

FPPU: Design and Implementation of a Pipelined Full Posit Processing Unit

Federico Rossi¹, Francesco Urbani¹, Marco Cococcioni¹, Emanuele Ruffaldi², Sergio Saponara¹

Abstract—By exploiting the modular RISC-V ISA this paper presents the customization of instruction set with posit[™] arithmetic instructions to provide improved numerical accuracy, well-defined behavior and increased range of representable numbers while keeping the flexibility and benefits of open-source ISA, like no licensing and royalty fee and community development. In this work we present the design, implementation and integration into the low-power Ibex RISC-V core of a full posit processing unit capable to directly implement in hardware the four arithmetic operations (add, sub, mul, div and fma), the inversion, the float-to-posit and posit-to-float conversions. We evaluate speed, power and area of this unit (that we have called Full Posit Processing Unit). The FPPU has been prototyped on Alveo and Kintex FPGAs, and its impact on the metrics of the full-RISC-V core have been evaluated, showing that we can provide real number processing capabilities to the mentioned core with an increase in area limited to 7% for 8-bit posits and to 15% for 16-bit posits. Finally we present tests on the use of posits for deep neural networks with different network models and datasets, showing minimal drop in accuracy when using 16-bit posits instead of 32-bit IEEE floats.

Index Terms—Posits, RISC-V, Accelerator, Arithmetic

I. INTRODUCTION

THE RISC-V (Reduced Instruction Set Computer - V) open instruction set architecture (ISA) [1], [2] is a modern, open-source ISA that is gaining popularity in recent years. RISC-V was designed to be a free and open alternative to proprietary ISAs such as ARM and x86, with the goal of fostering innovation and research in the computer architecture community. Unlike proprietary ISAs, RISC-V is not tied to any specific vendor or implementation and can be used in a wide range of devices, from microcontrollers to supercomputers.

Posit[™] numbers are a recent alternative to traditional floating-point arithmetic (IEEE 754 floating point numbers: binary32, binary64, ...). The posit number system, introduced in 2017 by J. L. Gustafson [3] and standardized in 2022 [4], is designed to provide improved numerical accuracy and range, thus maintaining a similar level of performance as traditional floating-point numbers but for reduced data bit-width. Posit numbers are represented by a fixed number of bits, just like floating-point numbers, but they use a different encoding scheme that allows for a larger range of representable numbers.

Posit numbers also provide a well-defined and predictable behavior when dealing with numbers close to zero. Multiple works proved the capabilities of Posit number to be a drop-in replacement of binary32 numbers for Deep Neural Networks (DNNs) [5]–[8]. In particular 8-bit posits demonstrated to be able to maintain similar accuracy in DNN tasks when compared to binary32 numbers and 16-bit posits can even outperform binary32. In this work we advance the combination of the RISC-V open instruction set and posit arithmetic to provide an efficient way to process real numbers in RISC-V cores already seen in [9], adding this capability to the small Ibex RISC-V core [10] without using a Floating Point Unit (FPU). With respect to the Light PPU [9] (which supported only float-to/from-posit conversions and it was essentially used to save storage of DNN weights), the Full Posit Processing Unit (FPPU) proposed in this work also provides hardware support to the four posit operations (sub, add, mul, div) and to the computation of reciprocals. This means that, beside storage, also computations in posit domain become fast, being native and not simulated. The RISC-V ISA is highly modular and customizable, making it well-suited for the integration of a posit processing unit. The modular nature of RISC-V allows to include only the instructions and functional units that are needed for posit arithmetic, without the need to include instructions or functional units that are not required (e.g. IEEE FPUs and related RV32F ISA). This can result in a smaller and more efficient processing core, as already proved in [11]–[13]. On the other hand, posit arithmetic can provide improved numerical accuracy, well-defined behavior, and increased range of representable numbers, making it well-suited for applications that require these features such as embedded systems, large data sets, and scientific computing.

A. Organization of the paper

Section II summarise differences and similarities of this work with other related researches. Section III describes posit numbers and outlines main properties. Section IV elaborates on the logical implementation of posit operations, focusing on the different aspects of division algorithms. Section V shows the design of the FPPU. Section VI describes the Instruction Set Architecture (ISA) extension for posit operations and the SW support for compilation. Section VII shows the design steps for the integration of the FPPU inside the low-power Ibex core. Section VIII characterizes the FPPU and Ibex components for area occupation, clock and power consumption. Discussions and Conclusions are given in the last section. The IP database for the posit unit is completely contained at: https://github.com/federicorossifr/ppu_public.

Work partially supported by H2020 projects EPI-SGA2 (grant no. 101036168, <https://www.european-processor-initiative.eu/>) and TEXTAROSSA (grant no. 956831, <https://textarossa.eu/>).

¹Federico Rossi, Francesco Urbani, Marco Cococcioni, Sergio Saponara are with University of Pisa, Department of Information Engineering, Pisa, Italy (e-mail: federico.rossi@ing.unipi.it, [name.surname]@unipi.it).

²Emanuele Ruffaldi is with Medical Microinstruments Inc., Italy (e-mail: emanuele.ruffaldi@mmimicro.com).

TABLE I
SIMILARITIES AND DIFFERENCES BETWEEN THIS WORK AND RELATED POSIT PROCESSING UNITS

HARDWARE STACK				
	This work	Percival [14]	Clarinet [12]	Pacogen [15]
Configurable posit size	✓	×	✓	×
High-precision posit support (32,64 bit)	✓	✓	✓	✓
Low-precision posit support (8,16 bit)	✓	×	✓	✓
Quire/Fused support (FMA)	✓	✓	✓	×
Dynamic power monitored	✓ (see [16])	×	×	×
RISC-V Integration	✓	✓	✓	×
Platform and soft-core independent	✓	×	✓	✓
SIMD/Vector operations	✓	×	×	×
SOFTWARE STACK				
High-level software support (e.g. posit library)	✓	×	×	×
Compiler independent instruction support	✓	×	N/A	N/A
Integration with Deep Neural Networks (DNNs) frameworks	✓	×	×	×

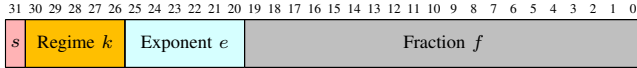


Fig. 1. Bit fields of a posit(32, 6) data type.

II. RELATED WORKS

There exists several hardware implementation of posit processing units, each one with different peculiarities in terms of configurable posit size, pipeline architecture, implemented operations, resource utilisation, RISC-V integration and software support. We summarise some of them, comparing with this work in Table I. We highlight the hardware characteristics in terms of flexibility of the posit configuration (e.g. number of bits and number of exponent bits), support for fused operations, vectorized operations and capabilities to integrate the unit inside a RISC-V soft-core. Furthermore we elaborate on the software stack, highlighting whether the unit has an high-level support in the software domain via libraries, compilers and computational libraries.

