

# Unified Posit/IEEE-754 Vector MAC Unit for Transprecision Computing

Luís Crespo, Pedro Tomás, *Senior Member, IEEE*,  
Nuno Roma, *Senior Member, IEEE*, and Nuno Neves, *Member, IEEE*

**Abstract**—Transprecision computing targets energy-efficiency with multiple floating-point modules with different precisions to suit application requirements. Variable-precision architectures aim at making a more efficient hardware resource utilization, but they often rely on the IEEE-754 standard, without low-precision arithmetic support. Alternatively, the Posit format is particularly well-suited for low-precision arithmetic. However, for higher precisions, hardware requirements become prohibitive. Accordingly, this paper proposes a new unified Posit/IEEE-754 Vector Multiply-Accumulate (VMAC) unit, comprising a vectorized variable-precision datapath with shared support for the Posit and IEEE-754 formats. A 28nm ASIC implementation resulted in 50% less area and 2.9× less power consumption than typical transprecision setups.

**Index Terms**—Floating-point arithmetic, Posit, IEEE-754, Variable-Precision, SIMD

## I. INTRODUCTION

TRANSPRECISION computing [1] has received a gradually increasing attention as a viable paradigm to cope with ever increasing performance and energy efficiency demands in modern computing systems. It is set on the principle that different applications have different precision requirements (e.g., while some physics simulations require higher than 64-bit precisions [2], deep learning applications sustain lower precisions with as little as 4 bits [3]), and that, as recent studies have shown [4–10], by lowering floating-point (FP) precision it is possible to gain straightforward acceleration and efficiency.

However, most transprecision hardware solutions [11] attempt to support different precisions by instantiating multiple arithmetic modules. This leads to an increased chip area and a waste of resources [12]. To tackle this issue, recent variable-precision arithmetic units [12–14] introduce dynamic datapaths that can operate in different precisions with the same hardware resources. To do so, they deploy a higher precision arithmetic logic (e.g., 32-bit) and allow parts of the circuit to be turned off to lower the operand precision (e.g., to 8-bit or lower [3]). While this approach provides for significant area reductions and enables straightforward Single-Instruction

Multiple-Data (SIMD) capabilities [12], existing solutions are often limited by their adoption of the IEEE-754 standard [13], whose lowest supported precision is only 16 bits.

Alternatively, some recent solutions [12, 14] adopt the Posit format [4], mainly since it allows parameterizable *precision* and dynamic range (*exponent size*). The Posit format is also interesting for fused operations, since it adopts an exact accumulator structure (quire) with enough precision to avoid overflow and accuracy losses [15]. While Posit-based implementations traditionally define and fix its parameters at design-time [8, 9, 16–18], it has been shown that it is possible to support runtime-configurable exponent sizes with minimal hardware overheads [10]. This allows making use of the entire representable dynamic range for a given posit precision by specifying the exponent size of the input values. In turn, it also provides the possibility to encode a larger dynamic range, capable of supporting (within the same hardware) both values with high decimal precisions and very large magnitude.

Nevertheless, while these features make posits well suited for low-precision and transprecision computing, the overheads associated with the quire becomes prohibitive when the precision increases [16, 17]. Hence, for a general-purpose context, it is desirable to maintain compatibility with the standard IEEE-754 format, as it still is the most established FP format.

This paper proposes a new Posit/IEEE-754 Vector Multiply-Accumulate (VMAC) unit for transprecision computing. Besides combining variable-precision arithmetic and SIMD capabilities, it takes a step further from existing solutions by deploying a unified support for the IEEE-754 and Posit formats. It introduces the following contributions and features:

- an efficient variable-precision FP multiply-accumulate (MAC) 32-bit architecture for transprecision computing;
- a unified FP arithmetic architecture compatible with both the IEEE-754 and the Posit formats with support for inter-format operation and conversion, which is also compatible with the existing RISC-V Vector (RVV) [19] and recently proposed RISC-V Posit extensions [20, 21];
- a fully vectorized datapath to efficiently make use of the released hardware resources in low-precision scenarios;
- SIMD decoding/encoding modules with shared support for FP vectors encoded with *i)* posit formats with configurable exponent size; *ii)* IEEE-754 standard and low-precision non-standard formats; and *iii)* multiple scalar/vector element precisions.

