# Hardware and Software Support for Mixed Precision Computing: a Roadmap for Embedded and HPC Systems

William Fornaciari
william.fornaciari@polimi.it
DEIB Politecnico di Milano & CINI
National Laboratory HPC-KTT
Milan, Italy

Giovanni Agosta
giovanni.agosta@polimi.it
DEIB Politecnico di Milano & CINI
National Laboratory HPC-KTT
Milan, Italy

Daniele Cattaneo
daniele.cattaneo@polimi.it
DEIB Politecnico di Milano
Milan, Italy

Lev Denisov
lev.denisov@polimi.it
DEIB Politecnico di Milano
Milan, Italy

Andrea Galimberti
andrea.galimberti@polimi.it
DEIB Politecnico di Milano
Milan, Italy

Gabriele Magnani
gabriele.magnani@polimi.it
DEIB Politecnico di Milano
Milan, Italy

Davide Zoni
davide.zoni@polimi.it
DEIB Politecnico di Milano
Milan, Italy

## ABSTRACT

Mixed precision is an approximate computing technique that can be used to trade-off computation accuracy for performance and/or energy. It can be applied to many error-tolerant applications, but manual precision tuning is both tedious and error-prone. Furthermore, the effectiveness of the technique heavily depends on hardware characteristics. Therefore, a hardware/software co-design approach is necessary for an effective exploitation of precision tuning opportunities offered by the applications. In this paper, we propose, based on the state of the art of precision tuning software and mixed precision hardware, a roadmap for the evolution of hardware designs and compiler-based precision tuning support, which is ongoing in the context of the European projects TEXTAROSSA and APROPOS.

## CCS CONCEPTS

• **Software and its engineering → Compilers**; • **Computer systems organization → Reduced instruction set computing**; **Parallel architectures**.

## KEYWORDS

Approximate computing, mixed precision, hardware/software co-design, High Performance Computing (HPC)

## 1 INTRODUCTION

Approximate Computing is an emerging class of optimization techniques to trade off computation accuracy for performance and energy [9]. In general, Approximate Computing techniques can be performed in a wide range of contexts, from hardware to application level, and can introduce approximation in several ways, ranging from faults induced by undervolting the system, as in near threshold computing [13] to skipping entire iterations of a loop, as in loop perforation techniques [25], in real-time optimization [21], and for on-chip communication [29].

In this paper, we focus on one particular Approximate Computing technique, namely Mixed Precision Computing [3]. Mixed Precision Computing aims at controlling the accuracy-performance/energy trade-off at a fine grain, by modifying the data types involved in each computation. Whenever a computation is introduced in an application source code, it is assigned a data type based on the programmer understanding of the semantics of the variables and constants involved as well as the available data types in the programming language. However, this choice usually produces a computation that significantly exceeds the precision needed for the actual ranges of values involved in it at run-time, since the programmer cannot be asked to fine-tune it, even when the programming language and the underlying instruction set architecture allow it. In the context of resource-constrained embedded systems, where computing resources are scarce and floating point units may not be available at all in low power micro-controllers, an expert needs to manually tune the code designed by the application programmer to produce an optimized version.

To address this issue, new hardware extensions are being designed, including support for new data types such as bfloat16, and

compiler-based tools have been developed to support the programmer in selecting the best solution for their applications. However, these developments have mostly progressed in parallel, and there is now need to combine them in a hardware/software co-design approach in other to reap the highest possible benefits from Mixed Precision Computing. Two European efforts have recently started with complementary goals addressing the problem stated above. EuroHPC TEXTAROSSA [1, 2] focuses on heterogeneous platforms for High Performance Computing (HPC), including both reconfigurable fabrics and general purpose accelerators (GPU), whereas MSCA-ITN APROPOS [19] is more oriented towards low-power embedded systems. In this paper, we show how the two projects will work synergically towards the development of both hardware and software components needed for effective Mixed Precision Computing across the Computing Continuum.

In Section 2 we review the state of the art of Mixed Precision Computing hardware and compiler support. In Section 3 and 4 we outline the research roadmaps of TEXTAROSSA and APROPOS, while in Section 5 we outline some conclusions.