III. POSIT NUMBERS

A posit number (see Fig. 1) is represented by a signed integer on 2's complement. It can be configured with the total number of bits N and maximum number of exponent bits ES . We define such a posit as Posit $\langle N, ES \rangle$. The format can have at most four fields: i) sign on 1-bit, ii) regime with a variable size (run-length encoded), iii) exponent with at most ES bits and iv) fraction with a variable length. An example of a posit number instance is shown in Fig. 2. Note that, if the regime fields is large enough, it is possible that the exponent field has less bits available than ES . In this case the actual exponent value is computed by padding zeroes to the right of the exponent bits in the format.

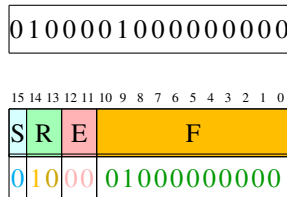


Fig. 2. An example of Posit configuration with 16 bits and 2 exponent bits. The associated real value to the shown Posit is: $+16^0 \times 2^0 \times (1 + 512/2048) = 1.25$.

The length l of the regime corresponds to the number of identical bits following the sign bit:

$$s, \underbrace{b_1, b_2, b_3, \dots, b_l}_{=b}, \underbrace{b_{l+1}, \dots}_{=\bar{b}} \quad (1)$$

The regime length is l . Depending on the value of the bit b , the regime value k will be computed as follows:

$$k := \begin{cases} l - 1 & \text{if } b = 1 \\ -l & \text{if } b = 0 \end{cases} \quad (2)$$

The regime value is a scale factor for a special constant, that depends on the posit configuration, called *useed*. The *useed* value is computed as follows:

$$\text{useed} := 2^{2^{ES}} \quad (3)$$

Hence, the real value r associated with a posit represented by the integer P on two's complement (with sign s) is computed as in Equation (4):

$$r := \begin{cases} 0 & \text{if } P = 0 \\ \text{NaN (Not a Real)} & \text{if } P = -2^{N-1} \\ (-1)^s \times \text{useed}^k \times 2^e \times \left(1 + \frac{f}{2^F}\right) & \text{otherwise} \end{cases} \quad (4)$$

The value F is the length of the fraction field. Note that there will be always an implicit one in front of the fraction (i.e. $1.f_1, f_2, \dots, f_F$), without any subnormal number differently from IEEE binary32 numbers.

In [4] there was reported an alternative version of the third expression of Equation (4), which is:

$$r = (1 - 3s + f/2^F) \times 2^{(1-2s) \times (2^{ES} \times k + e + s)} \quad (5)$$

IV. POSIT ARITHMETIC OPERATIONS AND IMPLEMENTATION

This Section elaborates on the logic implementation of posit arithmetic operations, deepening on the arithmetical correctness and numerical rounding of such operations. Each posit operation described here starts with a *decoding* and *input conditioning* phase, where the posit operands are decoded into the *sign*, *regime*, *exponent*, *significand* fields and decisions are made depending on few special cases (for example when

one of the operands is equal to 0 or NaR). After this phase, each posit operand is transformed into a general Floating-point Intermediate Format (FIR) $\langle s_f, te, f \rangle$, where s_f is the sign, te is the total exponent (without bias) and f is the fractional part of the significand (**note: the FIR representation is just a zero-cost transformation of a decoded posit and do not share any similarity with the IEEE 32-bit floating point (e.g. no NaN, Infs or subnormals):**).

$$p = (-1)^{s_f} \times 2^{te} \times (1.f).$$

While s_f and f are extracted from the posit decoding phase without modifications, the total exponent te is introduced as $te = 2^{ES} \times k + e$, in agreement with Equation (4). Of course, te can be both positive, zero or negative, while e is a non-negative integer. We will use this intermediate representation as an helper, to compute the intermediate result of operations between two posits. Then we will normalize the intermediate result, hence we will transform FIR result into a proper posit. It is during this last phase that the posit rounding mechanism comes into play.

A. Addition and Subtraction

Addition and subtraction operations share several portion of arithmetic and logic parts. In general we want a FIR-represented posit p_{out} such that:

$$\begin{aligned} & \underbrace{[(-1)^{s_1} \times 2^{te_1} \times (1.f_1)]}_{p_1} \pm \underbrace{[(-1)^{s_2} \times 2^{te_2} \times (1.f_2)]}_{p_2} = \\ & = \underbrace{(-1)^{s_{out}} \times (2^{2^{ES}})^{k_{out}} \times 2^{e_{out}} \times (1.f_{out})}_{p_{out}} \end{aligned}$$

By supposing that $|p_1| \geq |p_2|$ we can transform the previous equation by extracting a scale factor $b = te_1 - te_2 \geq 0$ from the exponents, obtaining the following expression:

$$p_{out} = (-1)^{s_1} \times 2^{te_1} \times \underbrace{[(1.f_1) \pm (1.f_2) \times 2^{-b}]}_{1.f_{out}}$$

After this step, we just need to normalize $f_{sum} = (1.f_1) \pm (1.f_2) \times 2^{-b}$ to obtain the final $1.f_{out}$ value. With the *addition* we need to ensure that the sum is strictly lower than 2, to be in a valid format. If f_{sum} overflows 2, we scale it back by one position, summing 1 to te_1 . On the other hand, with the *subtraction*, we can normalize the result by subtracting the number of leading zeroes in f_{sum} from te_1 and shifting the result of the same number of positions to the left. Note that, at the end of addition, when we shift right f_{sum} of one position, we discard the last bit. If the discarded bit is 1, we signal this outside, to preserve the information for final rounding phase.

B. Multiplication

As before, we have two posits: $p_1 = \langle s_1, k_1, e_1, f_1 \rangle$ and $p_2 = \langle s_2, k_2, e_2, f_2 \rangle$. The goal is finding the FIR tuple

$\langle s_{out}, te_{out}, f_{out} \rangle$ such that:

$$\begin{aligned} & [(-1)^{s_1} \times 2^{te_1} \times (1.f_1)] \times [(-1)^{s_2} \times 2^{te_2} \times (1.f_2)] = \\ & = (-1)^{s_{out}} \times \underbrace{(2^{2^{ES}})^{k_{out}}}_{2^{te_{out}}} \times 2^{e_{out}} \times (1.f_{out}) \end{aligned}$$

where, of course, $te_1 = 2^{ES} \times k_1 + e_1$ and $te_2 = 2^{ES} \times k_2 + e_2$. Simplifying the equation, we obtain the following:

$$(-1)^{s_1 \oplus s_2} \times [2^{te_1} \times 2^{te_2}] \times [(1.f_1) \times (1.f_2)]$$

where \oplus is the exclusive OR between the two sign bits s_1 and s_2 . As we can see, $te_{out} = te_1 + te_2$ while $1.f_{out}$ is the integer multiplication of the two fractions, readjusted for the normalization (along with the possible increment of te_{out}).