Finally, when implemented in a 28nm ASIC technology,

This work was partially supported by national funds through Fundação para a Ciência e a Tecnologia (FCT) under projects UIDB/50021/2020 and PTDC/EEI-HAC/30485/2017, and by funds from the European Union Horizon 2020 Research and Innovation programme under grant agreement No. 101036168.

All authors are with INESC-ID and Instituto Superior Técnico, Universidade de Lisboa, in Lisbon, Portugal (e-mail: luis.miguel.crespo@tecnico.ulisboa.pt, pedro.tomas@inesc-id.pt, nuno.roma@inesc-id.pt, nuno.neves@inesc-id.pt).

the proposed VMAC results in 50% less area and 2.9× less power to achieve the same multiple-precision functionality when compared with typical transprecision architecture topologies [11], while supporting a unique unified FP format.

## II. BACKGROUND

### A. IEEE-754 Standard

The IEEE-754 standard defines a FP number with sign (S), biased exponent (E) and mantissa (M), with value:

$$(-1)^S \times 2^{E-bias} \times 1.M, \quad (1)$$

where *bias* is the exponent bias value. Although the standard defines formats for half- (16-bit), single- (32-bit), and double-precision (64-bit), it does not define a low-precision 8-bit format. However, since the proposed architecture supports 8-bit posits, for comparison purposes, an 8-bit floating point format is adopted with 4 exponent bits and 3 mantissa bits.

### B. Posit Number System

The posit number system is defined by the pair  $\langle n, es \rangle$ , where  $n$  represents the word size (*precision*) and  $es$  is the maximum *exponent* size. Eq. 2 depicts the Posit encoding:

$$\underbrace{\begin{array}{c} \text{sign} \\ s \end{array}} \underbrace{\begin{array}{c} \text{regime} \\ r r \dots \bar{r} \end{array}} \underbrace{\begin{array}{c} \text{exponent} \\ e_0 e_1 \dots e_{es-1} \end{array}} \underbrace{\begin{array}{c} \text{fraction} \\ f_0 f_1 f_2 \dots \end{array}} \quad (2)$$

*n bits*

Similarly to floats, posits include the sign, exponent, and fraction, with an additional field called regime. Contrarily to floats, whenever the sign bit corresponds to a negative number, it is necessary to take the 2's complement before decoding the remaining fields. The regime is a variable-sized field, whose encoded value ( $k$ ) is given by the run-length of '0' or '1' bits.

Together with the exponent field, the  $k$  encoded value in the regime represents a scale factor of the represented value, equivalent to the exponent in floats. As a consequence of the variable-sized regime, the exponent and fraction contents are unknown before decoding the regime. In fact, depending on the run length, they can be partly (or fully) left out of the binary encoding. Hence, a posit number value is given by:

$$(-1)^{sign} \times 2^{exp+k2^{es}} \times 1.f \quad (3)$$

The Posit format has a single encoding for zero (000...0) and a single Not-a-Real (NaN) mathematical exception (100...0).

Additionally, it makes use of a 2's complement fixed-point accumulator (*quire*) based on the Kulisch accumulator, used to store sums of products of posits without rounding and accuracy loss. Naturally, the quire has a considerable hardware overhead. It is composed by 4 fields: sign, carry guard (cg), integer (int) and fraction (frac); and its size is given by:

$$quire\ size = 1 + cg + 2^{es+2} \times (n - 2) \quad (4)$$

Hence, the quire must be carefully dimensioned as it grows exponentially with the exponent size and precision [17].

