64	628	758	3	10	2.43 ns

(b) Pipelined multiplier						
	N	LUT	Reg.	DSP	cycles	delay
Posit	16	213	80	1	4	2.85 ns
	32	443	198	4	6	2.93 ns
	64	1140	811	16	12	4.10 ns
IEEE	16	189	122	1	4	2.69 ns
	32	381	246	2	6	2.74 ns
	64	783	801	9	8	2.67 ns
Xilinx Float	32	82	151	3	5	2.72 ns
	64	115	494	11	10	2.75 ns

with the expectations and with Table V.

#### F. The cost of supporting all rounding modes in IEEE

Tables V and VI report result for IEEE operators that only support round to nearest, ties to even. Another example of IEEE complexity that translates to very little hardware cost is the support of the 5 standard rounding modes. For instance, adding this support to the 32-bit adder increases its area from 434 to 458 LUTs and actually decreases the delay from 11.9 to 11.7ns (another synthesis artifact). It remains well below the posit cost.

#### G. Quire versus standard operations

Synthesis results for the quire are given in Table VII, where we use MarTo to write a C++ loop that performs the sum of 1000 products and return the result as a posit. They are compared to a similar loop using floating-point Kulisch accumulator, and using regular floating-point hardware.

Quire is presented in unsegmented (U) version along with two segmented versions (S32 and S64 for segments of 32 or 64 bits). For 32 bits, the unsegmented version is not able to achieve 3ns cycle time, due to the long carry propagation.

The Kulisch accumulator used in this paper also uses a 2's complement segmented accumulator [23, variant 3], but with a final conversion to float that is IEEE-compliant (round to nearest, ties to even). The implementation has been validated against MFPR [28] simulations.

A first observation is that, unsurprisingly, the cost and performance of a posit32 quire and a Kulisch accumulator for 32 bits floats are almost identical.

TABLE VII: Synthesis results for a sum of 1000 products (U: Unsegmented, S32 and S64: Segment sizes of 32 and 64 bits).

		LUT	Reg.	DSP	cycles	delay
quire 16	U	1200	1026	1	1019	2.70 ns
	S32	978	1062	1	1021	2.68 ns
	S64	1004	958	1	1019	2.36 ns
quire 32 (512 bits)	U	5884	6235	4	1031	3.65 ns
	S32	3641	7237	4	1040	2.89 ns
	S64	3513	5189	4	1033	2.78 ns
Kulisch 32 (559 bits)	S32	3624	7632	2	1034	2.937
	S64	3612	5165	2	1026	2.801
IEEE Float 32		840	711	2	6012	2.92 ns
IEEE Float 64		1798	1723	9	8015	3.33 ns
Xilinx Float 32		445	544	3	6008	2.72 ns
Xilinx Float 64		809	1386	11	8013	2.70 ns

Classically, using an exact accumulator consumes vastly more resources than using standard operators: a factor 10 for 32-bit floats (a smaller factor for posits, but only due to the higher cost of the standard operators). Such factors should not come as a surprise: the 512 bits of the posit32 quire are indeed 18 times the 27 bits of the posit32 significand. A claim of Kulisch is that the increase in accuracy justifies this cost.

Another advantage of exact accumulation is that it offers a latency reduction proportional to the latency of the floating-point or posit adder (here a factor 6-7). This is thanks to the fact that the accumulation loop is 1/ a fixed-point addition and 2/ exact, which offers opportunities to exploit more parallelism in its computations[17], [20].

Detailed synthesis results of the quire sub-components are given in Table VIII. The *quire addition* line reports the cost for the architecture of Figure 9, including the large shifter and the fixed-point accumulation loop. This component can be pipelined with an initiation interval of one cycle, in other words it accepts a new input every cycle. The two other lines describe the conversion of the quire result back to posit. The *carry propagation* consists in a loop component will actually be merged with the *quire addition* during synthesis, reducing its cost. However, there is an irreducible latency for the final carry propagation once the accumulation is over.

The latency overhead of the expensive conversion from quire to float or posits is easily amortized for large loops. However, it is also clear that a hardware quire will be very slow when used for small sequences of operations (e.g. fused multiply and add, complex arithmetic, small matrices used in graphical applications, additive range reduction for elementary functions, etc). For such small computations, we expect other accuracy-enhancing techniques [6] to remain competitive.

## VII. CONCLUSION

The purpose of this work is to compare the cost of hardware arithmetic operators for two competing number systems: the established IEEE-754 system and its posit challenger.