C. Division

Again, we have two posits $\langle s_{1,2}, k_{1,2}, e_{1,2}, f_{1,2} \rangle$; the goal is finding the tuple $(s_{out}, te_{out}, f_{out})$ such that the following holds:

$$\frac{[(-1)^{s_1} \times 2^{te_1} \times (1.f_1)]}{[(-1)^{s_2} \times 2^{te_2} \times (1.f_2)]} = (-1)^{s_{out}} \times \underbrace{(2^{2^{ES}})^{k_{out}}}_{2^{te_{out}}} \times 2^{e_{out}} \times (1.f_{out}). \quad (6)$$

Simplifying the left-hand side of (6), we obtain

$$\frac{(-1)^{s_1}}{(-1)^{s_2}} \times \frac{2^{te_1}}{2^{te_2}} \times \frac{(1.f_1)}{(1.f_2)} \quad (7)$$

which similarly to the multiplication, suggests that $s_{out} = s_1 \oplus s_2$ and $te_{out} = te_1 - te_2$. From the mathematical standpoint it comes off as not too dissimilar than a product. If we consider that $(1.f)$ does not belong to \mathbb{R} , but instead it belongs to \mathbb{Q} , we can multiply both numerator and denominator such that they are two integer numbers. Equation (8) shows an example of this simplification, with the result being the one of an integer division.

$$\frac{\overbrace{1.011 \dots 0001001}^{F \text{ bits}}}{\underbrace{1.110 \dots 0001111}_{F \text{ bits}}} \equiv \frac{1011 \dots 0001001 \times 2^{\cancel{F}}}{1110 \dots 0001111 \times 2^{\cancel{F}}}. \quad (8)$$

D. Result normalization

In the final stage we take the FIR output from the previous one and output a posit number (i.e. the result of the operation). The final posit $(s_{out}, k_{out}, e_{out}, f_{out})$ is computed from the resulting FIR (s, te, f) . Firstly, we split te into a posit regime k' and exponent e :

$$\begin{cases} k' \leftarrow \left\lfloor \frac{te}{2^{ES}} \right\rfloor \\ k_{out} \leftarrow \text{clip}(k') \\ e \leftarrow te - 2^{ES} \times k_{out}. \end{cases} \quad (9)$$

Note that $k' \neq k_{out}$ since it may result in a regime length higher than the maximum regime length for a posit $\langle N, x \rangle$, that is $N - 1$, including the stop bit. Typically this means that we need to *clip* the value for k' to the (maximum,

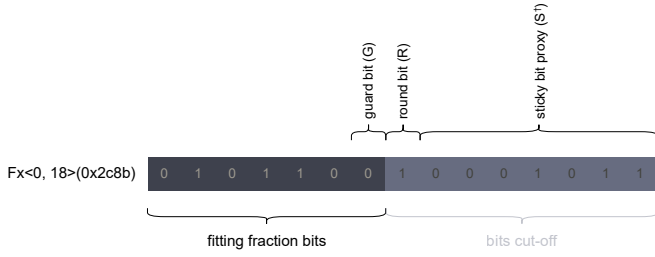


Fig. 3. Bit-string layout of the fraction before rounding, with (G, R, S) bits highlighted.

or minimum) k_{out} allowed. If the regime value is negative, its minimum value is $-(N - 1)$; if the regime is positive, its maximum value is $N - 2$, according to Eq. (2). After we compute the number of regime bits, the actual size of the exponent is inferred from the regime field size and the posit parameters; the fractional field size is computed as the remaining bits, if any. Finally, we need to accommodate for the fraction bits that fall off the posit size, if any. That turns out to be an implicit *rounding to lowest*. Depending on the adopted rounding scheme, a few more operations must be considered. Assuming the selected rounding scheme is *round to nearest even* (the only one considered in the posit standard [4]), we consider 3 bits from the fraction (Figure 3):

- guard bit (G): the least significant bit of the sequence of digits that fit the fraction field of the final posit,
- round bit (R): the most significant bit of the sequence of discarded digits,
- sticky bit (S): the *or-reduction*¹ of the sequence of bits to the right of the round bit (i.e. the discarded bits).

These bits will be used to complete the rounding of the posit, applying the *round to nearest even* policy as described in [4].

V. MULTI-STAGE FULL POSIT PROCESSING UNIT

This section presents the design of the FPPU. Similarly to previous section, we identify 3 main stages of execution inside the processing unit: i) decoding and input conditioning, ii) actual computation, iii) normalization and encoding of the result. To limit the length of pure combinatorial paths the FPPU has 4 pipelined stages corresponding to these 3 main stages of execution, with the computation phase split into two to take into account for the longer path in the division logic. Figure 4 shows a top-level schematic of the FPPU. The FPPU control unit has the task of signaling whether the current output is valid (i.e. 4 stages have been traversed by an instruction) or not. Figure 5 shows an example of interaction with the FPPU with the submission of an operation identified by OP and the two operands P1, P2. When the operands are ready, `valid_in` is set to active and after 3 cycles the FPPU produces a valid result PO signaled by `valid_out`.

¹adopting the Verilog nomenclature, the *or-reduction* operator (`|`) applies the bitwise inclusive *or* to the elements of a vector and returns a scalar.

A. The division algorithm

This Section highlights some properties of the division algorithm implemented in the FPPU. In Section IV we stated that the very last step in posit division is the integer division of the fractional parts. This operation is not as easy as multiplying or sum two integer numbers. Integer division is a problem that can be tackled in three main different ways [17]:

- Digit Recurrence algorithms: similar to pen and paper algorithm when dividing numbers. It produces a certain number of digits at each step by i) determining the next quotient digit(s), ii) multiplying such digits by the divider and iii) subtracting this result from the current remainder.
- Polynomial approximation: the division $\frac{x}{y}$ is computed multiplying x with an approximation of $\frac{1}{y}$ using a polynomial expression.
- Iterative approximation: the reciprocal approximation is computed iteratively (e.g. Newton-Raphson).

Our architecture combines the last two families of algorithms to provide an accurate approximation of the reciprocal to be used in the subsequent multiplication. In particular, we start from the idea of Chebyshev polynomials [18] to provide an approximation of the reciprocal across the interval (0.5, 1). The expression for a third-order Chebyshev polynomial that approximate the $\frac{1}{x}$ function is:

$$f(x) = 5.65685 - 11.75737x + 10.64818x^2 - 3.54939x^3 + \mathcal{O}(x^4). \quad (10)$$

The issue in using this approximation is that it requires a sequence of 6 fixed-point multiplications. An optimized approach is present in [19], where the reciprocal approximation is computed as in Algorithm 1, where k_1, k_2 are two parameter that we are going to evaluate to control the accuracy. This formulation has only 2 fixed point multiplication, since the $e \cdot 4$ can be implemented with a logical shift left of two positions. Expanding such routine in a polynomial we obtain

Algorithm 1 Reciprocal approximation of x from [19].