## III. POSIT/IEEE-754 VMAC ARCHITECTURE

### A. Overview

The proposed VMAC architecture takes a step further from existing multiple-precision arithmetic units, not only by combining variable-precision arithmetic and dynamic vectorization capabilities, but also by providing an unified support for the Posit and IEEE-754 FP formats. Accordingly, it features:

**1 Posit-based Variable-Precision Structure:** All modules of the 32-bit posit fused MAC datapath are designed to easily allow an adaptation of their arithmetic precision (at runtime), supporting 32, 16, or 8-bit operations (as illustrated in Fig. 1.A). To mitigate the hardware overheads associated with the quire, the proposed unit only provides an exact accumulation for low-precision scenarios with standard 8-bit posits ( $es = 2$ ) [15], by using a quire of 128-bits (as opposed to 512 bits for 32-bit posit accumulation). Hence, a scale factor value is paired with the quire to ensure the correct representation of the accumulations for all the supported precisions.

**2 Dynamic Vectorization:** All arithmetic modules are fully vectorized and configurable at runtime to support 1x32-bit, 2x16-bit, and 4x8-bit vector operations using the same hardware (see Fig. 1.B). Hence, the resources released when precision is reduced provide support for parallel computations, offering increased throughput. To support vectorization, the 32-bit input vectors are decoded into three unified vector formats that gather the sign ( $s$ ), scaling factor (or exponent -  $sf$ ), and fraction ( $f$ ) of each vector element, according to the  $(-1)^s \times 2^{sf} \times 1.f$  generic exponential format (see Fig. 1.C).

**3 Variable-Exponent Posit Configuration:** Posit exponent size can be defined at runtime (instead of being fixed at design-time), allowing most of the dynamic range for a given precision to be representable. Since the quire is already paired with a scaling factor (see Fig. 1.C), the arithmetic logic can support dynamic ranges larger than those that can be represented by the quire. Accordingly, it is only necessary to include a set of shifters to decode/encode the posit format according to the configured exponent size (described below).

**4 FP Format Unification:** While the Posit and IEEE-754 formats are fundamentally different in their representation, after decoded, both represent a FP number in the generic exponential format. As such, the logic to perform multiplica-

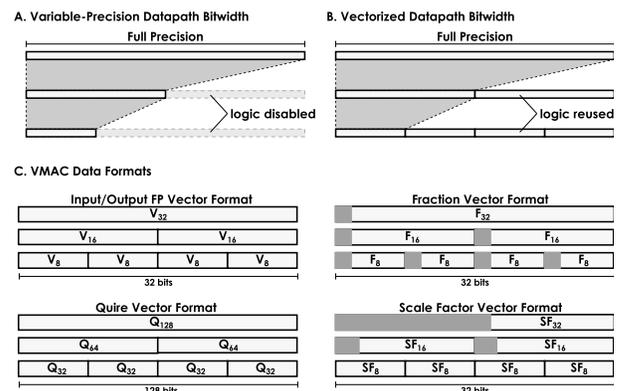


Fig. 1. Proposed VMAC (A) variable-precision and (B) vector datapath configuration schemes, together with the (C) encoded/decoded FP and quire vector data formats. Grey areas represent unused bits.