## 2 STATE OF THE ART

In this section, we provide an overview of the most relevant hardware and compiler methods for mixed-precision found in the literature.

### 2.1 Transprecision and Mixed Precision Hardware

Architectures targeting mixed-precision floating-point arithmetic must provide fast, energy-efficient, and area-efficient support for carrying out computations on a variety of floating-point formats.

The *FloPoCo* framework can generate hardware accelerators for basic as well as more complex, non-standard floating-point arithmetic operations on FPGA targets [12]. The designer can configure at design time the precision of the generated cores, which can then be integrated into a general-purpose CPU in the form of a floating-point unit (FPU) or employed as a co-processor. While *FloPoCo*-generated cores can compute complex operations far more efficiently than as a sequence of standard FPU operations in a baseline CPU, the high throughput of those accelerators is countered by a high usage of FPGA resources, which makes such cores suitable to high-end embedded systems and more complex computing platforms.

The fused multiply–accumulate (FMAC) unit for transprecision computing presented in [15], meant for ASIC designs, can compute multiple low-precision floating-point operations at the same time in a SIMD fashion. Meant to be integrated in HPC solutions, the FMAC unit adds, for each low-precision result, a bit acting as a flag for its accuracy, signaling whether the corresponding operation has to be computed again at a higher precision in order to achieve the desired accuracy. Significant compiler changes must be implemented to support such feature.

[27] introduced a transprecision FPU that was integrated within the PULPino RISC-V-based open-source microcontroller, meant for ultra-low-power applications. The proposed FPU can handle 32-, 16-, and 8-bit floating-point formats on the same datapath in a packed SIMD fashion, thus providing hardware support for transprecision

when coupled with a software framework that can explore and tune the precision and dynamic range of floating-point variables. The concurrent support for three floating-point formats on the same FPU is countered by the complexity and high resource utilization of the packed SIMD datapath.

The mixed-precision floating-point unit introduced in [32] targets FPGA chips and is meant to be suitable for embedded systems platforms, such as the SoC supporting the RISC-V instruction set in which it was integrated for evaluation purposes [24]. Each type of operation can be implemented with a different floating-point format selected at design time according to the accuracy and performance requirements and resource constraints. Once instantiated, the supported formats can not be modified at run time, unless performing reconfiguration on FPGA targets. No changes are needed to the compiler, which can still work with standard *float32* variables. The supported floating-formats are 32-bit ones with any number of mantissa bits ranging from 1 to 23 and with the same 8-bit exponent length as the IEEE 754 *float32* floating-point format.

While fixed-point operations can in general be computed as sequences of integer arithmetic and shift operations, mixed-precision computing making use of fixed-point arithmetic must also be supported at the hardware level to provide more effective performance. Few state-of-the-art solutions implement therefore dedicated instruction set extensions providing fixed-point operations coupled with the corresponding hardware architecture to execute them.

[16] added support for the RISC-V P extension to the RISC-V 64-bit CVA6 processor [28]. The RISC-V P extension [22] extends the RISC-V instruction set architecture (ISA) with support for packed SIMD instructions, including fixed-point instructions for the Q1.63, Q1.31, Q1.15, and Q1.7 fixed-point formats, i.e., formats with 1 integer bit and 63, 31, 15, and 7 fractional bits, respectively.

The mixed-precision hardware support for fixed-point arithmetic introduced in [30, 31] implements instead a custom extension for the RISC-V ISA to enable the execution of fixed-point multiplications and divisions with 32-bit fixed-point formats selected at run time. The custom instructions, whose support must be added at the compiler level, encode indeed, in addition to the two operands, the position of their decimal point, i.e., the number of their integer and fractional bits. The modifications to be applied to the baseline ALU implementing integer multiplication and division instructions are minimal, also in terms of FPGA resource utilization, making the proposed fixed-point hardware support suitable even for constrained platforms such as embedded systems at the edge.