Require: x, k_1, k_2

Ensure: $y = \frac{1}{x}$

$b \leftarrow k_1 - x$

$c \leftarrow x \cdot b$

$d \leftarrow k_2 - c$

$e \leftarrow d \cdot b$

$y \leftarrow e \cdot 4$

an expression similar to (10):

$$f(x, k_1, k_2) = 4k_1k_2 - 4(k_1^2 + k_2)x + 8k_1x^2 - 4x^3. \quad (11)$$

Instead of using the k_1, k_2 values proposed in [19] we set up a minimization problem of the error function e^2 of the reciprocal approximation. In particular we defined this function as:

$$e^2(k_1, k_2) = \int_{1/2}^1 rerr^2(x, k_1, k_2) dx, \quad (12)$$

where *rerr* is the relative error between the approximated and the exact inverse function. We then use this function to solve

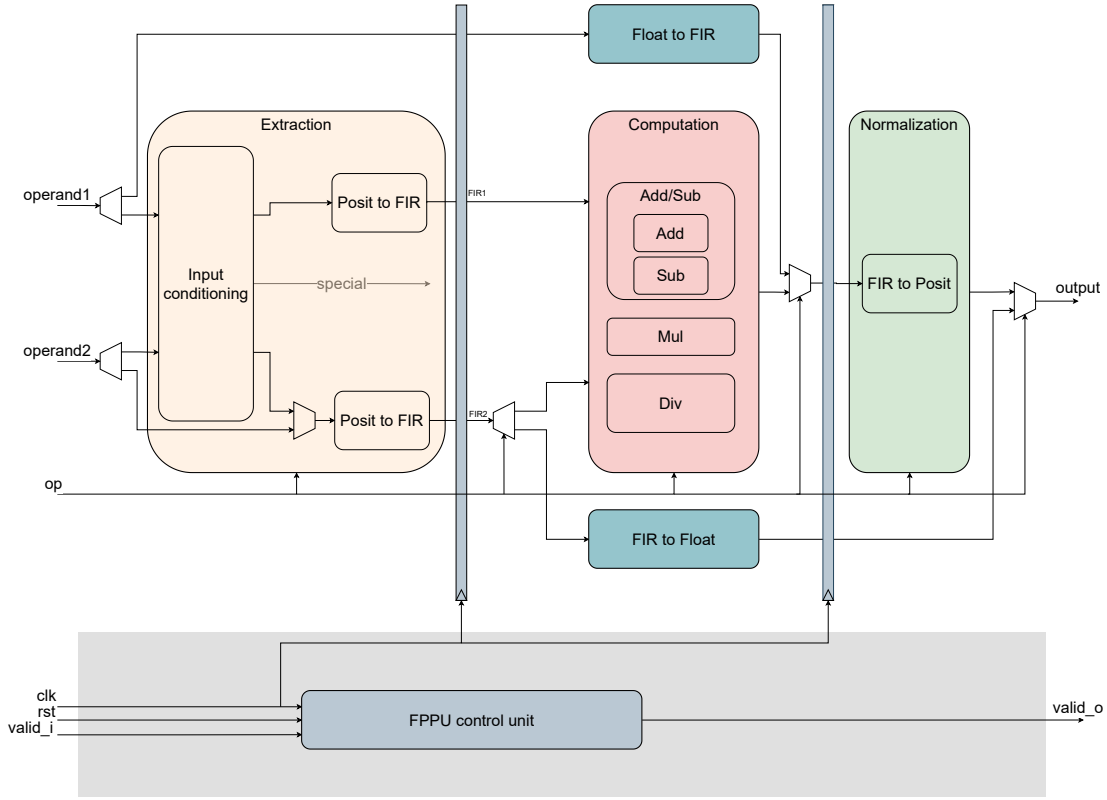


Fig. 4. 3-stage Full PPU with control unit.

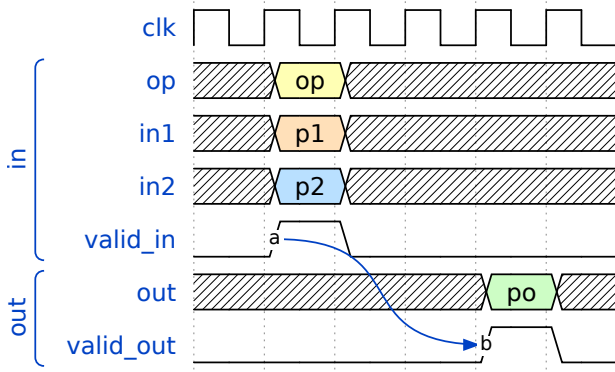


Fig. 5. Example of interaction with the FPPU.

TABLE II
PERCENTAGES OF POSITS $P(N, ES)$ INEXACT DIVISION RESULTS.
PACoGEN VERSION [11] VS PROPOSED.

N	ES	PACoGen				proposed	
		IN	OUT	NR	wrong [%]	NR	wrong [%]
8	0	8	9	0	4.8	1	1.4
8	1	8	9	0	5.4	1	1.2
8	2	8	9	0	9.3	1	2.1
8	3	8	9	0	13.5	1	4.2
8	4	8	9	0	16.4	1	7.5
16	0	8	9	1	10.0	1	1.5
16	1	8	9	1	10.0	1	0.6
16	2	8	9	1	8.8	1	0.5
16	3	8	9	1	9.0	1	0.1

the following problem:

$$(k_{1_{opt}}, k_{2_{opt}}) : \min\{e^2(k_1, k_2)\}. \quad (13)$$

The obtained solution gives $k_{1_{opt}} = 1.4567844114901045$ and $k_{2_{opt}} = 1.0009290026616422$, yielding a 36.4% improvement over [19]. We can pair this solution with a round of Newton Raphson using the output of this stage as starting condition to refine the accuracy of the approximation.

We compared the accuracy of our solution to the solution of PACoGen [11], that employs a pre-computed look-up table for reciprocal approximation and as input of a round of Newton-Raphson. Table II reports this comparison. In the table, the **IN** column is the number of fraction bits from the posits used to index the LUT, while **OUT** is the number of bits

of the reciprocal approximation of the fraction; **NR** indicates the number of Newton-Raphson rounds used and **wrong [%]** states the error percentage of the division results when compared to a software golden model for posit computation.