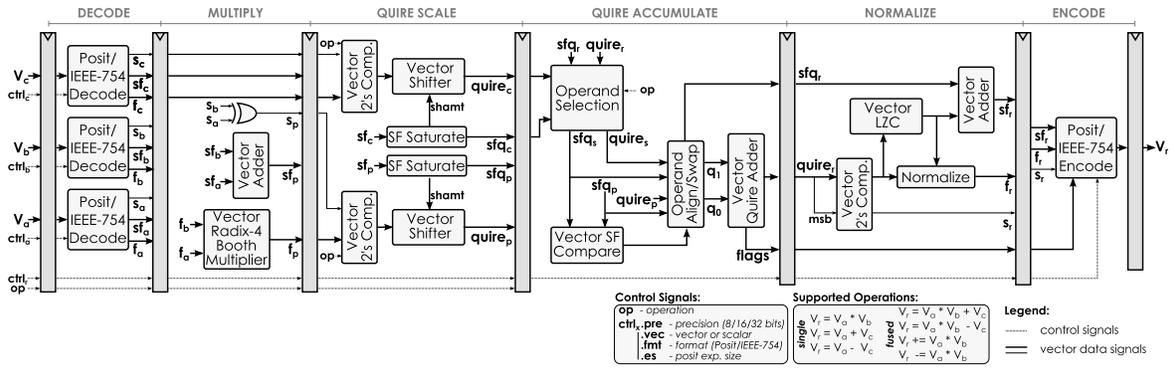


Fig. 2. Proposed Posit/IEEE-754 VMAC unit architecture diagram.

tion and addition/subtraction is virtually the same. Conversely, to add IEEE-754 support in a Posit base architecture, it is only necessary to add minimal decoding/encoding logic and detection for mathematical exceptions (nonexistent in the Posit format). For 8-bit precision operations, the 8-bit minifloat variant was adopted to match the equivalent Posit precision.

**5 Inter-format Operations and Conversions:** The introduced unified FP format allows the proposed unit to perform inter-format operations between equivalent Posit and IEEE-754 precisions. Since the unit's internal representation is compatible with both formats, it is only necessary to decode each operand according to their specific format (controlled by dedicated configuration signals - see below). Similarly, the format of the output can also be configured independently of the input formats, enabling straightforward format conversions.

### B. Proposed Architecture

The proposed VMAC unit (depicted in Fig. 2) comprises a fully pipelined architecture, supporting vector variable-precision FP addition, subtraction, and multiplication, together with fused multiply-add and multiply-accumulate operations. The unit deploys a 32-bit SIMD datapath with unified support for Posit and IEEE-754 FP formats, implemented by a 6-stage pipeline : i) Decode; ii) Multiply; iii) Quire Scale; iv) Quire Accumulate; v) Normalize; and vi) Encode. The unit accepts three input vector operands ( $V_a$ ,  $V_b$  and  $V_c$ ), and outputs one result vector ( $V_r$ ), and is capable of operating with 32/16/8-bit scalar values or with 2x16/4x8-bit vectors.

The following paragraphs detail each of the pipeline stages.

**Unified Decode:** The Decode stage comprises three equivalent vectorized decoding modules (one for each input value - see Fig. 3.A), each containing the necessary logic to decode either the Posit and IEEE-754 formats, to their  $s$ ,  $sf$ , and  $f$  fields. The FP format and precision are selected according to a set of control signals paired with the input value (see Fig. 2).

For the IEEE-754 format, the three fields are extracted and a bias is subtracted from the exponent value, according to the configured precision. Conversely, for the Posit format, the 2's complement is applied to the input value according to the sign bit. Next, the regime run-length is decoded by means of a leading zero counter (LZC) (if it starts with '1' the value is first inverted). Then,  $k$  is calculated and the regime is left-shifted out according to the zero count, leaving the exponent and fraction. This value is then shifted again by  $es$ , to split the

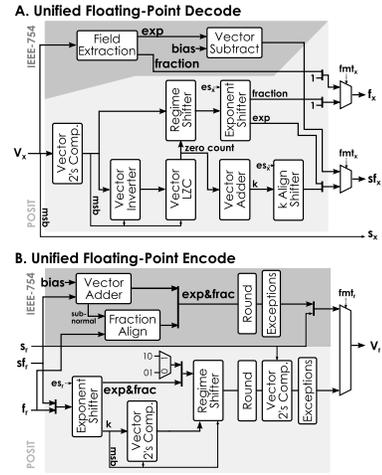


Fig. 3. Unified FP (A) decoding and (B) encoding modules, providing simultaneous support for Posit and IEEE-754 vector formats.

exponent and the fraction. The  $k$  value is then shift-aligned according to  $es$  and concatenated with the exponent to obtain  $sf$ . Finally, a '1' bit is added to the fraction to obtain  $f$ .