### 2.2 Precision Tuning support at compiler level

Many of the precision tuning tools are implemented as a step in the program compilation process. There are multiple benefits to doing it this way: (i) operating on the lower level of granularity exposes more opportunities for optimizations, (ii) access to information about the target hardware allows to tailor the program to the specifics of that hardware (e.g. supported floating-point types, operations performance, etc.), (iii) it is non-intrusive and transparent for the programmer, (iv) it benefits from highly developed compiler ecosystem. In practical terms, precision tuning tools used to optimize real applications need to satisfy the following requirements: (i) be able to optimize programs containing loops, conditionals and

**Table 1: Precision tuning tools summary**

| Tool | Validation | Input Language | Algorithm |
|------|-----------|----------------|-----------|
| TAFFO | Static | C, C++ | Interval Arithmetic, Affine Arithmetic, ILP |
| Rosa | Static | Scala | Interval Arithmetic, Affine Arithmetic, SMT |
| Daisy | Static | Scala, C | Interval Arithmetic, Affine Arithmetic, SMT, rewriting rules |
| Precimonious | Dynamic | C, C++ | Delta-Debugging |
| FloatSmith | Dynamic | C, C++ | Algorithmic Differentiation, Delta-Debugging, Hierarchical Composition |

memory operations, (ii) be able to work with programs written in commonly-used programming language, (iii) support wide variety of execution platforms, (iv) be able to work with a modern compilation ecosystem. Very few tools discussed currently in the literature satisfy these points.

TAFFO [4] is a precision tuning tool implemented as a plugin for LLVM compiler framework. It works with programs written in C/C++ and compiles them into transprecision binaries trading off accuracy of result for the execution time. TAFFO requires programmer to annotate input variables with the dynamic intervals of their values. It then uses Interval Arithmetic to derive the dynamic intervals of other variables, and Affine Arithmetic to estimate the errors. Provided the information about the supported types and the speed of the floating- and fixed-point operations TAFFO uses Integer Linear Programming (ILP) model to select the most optimal type allocation for the particular execution platforms [5], allowing it to target a wide variety of systems ranging from HPC to embedded. It supports fixed-point [8] and multiple floating-point [5] formats. TAFFO can perform static precision tuning as well as run-time optimization [7].

Rosa [11] is a source-to-source compiler precision tuning tool for programs written in Scala. It requires programmer to use a special Real type together and preconditions and precision requirements on functions. From that information Rosa derives the type allocation that satisfies the requirements using Interval and Affine Arithmetics and SMT (Satisfiability modulo theories) Solver. It does not support loops and conditionals. Daisy [10] is a source-to-source compiler precision tuning tool that extends Rosa for programs written in Scala and C. It uses genetic algorithm to explore rewriting rules that may improve accuracy of the program.

FloatSmith [17] is a source-to-source compiler for precision tuning for programs written in C/C++. It integrates previously existing tools to create a complete precision tuning pipeline. FloatSmith requires programmer to annotate the variables that need to be tuned with their error thresholds. It uses Algorithmic Differentiation to statically estimate the error introduced by a type change. It uses dynamic evaluation of floating-point types configurations with different search strategies: delta-debugging, hierarchical, and composition of the successful configurations.

Precimonious [23] is an LLVM-based precision tuning tool for programs written in C/C++. It explores the configuration space for the types used in the program and finds the best within the given error threshold given as annotations in the program. It uses delta-debugging algorithm for more efficient exploration and tests configurations by running the selected configurations with the inputs provided by the programmer. It uses LLVM version 3, which limits its usefulness for the modern applications.