VI. RISC-V ISA EXTENSION AND COMPILER SUPPORT

This section briefly presents the RISC-V ISA extension for Posits and the related compiler support. We decided to reuse the existent RISC-V registers (i.e. $x1 \dots x31$). Since posits can be treated as signed integers inside the architecture, we followed the RV32I base integer instruction set from the RISC-V standard for arithmetic operations. In particular, ADD, SUB, MUL, DIV posit operations are encoded as *R-type* instructions (see [1]). We leveraged the custom opcode

TABLE III
INSTRUCTION LISTING FOR POSIT ARITHMETIC ISA EXTENSION

funct7	rs2	rs1	funct3	rd	opcode	R-TYPE
1100000	rs2	rs1	000	rd	0001011	PADD
1101010	rs2	rs1	001	rd	0001011	PSUB
1100000	rs2	rs1	010	rd	0001011	PMUL
1100000	rs2	rs1	100	rd	0001011	PDIV
rs3 — 00	rs2	rs1	000	rd	0101011	PFMADD

space 0x0B from the RISC-V standard to add the new instructions to the ISA. In Table VI we reported the listing of posit instructions used for the RISC-V ISA extension proposed in this work. Besides the posit arithmetic operations we also added conversion instructions between posits and binary32 numbers, so that we can enable the use of binary32 numbers as frontend while maintaining posit computation as backend in HW. To enable the use of the novel instructions in SW, we provided a set of intrinsic functions to map high-level C/C++ function call to the underlying machine code. Doing this, we do not need to change the RISC-V compiler but we generate the correct assembly for the posit ISA at compile time.

We report an example of intrinsic and relative call from C code in Listing 1. The `register` keyword suggests the compiler to put the values of the input operands and the result in three registers, since the instruction type is R-type. At lines 6, 7, 8 we set up the `opcode`, `funct3`, `funct7` parameters for the specific operation (this is the only part that varies between different operations).

```

1  posit_t padd(long a, long b) {
2      register int p1 asm("a1") = a;
3      register int p2 asm("a2") = b;
4      register int result asm("a0");
5      __asm__(
6          ".set op,0xb\n"
7          ".set of1,0\n"
8          ".set of2,0x6A\n"
9          ".byte op|(((r[result]>>1)&0xF)|"
10         "((r[result]>>1)&0xF)|"
11         "(of1<<4)|((r[1]&1)&1)<<7),\n"
12         "(((r[2]&0xF)<<4)|((r[1]>>1)&0xF),\n"
13         "(((r[2]>>4)&0x1)|((of2<<1)&1))\n"
14         : [result] "=r"(result)
15         : "r"(p1), "r"(p2), "[result]"(result));
16      return result;
17  }
```

Listing 1. Example of intrinsic for a posit addition operation.

A more complex example is shown in Listing 2 and 3 with the implementation of a simple square matrix multiplication and a 3×3 convolution.

```

1  void gemm(posit_t *a,
2            posit_t *b,
3            posit_t *c,
4            int n) {
5      for (int i = 0; i < n; i++) {
6          for (int j = 0; j < n; j++) {
7              posit_t sum = 0;
8              for (int k = 0; k < n; k++) {
9                  sum = padd(sum, pmul(a[i*n+k], b[k*n+j]));
10             }
11             c[i * n + j] = sum;
12         }
13     }
```

```

14 }
```

Listing 2. Square matrix-matrix multiplication using posit intrinsics.

```

1  void conv3x3(posit_t *a,
2               posit_t *f,
3               posit_t *c,
4               int n) {
5      for (int i = 0; i < n; i++) {
6          for (int j = 0; j < n; j++) {
7              posit_t sum = 0;
8              for (int k = 0; k < 3; k++) {
9                  for (int l = 0; l < 3; l++) {
10                     sum = padd(sum, pmul(a[(i+k)*n+j+l], f[k*3+l]));
11                 }
12             }
13             c[i * n + j] = sum;
14         }
15     }
16 }
```

Listing 3. 3×3 convolution using posit intrinsics.

We cover more examples and tests in Section VII with the verification and validation of the FPPU inside the Ibex Core.

VII. THE IBEX CORE AND FPPU INTEGRATION

The Ibex Core [10] is a small 32-bit RISC-V core with a 2-stage pipeline. It supports the Integer (I) or Embedded (E), Integer Multiplication and Division (M), Compressed (C), and B (Bit Manipulation) extensions. Since it does not employ a floating point unit, it is particularly useful to test the impact of adding real number arithmetic using the FPPU.

We put the FPPU alongside the Arithmetic Logic Unit (ALU) inside the execution stage of the Ibex pipeline. Since we did not add specific registers for the posit operations, there was no need to modify the *register file* but only the *decoder* module by adding the instructions designed in Section VI.

After the integration, we tested and validated the integration exploiting the RTL simulation with binaries compiled for the RV32IM architecture as shown in Fig. 6. Using the instruction tracer inside the Ibex core we managed to dump all the instructions executed by the core during the RTL simulation, including the newly added posit instructions. We fed the output of the tracer to a trace parser (written in python language) to evaluate the output of each posit instruction, comparing it to the same software golden model used for validating the standalone FPPU. Furthermore, we used the trace parser to assess the accuracy of each operation, comparing the result to the IEEE binary32 correspondent operation result.

A. Integration tests and validation

To test the compliance of the FPPU HDL design vs. a posit golden model, we tested the overall system with different DNN kernels on 32×32 matrices (i.e. size of images for MNIST/CIFAR10 datasets). In detail, we tested a matrix-matrix multiplication, a 3×3 convolution and a 4×4 average pooling. Each test was run on the 8-bit FPPU and on the 16-bit FPPU. We collected the instruction traces and fed them to the trace parser we described before. We then collected two sets of results: i) the accuracy w.r.t the posit golden model and ii) the normalized mean error w.r.t the same operation executed with binary32 format. Let r_p^i, r_f^i be, the i -th result

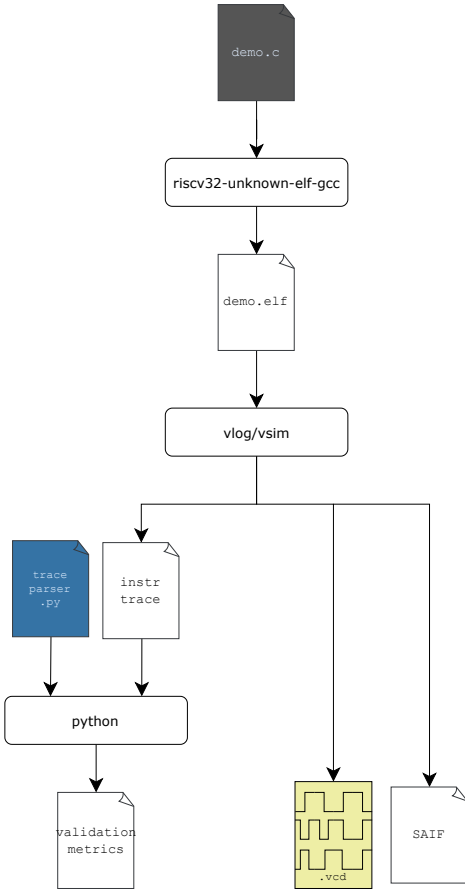


Fig. 6. Workflow diagram for compiling, simulating and validating the instruction set added to the RISC-V ISA.