**Multiplication:** The Multiply stage implements a variable-precision vector FP multiplier, operating the decoded  $V_a$  and  $V_b$  values while propagating  $V_c$  to the next stage. To do so, the product of the fractions is performed with a 4x4 structure of 8-bit radix-4 Booth multipliers, generating 16 partial products in carry-save format, accumulated in a Wallace tree-like structure, resulting in a 64-bit value. Variable-precision and/or vectorization are applied by only enabling the required multipliers. The scale factor vectors are added in a vectorized carry-lookahead adder, capable of breaking its carry-chain to perform lower-precision parallel additions. The sign vector results from a bitwise XOR of the input vectors.

**Quire Scale and Accumulate:** To mitigate the critical path associated with the quire processing structure, the Quire module is subdivided into two pipeline stages: Scale and Accumulate. The operands ( $V_c$  and the product from the previous stage) are first converted to a 128-bit fixed-point quire format vector (in accordance with the features discussed in Section III-A), paired with a scale factor vector (see Fig. 1.C).

The conversion is done in the Scale stage, by first taking the 2's complement of the fraction vectors and sign-extending them according to the precision. Next, the fraction is aligned to the quire fixed-point format, with a vectorized barrel shifter.

TABLE I  
COMPARISON OF THE PROPOSED VMAC WITH THE STATE-OF-THE-ART.

UNIT	NUM. BITS	PIPEL. STAGES	ASIC TECH.	DELAY (ns)	AREA ( $\mu\text{m}^2$ )	POWER (mW)	PERF. (GOPS)	AREA EFF. ( $\times 10^{-6}$ GOPS/ $\mu\text{m}^2$ )	EDP ( $\times 10^{-22}$ J.s)
Ref. Posit Std. MAC	8	5	28 nm	0.65	7598	21	1.54	202.4	0.89
Ref. Posit Std. MAC	16	5	28 nm	0.8	17384	47	1.25	71.91	3.01
Ref. Posit Std. MAC	32	5	28 nm	0.91	39767	108	1.10	27.63	8.94
<b>Proposed VMAC</b>	<b>8/16/32</b>	<b>6</b>	<b>28 nm</b>	<b>1.5</b>	<b>51563</b>	<b>99</b>	<b>2.7/1.3/0.7</b>	<b>51.7/25.7/12.9</b>	<b>5.6/11.1/22.3</b>
Posit DFMA [10]	32	5	45 nm	1.5	112350	370	0.67	5.95	83.25
FP VFMA [13]	16/32/64	3	90 nm	1.5	180610	44	2.7/1.3/0.7	14.8/7.4/3.7	2.5/4.9/9.9
Posit VMULT [14]	8/16/32	-	90 nm	2.3	91861	64	1.7/0.9/0.4	18.9/9.5/4.7	8.5/16.9/33.9

Given that the quire size is limited to constrain hardware resources, the shifting amount is calculated from the scale factor dynamic range with the aid of a *saturation module* (see Fig. 2), which saturates the shifting amount (adjusting the scale factor accordingly) whenever the fixed-point value overflows.

The *Accumulate* stage is responsible for implementing the quire arithmetic logic, either by adding the values obtained from the previous stage, or by accumulating the saved quire value with one of such values. As such, the required operands are first selected and then aligned according to their scale factors. This is done by typical FP alignment logic, with the aid of a vectorized right barrel shifter, while any discarded bits are condensed in a sticky vector. Finally, the mathematical exception flags are generated (both for the Posit and IEEE-754 formats), and the result is saved in the pipeline registers.

**Normalization:** The *Normalize* stage is responsible for re-normalizing the quire and extracting the *s*, *sf*, and *f* vectors. First, the sign vector is extracted from the MSB of each quire vector element, which allows converting the quire to an unsigned value. Next, the number of shift positions required to normalize the quire is obtained with a vectorized LZC. The obtained zero count is used by a vectorized left shifter to align the unsigned quire vector. Any discarded bits are condensed in a sticky vector. Finally, the scale factor is obtained by adding the quire scale factor and the obtained zero count.