This comparison is performed thanks to a library of operators for the two systems, providing hardware description that

TABLE VIII: Detailed synthesis results of hardware posit quire (U: Unsegmented, S32 and S64: Segment sizes of 32 and 64 bits).

(a) Posit 16						
		LUT	Reg.	Cycles	Delay (ns)	
Quire addition	U	618	885	4	2.576	
	S32	403	585	3	1.886	
	S64	444	606	3	1.984	
Carry prop.	S32	6*	390	3	1.539	
	S64	2*	261	2	1.651	
Quire to posit		480	166	3	2.735	
(b) Posit 32						
		LUT	Reg.	Cycles	Delay (ns)	
Quire addition	U	3609	4986	7	3.212	
	S32	1305	2265	3	2.791	
	S64	1389	2276	3	2.791	
Carry prop.	S32	281*	2874	8	2.851	
	S64	189*	2391	7	2.183	
Quire to posit		1845	1457	17	2.878	

LUT count for quire carry propagation is very low as the HLS tool is able to detect the possibility of fast carry logic utilization. Indeed, the carry propagation consists in adding 0 *nbanks* times to the quire, but there is no need to store zero in the LUTs.

are state-of-the art for posit, and high quality (if not state of the art) for IEEE-754. This open-source library is provided as header-only templated C++, designed for modern High-Level Synthesis tools.

Posit-to-posit operators are shown to be significantly more expensive, both in terms of resources and delay, than IEEE operators for the same input width. For instance, addition and multiplication on 32-bit standard posits require about 50% more hardware and about 50% more delay than standard-compliant addition of binary32 floats. This overhead should be put in balance with the increased accuracy sometimes offered by posits. On the example of 32-bit formats, posits offer up to 3 extra bits of accuracy (a 11% improvement) in a limited domain of exponents (while degrading the accuracy outside of this domain due to tapered precision).

An original alternative implementation of posits is proposed: it keeps posits decoded in a wider intermediate format to avoid some of the posit encoding overhead. This alternative leads to operations that are comparable in delay to IEEE floats, but at a higher cost, all the more as it requires wider internal registers which also have a system-wide cost.

This article also provides and compares exact accumulators in both systems, without a clear advantage on a side or the other.

If there is a take-away message in this study, it would be that the indisputable complexity of the IEEE-754 standard, much attacked by posit proponents, does not necessarily translate into expensive hardware. Among the features that the posit system discards as useless, most (in particular overflow man-

agement, NaNs, and directed rounding mode) were designed to be implementable at very little cost. The only really expensive feature is subnormal support, due to rounding happening in a variable position of a bit vector. Posit arithmetic, despite the simplicity and elegance of the number system, involves such variable-position rounding, and therefore entail an overhead that is comparable in nature to subnormal support.

This work has framed baseline posit implementations. On this basis, it is possible to consider many optimizations studied for floating-point operators (such as dual-path architectures, leading zero anticipation, or various forms of hardware speculation). These optimizations will improve delay, but at the expense of area.

Before that, future work includes completing the library with missing operations (starting with division and square root). HLS has the potential of making it very easy to study, at the application level, the impact of number systems on cost, performance, and accuracy. This is the long-term goal of the library presented here.

### Acknowledgements

This work was partly funded by the Imprenum project of Agence Nationale de la Recherche. Many thanks to Orégane Desrentes for her corrections to some of the figures.