TABLE IV

NORMALIZED MEAN ERROR OF FPPU OPERATIONS VS. CORRESPONDENT BINARY32 OPERATION IN SEVERAL LINEAR ALGEBRA TASKS.

	Conv (3x3 filter)		GEMM		Average Pooling 4x4	
	p(8,0)	p(16,2)	p(8,0)	p(16,2)	p(8,0)	p(16,2)
p.mul	0.042	0.004	0.019	0.003	-	-
p.add	0.025	0.0004	0.016	0.0007	0.019	0.0002
p.div	-	-	-	-	0.002	0

for, respectively, the posit and the binary32 operation, we computed the normalized mean error \overline{e}_{op} for a given operation op as follows:

$$\overline{e}_{op} = \frac{1}{N} \sum_i^N \left| \frac{r_p^i - r_f^i}{r_f^i} \right|.$$

Furthermore, we report more complex accuracy tests done with the posit format with the following benchmarks:

- Small datasets (MNIST, GTSRB and CIFAR-10) on LeNet5 convolutional neural network [20]–[22]
- Big datasets (ImagenetV2 [23], VOC2007 [24]) on complex DNN models (EfficientNetB0 [25] and Single Shot Detector SSD300).

Table IV shows the normalized mean error of posit operations compared to the correspondent binary32 ones. Figures 7 and 8 show accuracy performance of posit employed in several neural network tasks.

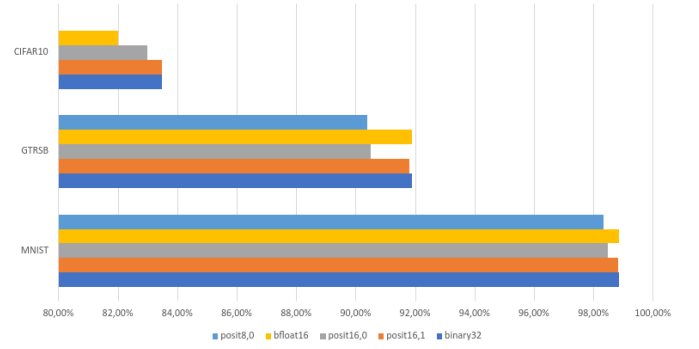


Fig. 7. Accuracy comparison between 8-bit, 16-bit posit formats and 32-bit IEEE binary32 on the LeNet-5 convolutional neural network.

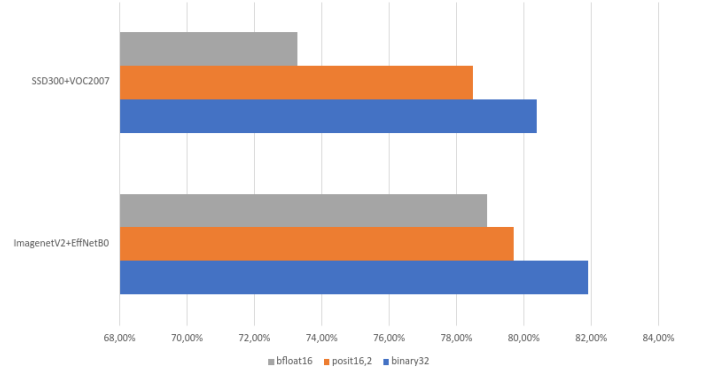


Fig. 8. Accuracy comparison between 16-bit formats (posit, bfloat16) and 32-bit IEEE binary32 on complex DNN task.

VIII. FPPU AND IBEX CHARACTERIZATION

In this section we characterize the FPPU by implementing it, alone or integrated within the Ibex RISC-V core, targeting Xilinx Alveo U280 (xcu280-fsvh2892-2L-e). From the synthesis process we extracted area and power metrics. Furthermore, we instantiated the Pulpino SoC [26] with the CV32E40P core [10], that embeds a 32-bit floating point unit, in order to compare the area occupation of the operation implementations in both the solutions. Fig. 9 shows the area occupation of the arithmetic logic unit (ALU) and the FPPU relatively to the total area of the Ibex core. As we can see, with an 8-bit PPU we are able to provide real number arithmetic capabilities to the core with an area cost that is less than the cost of the original Ibex ALU. Fig. 10 shows the comparison between FPPUs and FPU in terms of area occupation of the logic related to ADD, MUL and DIV operations. As we can see, using half of the bits for the representation results in a area cost that is less than half of the area for the 32-bit counterpart. The area cost for the comparison is not reported, since posits can be compared as signed integers, while binary32 numbers require dedicated circuits. We exploited the SAIF file produced at the end of the workflow shown in Fig. 6 to estimate the dynamic power of the FPPU component (see also [16]) in the four different operations. Table V reports the details for the four arithmetic operations. The results refer to Alveo FPGA implementation, and the dynamic power at a clock frequency of 20 MHz (with a maximum achievable frequency of 100

TABLE V
DETAIL OF MEAN AND STANDARD DEVIATION OF DYNAMIC POWER
MEASUREMENT OF THE FPPU COMPONENT.

	8-bit FPPU (mW)	16-bit FPPU (mW)
ADD	< 1	1
SUB	< 1	1
MUL	< 1	1
DIV	1	2

MHz) is below one mW for each 8-bit posit operation. The maximum FPPU peak throughput at 100 MHz is 33 MOps/s considering a 30 ns latency for the FPPU component with three pipeline stages.

A. SIMD configuration

Considering that we are under-using the 32-bit registers when employing 8-bit or 16-bit posits, we can think of increasing the number of FPPUs that elaborate in parallel, similarly to a Single Instruction Multiple Data (SIMD) paradigm. This can be done transparently to the instruction caller and with the same opcode. When using only one FPPU we just need to put the posit arguments in 8 (or 16) least significant bits. If we want to compute more posit operations in parallel, we would just need to put another three posits in the remaining 24 bits of the register (or another one posit in the remaining 16 bits). This can be easily implemented by reproducing the same FPPU 2 or 4 times (respectively, for 16-bit posits or for 8-bit posits). Then we can feed the same operand, valid, reset and clock inputs to all the FPPUs, while the two operands are constituted by portions of the source registers. The output is then concatenated from all the FPPU outputs into the destination register. We have verified that this approach would increase the throughput of the FPPU to 132 MOps/s for 8-bit posits and to 66 MOps/s for 16-bit posits. Of course the main difficulties arise from the software side, where we need to change the C++ source code (or any other programming language source code) to take advantage from this extended feature, and this can be challenging for complex applications. Indeed, we need either a compiler supporting automatic vectorization of the code or a partial rewrite of mathematical kernels to explicitly support the SIMD paradigm.