**Unified Encode:** Finally, the *Encode* stage provides the necessary logic for encoding the output vectors to Posit and IEEE-754 formats (see Fig. 3.B). The logic is fully vectorized and translates the *s*, *sf*, and *f* vectors of the result to the selected FP format vector. For the IEEE-754 format, the bias is added to the scale factor (according to the precision) and the resulting value is verified, adjusting the fraction for subnormal numbers. Afterwards, the fields of each vector element are concatenated and the fraction is rounded. The output result is selected between the rounded result, zero, infinity or canonical NaN, according to the flags generated by previous stages.

For the Posit format, *sf* and *f* are first concatenated and then right shifted, according to *es*, to obtain *k*. The *k* value's 2's complement is taken and the regime is shifted-in to *sf* and *f*, according to *k*'s sign. The resulting binary value is then rounded and the 2's complement is taken according to *s*.

#### IV. IMPLEMENTATION RESULTS

The proposed Unified Posit/IEEE-754 VMAC unit was described in RTL and synthesized for 28nm UMC standard cell technology, targeting an operating frequency of 667 MHz, under typical operating conditions (1.05 V, 25° C). Chip area

TABLE II  
AREA BREAKDOWN FOR THE PROPOSED VMAC AND ITS COMPONENTS.

	Pipeline Stage	Area ( $\mu\text{m}^2$ )	Power (mW)
<b>Decode</b>	Posit	2430	4.6
	IEEE-754	263	0.8
<b>Multiply</b>		15917	29.5
<b>Quire</b>		13711	27.9
<b>Normalize</b>		6418	11.9
<b>Encode</b>	Posit	3515	6.3
	IEEE-754	2340	4.6
<b>Total</b>	<b>VMAC</b>	<b>51563</b>	<b>99</b>

and power estimation results were obtained with Cadence Genus 19.11 and presented in Tables I and II. Energy efficiency was calculated using the energy-delay product (EDP), by considering energy consumption and latency. Its operation was validated with testing vectors generated with Sigmoid Numbers julia library [22] and TestFloat [23].

To establish reference architectures, three fixed Posit fused multiply-accumulate (MAC) architectures with 8, 16, and 32 bits were also implemented, with fixed 2-bit exponent size, as per the latest Posit standard [15]. These precisions imply the use of 128, 256, 512-bit quires, respectively. The proposed design was also compared with state-of-the-art variable-precision units, including a 64-bit IEEE-754 variable-precision fused multiply-add (FMA) [13] (VFMA), a 32-bit Posit variable-precision multiplier [14] (VMULT), and a 32-bit Posit dynamic FMA [10] with configurable exponent size (DFMA).

When compared with the reference 32-bit Posit MAC unit, it is observed that despite the introduced variable-precision and unified FP, the proposed VMAC only presents a 30% chip area increase, while showcasing a similar power consumption. In fact, the VMAC is more area- and energy-efficient (see Table I) since it reuses hardware resources to deploy variable-precision vectorization, increasing throughput. Also, although a higher latency was expected due to the increased complexity, the critical path is still majorly mitigated by limiting the size of the VMAC quire to 128 bits (as opposed to the reference 512-bit quire). This is also evident when comparing it with the DFMA [10], which also adopts a 512-bit quire. Additionally, as opposed to the VMAC, the DFMA [10] presents a fixed-precision datapath, unsuited for transprecision computation.

Since the state-of-the-art variable-precision solutions were implemented with distinct technology processes (90nm vs. 28nm), scaled area and power estimations were obtained for 28nm technology with the DeepScaleTool [24]. When considering the estimated results (see Table III), it is observed that although the proposed VMAC presents an increased resource utilization when compared to the VFMA [13] (IEEE-754), the difference is easily explained by two main factors: *i*) the

TABLE III

COMPARISON OF THE PROPOSED VMAC WITH THE STATE-OF-THE-ART (SCALED TO 28 nm TECHNOLOGY WITH DEEPSCALETOOL [24]).