### REFERENCES

- [1] "IEEE standard for floating-point arithmetic," IEEE 754-2008, also ISO/IEC/IEEE 60559:2011, Aug. 2008.
- [2] J. L. Gustafson and I. T. Yonemoto, "Beating floating point at its own game: Posit arithmetic," *Supercomputing Frontiers and Innovations*, vol. 4, no. 2, pp. 71–86, 2017.
- [3] P. W. Group, "Posit standard documentation," Jun. 2018, release 3.2-draft.
- [4] Z. Carmichael, H. F. Langroudi, C. Khazanov, J. Lillie, J. L. Gustafson, and D. Kudithipudi, "Performance-efficiency trade-off of low-precision numerical formats in deep neural networks," in *Proceedings of the Conference for Next Generation Arithmetic*. ACM, 2019, pp. 3:1–3:9.
- [5] P. Lindstrom, S. Lloyd, and J. Hittinger, "Universal coding of the reals: alternatives to IEEE floating point," in *Proceedings of the Conference for Next Generation Arithmetic*. ACM, 2018, p. 5.
- [6] F. De Dinechin, L. Forget, J.-M. Muller, and Y. Uguen, "Posits: the good, the bad and the ugly," in *Proceedings of the Conference for Next Generation Arithmetic*. ACM, 2019, p. 6.
- [7] N. Buoncristiani, S. Shah, D. Donofrio, and J. Shalf, "Evaluating the numerical stability of posit arithmetic," in *International Parallel and Distributed Processing Symposium (IPDPS)*. IEEE, 2020, pp. 612–621.
- [8] J.-M. Muller, N. Brunie, F. de Dinechin, C.-P. Jeannerod, M. Joldes, V. Lefèvre, G. Melquiond, N. Revol, and S. Torres, *Handbook of Floating-Point Arithmetic, 2nd edition*. Birkhauser Boston, 2018.
- [9] D. R. Lutz, "ARM floating-point 2019: Latency, area, power," in *26th Symposium on Computer Arithmetic*. IEEE, 2019, pp. 69–76.
- [10] A. Podobas and S. Matsuoka, "Hardware implementation of POSITs and their application in FPGAs," in *International Parallel and Distributed Processing Symposium Workshops*. IEEE, 2018, pp. 138–145.
- [11] R. Chaurasiya, J. Gustafson, R. Shrestha, J. Neudorfer, S. Nambiar, K. Niyogi, F. Merchant, and R. Leupers, "Parameterized posit arithmetic hardware generator," in *36th International Conference on Computer Design (ICCD)*. IEEE, 2018, pp. 334–341.
- [12] M. K. Jaiswal and H. K.-H. So, "Pacogen: A hardware posit arithmetic core generator," *IEEE Access*, vol. 7, pp. 74 586–74 601, 2019.
- [13] U. W. Kulisch, *Advanced Arithmetic for the Digital Computer: Design of Arithmetic Units*. Springer-Verlag, 2002.
- [14] M. D. Ercegovic and T. Lang, *Digital Arithmetic*. Morgan Kaufmann, 2004.
- [15] —, *Division and Square Root: Digit-Recurrence Algorithms and Implementations*. Kluwer Academic Publishers, Boston, 1994.
- [16] H. Zhang and S. Ko, "Design of power efficient posit multiplier," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 67, no. 5, pp. 861–865, 2020.
- [17] N. Brunie, "Modified Fused Multiply and Add for exact low precision product accumulation," in *24th Symposium on Computer Arithmetic (ARITH-24)*. IEEE, Jul. 2017.
- [18] J. Johnson, "Rethinking floating point for deep learning," Facebook, arXiv:1811.01721, 2018.
- [19] F. de Dinechin, B. Pasca, O. Creț, and R. Tudoran, "An FPGA-specific approach to floating-point accumulation and sum-of-products," in *Field-Programmable Technologies*. IEEE, 2008, pp. 33–40.
- [20] Y. Uguen, F. de Dinechin, V. Lezard, and S. Derrien, "Application-specific arithmetic in high-level synthesis tools," *ACM Transactions on Architecture and Code Optimization*, vol. 17, no. 1, 2020.
- [21] M. B. Taylor, "Is dark silicon useful? harnessing the four horsemen of the coming dark silicon apocalypse," in *Design Automation Conference*. ACM, 2012.
- [22] H. H. Saleh and E. E. Swartzlander, "A floating-point fused dot-product unit," in *International Conference on Computer Design (ICCD)*, 2008, pp. 426–431.
- [23] Y. Uguen and F. de Dinechin, "Design-space exploration for the Kulisch accumulator," Mar. 2017, working paper or preprint. [Online]. Available: <https://hal.archives-ouvertes.fr/hal-01488916>
- [24] L. Forget, Y. Uguen, F. de Dinechin, and D. Thomas, "A type-safe arbitrary precision arithmetic portability layer for HLS tools," in *International Symposium on Highly Efficient Accelerators and Reconfigurable Technologies*, Nagasaki, Japan, Jun. 2019, pp. 1–6.
- [25] F. Xiao, F. Liang, B. Wu, J. Liang, S. Cheng, and G. Zhang, "Posit arithmetic hardware implementations with the minimum cost divider and squareroot," *Electronics*, vol. 9, no. 10, 2020.
- [26] S. Bansal, H. Hsiao, T. Czajkowski, and J. H. Anderson, "High-level synthesis of software-customizable floating-point cores," in *Design, Automation & Test in Europe*. IEEE, 2018, pp. 37–42.
- [27] D. Thomas, "Templatized soft floating-point for high-level synthesis," in *27th Annual International Symposium on Field-Programmable Custom Computing Machines*. IEEE, 2019.
- [28] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélicier, and P. Zimmermann, "MPFR: A multiple-precision binary floating-point library with correct rounding," *ACM Transactions on Mathematical Software*, vol. 33, no. 2, p. 13, 2007.