IX. DISCUSSIONS AND CONCLUSIONS

In this work we presented the design and integration of a full posit processing unit inside the low-power Ibex RISC-V core extending the RISC-V ISA to support posit arithmetic operations (add/sub, mul, div, posit-to-float/float-to-posit and FMA) as well as vectorized posit operations in a SIMD configuration of 4 units for 8-bit posits and 2 units for 16-bit posits. We provided several integration and validation test against both a posit golden model and a series of baseline binary32 neural network models (LeNet-5, EfficientNet, SSD300), showing at most a degradation of 4% in terms of inference accuracy on the object detection tasks (PascalVOC), while obtaining negligible drops in accuracy for 16-bit posits in the other simpler tasks (MNIST, CIFAR10, GTSRB). With respect to other works in literature like PaCoGen [11], this work allows for an improved

division algorithm in terms of accuracy. We characterized the unit on power and area comparing in particular the occupation of the unit in relation to the the full-core, showing that the posit-8 FPPU complexity is even lower than the Ibex ALU. The increase in area occupation of the FPPU is 7% for 8-bit posits and 15% for 16-bit posits. When compared to the FPU used by the CV32E40P RISC-V core [10], the FPPU area occupation is less than half for 16-bit posits (for the same accuracy of binary-32) and 1 order of magnitude lower for 8-bit posits (for an accuracy loss within 4%). When compared to work in [9] it has to be noted that the Light PPU gives to a RISC-V core the support to Posit with a limited complexity overhead and no speed overhead, but the support is limited to float-to/from-posit conversion. Hence, [9] is not effective for Posit computation since it requires that posit are converted in float, then processed by the FPU and then the results converted back in posit. Instead, the FPPU gives direct HW support to posit operations (sub, add, mul, div, inversion, conversion) so that computation in posit domain becomes efficient vs floating-point representation. Since the FPPU replaces in RISC-V the FPU the complexity overhead of the FPPU is lower than the overhead of Light PPU plus the FPU.

ACKNOWLEDGMENTS

Work partially supported by H2020 projects EPI-SGA2 (grant no. 101036168) and TEXTAROSSA (grant no. 956831, <https://textarossa.eu/>). In addition we wish to thank the Italian Ministry of Education and Research (MIUR) for funding this research within the framework of the FoReLab project (Departments of Excellence).

REFERENCES

- [1] "RISC-V ISA Specification," <https://riscv.org/specifications/isa-spec-pdf/>, (Accessed 2020-03-11).
- [2] "RISC-V History," <https://riscv.org/risc-v-history/>, Accessed May 28th, 2020.
- [3] 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.
- [4] Posit Working Group, "Standard for Posit(tm) arithmetic (2022)," https://posithub.org/docs/posit_standard-2.pdf, 2022.
- [5] M. Cococcioni, F. Rossi, E. Ruffaldi, and S. Saponara, "Small reals representations for deep learning at the edge: A comparison," in *Next Generation Arithmetic*, J. Gustafson and V. Dimitrov, Eds. Cham: Springer International Publishing, 2022, pp. 117–133.
- [6] —, "A novel posit-based fast approximation of elu activation function for deep neural networks," in *2020 IEEE International Conference on Smart Computing (SMARTCOMP)*, 2020, pp. 244–246.
- [7] 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)*, 2019, pp. 1421–1426.
- [8] —, "Deep positron: A deep neural network using the posit number system," in *2019 Design, Automation Test in Europe Conference Exhibition (DATE)*, March 2019, pp. 1421–1426.
- [9] M. Cococcioni, F. Rossi, E. Ruffaldi, and S. Saponara, "A lightweight posit processing unit for RISC-V processors in deep neural network applications," *IEEE Transactions on Emerging Topics in Computing*, pp. 1–1, 2021.
- [10] P. Schiavone, F. Conti, D. Rossi, M. Gautschi, A. Pullini, E. Flamand, and L. Benini, "Slow and steady wins the race? A comparison of ultra-low-power RISC-V cores for internet-of-things applications," in *27th International Symposium on Power and Timing Modeling, Optimization and Simulation (PATMOS)*, 2017, pp. 1–8.
- [11] 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.

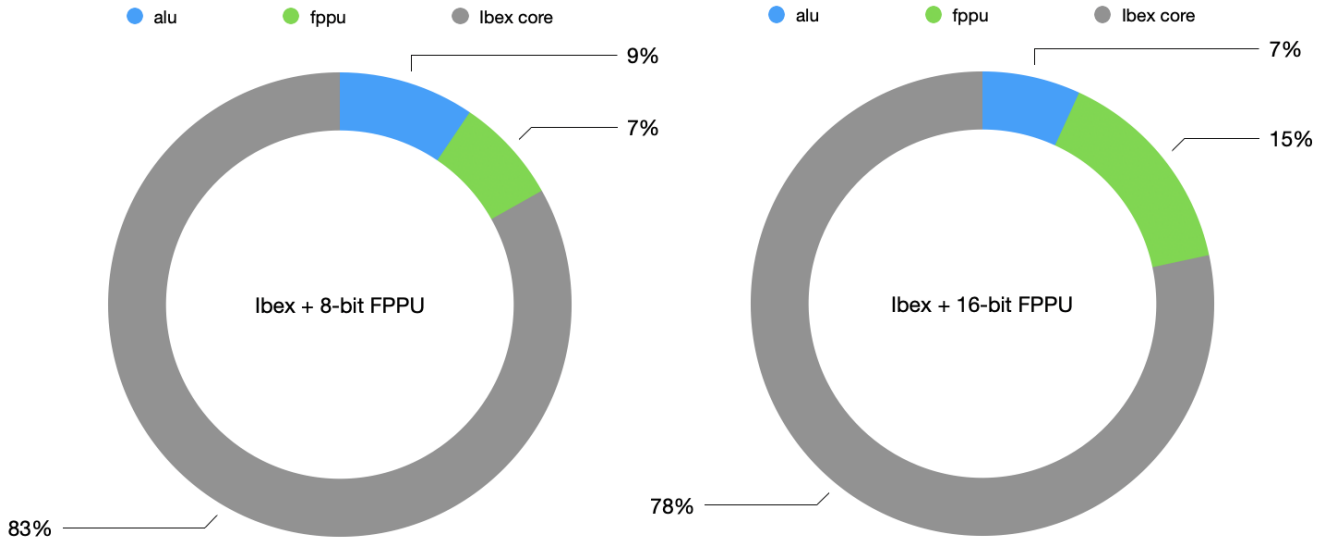


Fig. 9. Percent Area utilization (LUTs) of the FPPU with Posit (8, 2) (left) and Posit (16, 2) (right) and the other components of Ibex.