UNIT	NUM. BITS	PIPEL. STAGES	RESULT SOURCE	DELAY (ns)	AREA ( $\mu\text{m}^2$ )	POWER (mW)
<b>Proposed VMAC</b>	<b>8/16/32</b>	<b>6</b>	<b>Synthesis</b>	<b>1.5</b>	<b>51563</b>	<b>99</b>
Posit DFMA [10]	32	5	Estimated	1.24	39324	266
FP VFMA [13]	16/32/64	3	Estimated	0.77	16044	21
Posit VMULT [14]	8/16/32	-	Estimated	1.18	8160	31

Posit decoding, encoding and quire modules account for more than 50% of the total area and power of the VMAC (see Table II); and *ii*) the VFMA [13] unit is not fully pipelined, requiring much less area that is often related to registers. The latter is also true when comparing with VMULT [14], by only considering the multiplier stage of the VMAC (see Tables III and II). However, the VMAC presents a much higher functionality than previous solutions.

In particular, while the VFMA [13] presents a variable-precision architecture, with SIMD capabilities, it is bound by its sole adoption of the IEEE-754 format, and cannot perform 8-bit low-precision operations. Contrarily, while also providing IEEE-754 support, the proposed VMAC leverages the Posit format to perform low-precision operations with a configurable dynamic range. Hence, the VMAC shows a much higher flexibility and is better suited for low-precision computation scenarios. Conversely, while the more recent VMULT [14] presents variable and low-precision capabilities (similar to the proposed VMAC), it only implements the multiplier datapath and lacks the same flexibility of the VMAC in what concerns the configurable exponent size and IEEE-754 format support.

Furthermore, the benefits of the proposed VMAC are also evidenced when considering throughput and energy efficiency. Despite the implicit logic increase, necessary to implement its unified format and variable-precision datapath, the VMAC still presents similar throughput and energy efficient when compared to the state-of-the-art solutions (see Table I).

Finally, when considering the integration of the proposed VMAC in a typical transprecision architecture [11] to support multiple-precision datapaths, it is estimated to require 50% less area and  $2.9\times$  less power than a combination of the considered Posit MAC reference units with the same precision mix (i.e.,  $4\times 8\text{-bit} + 2\times 16\text{-bit} + 1\times 32\text{-bit}$  MACs). Moreover, the VMAC offers increased flexibility, by supporting a unified FP format.

## V. CONCLUSIONS

This paper proposes a new unified Posit/IEEE-754 Vector Multiply-Accumulate (VMAC) unit architecture for transprecision computing. It offers a variable-precision datapath with SIMD capabilities with a unified support for the Posit and IEEE-754 FP standards. Accordingly, it is capable of performing low- and high-precision Posit operations (with dynamic exponent size), while maintaining compatibility with the standard IEEE-754 format. A 28nm ASIC implementation resulted in 50% less area and  $2.9\times$  less power when compared with a typical transprecision system topology. In the future, we will consider deploying the proposed VMAC in transprecision acceleration platforms and the development of RISC-V extensions to support its inclusion in existing processors.