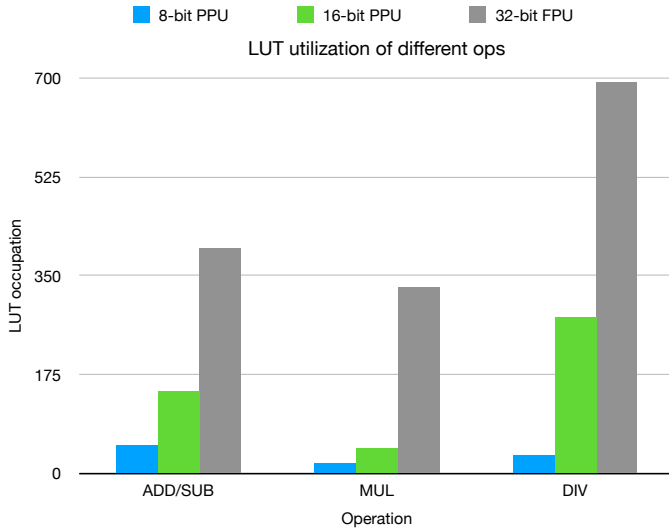


Fig. 10. Comparison between absolute area utilization (LUTs) of 8-bit, 16-bit FPPU and 32-bit FPU operations.

- [12] N. N. Sharma, R. Jain, M. M. Pokkuluri, S. B. Patkar, R. Leupers, R. S. Nikhil, and F. Merchant, "CLARINET: A quire-enabled RISC-V-based framework for posit arithmetic empiricism," *Journal of Systems Architecture*, vol. 135, p. 102801, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1383762122002867>
- [13] L. Crespo, P. Tomás, N. Roma, and N. Neves, "Unified posit/ieee-754 vector mac unit for transprecision computing," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 69, no. 5, pp. 2478–2482, 2022.
- [14] D. Mallasén, R. Murillo, A. A. D. Barrio, G. Botella, L. Piñuel, and M. Prieto-Matías, "Percival: Open-source posit risc-v core with quire capability," *IEEE Transactions on Emerging Topics in Computing*, vol. 10, no. 3, pp. 1241–1252, 2022.
- [15] 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.
- [16] M. Piccoli, D. Zoni, W. Fornaciari, G. Massari, M. Cococcioni, F. Rossi, S. Saponara, and E. Ruffaldi, "Dynamic Power Consumption of the Full Posit Processing Unit: Analysis and Experiments," in *14th Workshop on Parallel Programming and Run-Time Management Techniques for Many-Core Architectures and 12th Workshop on Design Tools and Architectures for Multicore Embedded Computing Platforms (PARMA-DITAM 2023)*, ser. Open Access Series in Informatics (OASISs), J. a. Bispo, H.-P. Charles, S. Cherubin, and G. Massari, Eds., vol. 107. Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023, pp. 6:1–6:11. [Online]. Available: <https://drops.dagstuhl.de/opus/volltexte/2023/17726>
- [17] J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres, "Hardware Implementation of Floating-Point Arithmetic," in *Handbook of Floating-Point Arithmetic*, J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres, Eds. Boston: Birkhäuser, 2010, pp. 269–320. [Online]. Available: https://doi.org/10.1007/978-0-8176-4705-6_9
- [18] N. Hale, *Chebyshev Polynomials*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2015, pp. 203–205. [Online]. Available: https://doi.org/10.1007/978-3-540-70529-1_126
- [19] "Reciprocal Approximation / Aliaksei Chapyzenka / Observable." [Online]. Available: <https://observablehq.com/@drom/reciprocal-approximation>
- [20] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov 1998.
- [21] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [22] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, "The German Traffic Sign Recognition Benchmark: A multi-class classification competition," in *In Proc. of the IEEE International Joint Conference on Neural Networks (IJCNN'11)*, 2011, pp. 1453–1460.
- [23] B. Recht, R. Roelofs, L. Schmidt, and V. Shankar, "Do imagenet classifiers generalize to imagenet?" 2019.
- [24] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results," <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
- [25] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," 2020.
- [26] M. Gautschi, P. Schiavone, A. Traber, I. Loi, A. Pullini, D. Rossi, E. Flamand, F. K. Gürkaynak, and L. Benini, "Near-threshold RISC-V core with DSP extensions for scalable IoT endpoint devices," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 25, no. 10, pp. 2700–2713, 2017.



Federico Rossi is a PhD student of the Information Engineering Department at University of Pisa. In 2019 he received his Master Degree in Computer Engineering *magna cum laude*. He is currently involved in the European Processor Initiative (EPI2), the TextaRossa and EuPilot european projects. His research topics include alternative real number representations and their applications to Deep Neural Networks for the automotive environment.



Sergio Saponara (SM'13) is Full Professor of Electronics at University of Pisa, where he got Master degree cum laude and Ph.D. degree. In 2012 he was a Marie Curie Research Fellow in IMEC. He is an IEEE Distinguished Lecturer and co-founder of special interest group on IoT of both IEEE CAS and SP societies. He is the director of I-CAS lab, of Crosslab Industrial IoT, of the Summer School Enabling Technologies for IoT. He is associate editor of several IEEE and Springer Journals. He co-authored more than 300 scientific publications and 18 patents.

He is the leader of many funded projects by EU and by companies like Intel, Magneti Marelli, Ericsson and PPC.



Francesco Urbani is a research fellow at the Department of Information Engineering at the University of Pisa. He received his Master's Degree in Electronic Engineering in 2022. He is currently involved in the European Processor Initiative (EPI) SGA-2 project. His research topics include alternative number representations for AI applications on the edge.



Marco Cococcioni (SM'12) received the Laurea degree in 2000 and the Diploma degree in 2001 in Computer Engineering from University of Pisa and Scuola Superiore S. Anna, respectively, both with *magna cum laude*. In 2004 he earned the Ph.D. degree in Computer Engineering at the University of Pisa. After working as a post-doc in the same department, in 2010-2011 he spent two years as Senior Visiting Scientist at the NATO Undersea Research Centre (now CMRE) in La Spezia, Italy.

For his collaboration with CMRE he obtained the NATO Scientific Achievement Award in 2014. Since 2016 he is an Associate Professor at the Department of Information Engineering of the University of Pisa. He is in the editorial board of several journals indexed by Scopus. He is member of three IEEE task forces: Genetic Fuzzy Systems, Computational Intelligence in Security and Defense, and Intelligent System Application. Prof. Cococcioni has co-authored more than 100 contributions to international journals and conferences and he is a Senior Member of both IEEE and ACM (Association for Computing Machinery). He has been involved in five H2020 European Projects (EPI SGA-1, EPI SGA-2, TextaRossa, EuPilot, and EUPEX).



Emanuele Ruffaldi (SM'18) is senior software engineer at Medical Microinstruments Inc. working on robotic assisted microsurgery. Formerly he has been Assistant Professor at Scuola Superiore Sant'Anna in the Perceptual Robotics laboratory, Pisa, Italy. His research interests are in the field of machine learning for HRI and embedded artificial intelligence. He is Senior IEEE Member and has served IEEE as Publicity Chair for the Haptics TC.