## REFERENCES

- [1] A. C. I. Malossi *et al.*, "The transprecision computing paradigm: Concept, design, and applications," in *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1105–1110, IEEE, 2018.
- [2] M. Klöwer, P. D. Düben, and T. N. Palmer, "Posits as an alternative to floats for weather and climate models," in *Proceedings of the Conference for Next Generation Arithmetic 2019*, pp. 1–8, 2019.
- [3] X. Sun, N. Wang, C.-Y. Chen, J. Ni, A. Agrawal, X. Cui, S. Venkataramani, K. El Maghraoui, V. V. Srinivasan, and K. Gopalakrishnan, "Ultra-low precision 4-bit training of deep neural networks," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [4] 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.
- [5] G. Raposo, P. Tomás, and N. Roma, "Positnn: Training Deep Neural Networks with Mixed Low-Precision Posit," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7908–7912, IEEE, 2021.
- [6] N. Wang, J. Choi, D. Brand, C.-Y. Chen, and K. Gopalakrishnan, "Training deep neural networks with 8-bit floating point numbers," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 7686–7695, 2018.
- [7] M. K. Jaiswal and H. K.-H. So, "Pacogen: A hardware posit arithmetic core generator," *IEEE Access*, vol. 7, pp. 74586–74601, 2019.
- [8] Z. Carmichael, H. F. Langroudi, C. Khazanov, J. Lillie, J. L. Gustafson, and D. Kudithipudi, "Deep positron: A deep neural network using the posit number system," in *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1421–1426, IEEE, 2019.
- [9] H. Zhang, J. He, and S.-B. Ko, "Efficient posit multiply-accumulate unit generator for deep learning applications," in *2019 IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1–5, IEEE, 2019.
- [10] N. Neves, P. Tomás, and N. Roma, "Dynamic Fused Multiply-Accumulate Posit Unit with Variable Exponent Size for Low-Precision DSP Applications," in *2020 IEEE Workshop on Signal Processing Systems (SiPS)*, pp. 1–6, IEEE, 2020.
- [11] G. Tagliavini, S. Mach, D. Rossi, A. Marongiu, and L. Benini, "A transprecision floating-point platform for ultra-low power computing," in *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1051–1056, IEEE, 2018.
- [12] N. Neves, P. Tomás, and N. Roma, "Reconfigurable Stream-based Tensor Unit with Variable-Precision Posit Arithmetic," in *2020 IEEE 31st International Conference on Application-specific Systems, Architectures and Processors (ASAP)*, pp. 149–156, IEEE, 2020.
- [13] H. Zhang, D. Chen, and S.-B. Ko, "Efficient multiple-precision floating-point fused multiply-add with mixed-precision support," *IEEE Transactions on Computers*, vol. 68, no. 7, pp. 1035–1048, 2019.
- [14] H. Zhang and S.-B. Ko, "Efficient multiple-precision posit multiplier," in *2021 IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1–5, IEEE, 2021.
- [15] P. W. Group, "Posit Standard Documentation," *Rel. 4.12-draft*, Jul. 2021.
- [16] 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 2019*, pp. 1–10, 2019.
- [17] F. d. D. Luc Forget, Yohann Uguen, "Hardware cost evaluation of the posit number system," in *Compas'2019 - Conférence d'informatique en Parallélisme, Architecture et Système*, pp. 1–7, Jun 2019.
- [18] S. Jean, A. Raveendran, A. D. Selvakumar, G. Kaur, S. G. Dharani, S. G. Pattanshetty, and V. Desalphine, "P-fma: A novel parameterized posit fused multiply-accumulate arithmetic processor," in *2021 34th International Conference on VLSI Design and 2021 20th International Conference on Embedded Systems (VLSID)*, pp. 282–287, IEEE, 2021.
- [19] A. Waterman and K. Asanovic, "RISC-V 'V' Vector Extension," 2019.
- [20] S. Tiwari, N. Gala, C. Rebeiro, and V. Kamakoti, "Peri: A configurable posit enabled risc-v core," *ACM Transactions on Architecture and Code Optimization (TACO)*, vol. 18, no. 3, pp. 1–26, 2021.
- [21] R. Jain, N. Sharma, F. Merchant, S. Patkar, and R. Leupers, "Clarinet: A risc-v based framework for posit arithmetic empiricism," *arXiv preprint arXiv:2006.00364*, 2020.
- [22] I. Yonemoto, "'sigmoid numbers'," [Online]. <https://github.com/interplanetary-robot/SigmoidNumbers>, 2018.
- [23] J. Hauser, "Berkeley testfloat," [Online]. Available: <http://www.jhauser.us/arithmetic/TestFloat.html>, 2018.
- [24] S. Sarangi and B. Baas, "Deepscaletool: A tool for the accurate estimation of technology scaling in the deep-submicron era," in *2021 IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1–5, IEEE, 2021.