**Figure 1: Example of the mixed-precision floating-point architecture, in a configuration where additions-subtractions are performed on *float32* operands and multiplications and divisions are computed on *bfloat16* operands [32]**
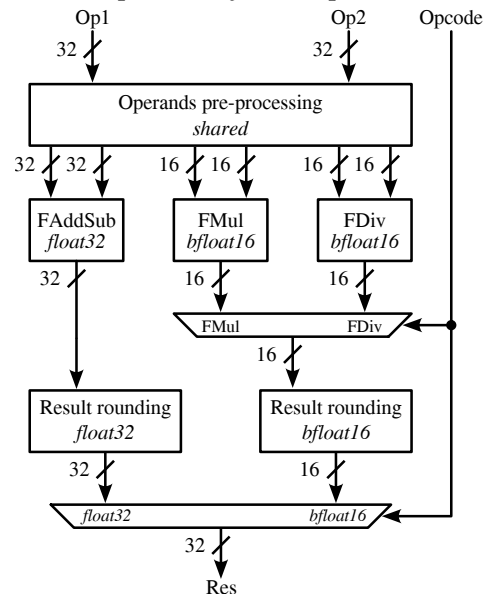


Table 1 summarises the relevant tool properties discussed in this section. For the more detailed discussion of the precision tuning tools we direct the reader to the survey: Cherubin and Agosta [6].

## 3 MIXED PRECISION HARDWARE IN TEXTAROSSA

In this section, we report the advances proposed in TEXTAROSSA towards mixed-precision hardware support.

### 3.1 Mixed-precision floating-point hardware support

The mixed-precision floating-point architecture [32] allows delivering floating-point hardware support where the precision of each type of operations, e.g., additions/subtractions, multiplications, and divisions, can be independently selected at design time depending on the target applications, on the accuracy requirements, and on the resource utilization constraints. Such flexibility in the supported floating-point formats still maintains a common dynamic range, in particular, the one of the standard IEEE 754 *float32* format, across

the entire FPU, providing two main advantages. On the one hand, the fixed dynamic range, i.e., the fixed number of bits encoding the floating-point exponent, simplifies the interoperability between the different floating-point data formats. On the other hand, the ensuing less complexity in the hardware architecture allows further optimizing the trade-off between efficiency and area.

Remarkably, when executing critical applications which mandate a higher accuracy than the one provided by a FPU implementing some combination of reduced-precision operations, resorting to the corresponding soft-float function calls can still guarantee the precision of *float32* computations, albeit at the cost of a reduction of performance. Such possibility can be exploited at the compiler level by selectively converting floating-point operations for which the HW support has reduced precision into soft-float function calls, which make use of the integer arithmetic resources of the CPU.

On the contrary, no changes or modifications must be applied on the compiler side to deal with the floating-point formats, possibly different from the standard *float32* one, employed by the different operations within the mixed-precision FPU. When low-precision formats are used, the operands are truncated to the desired precision and the ensuing result is extended by setting the least significant bits to 0s at the hardware level within the FPU, without any intervention required at the compiler or application level.

An instance of the mixed-precision FPU implementing *bfloat16* multiplications, conversions, and comparisons and *float32* additions/subtractions and divisions was shown to occupy 21% less resources and providing a 19% EDP improvement compared to a reference state-of-the-art FPU [20] while maintaining an average accuracy error below 3%.
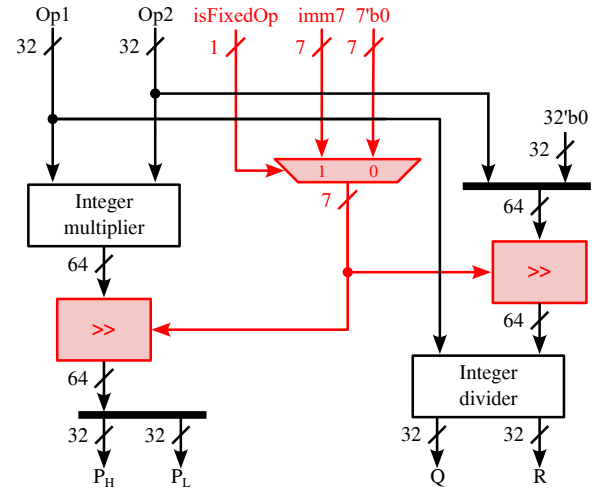
Figure 1 depicts an example configuration of our mixed-precision FPU, where additions and subtractions are implemented in the same functional unit according to the standard *float32* format, while multiplications and divisions are computed between *bfloat16* operands. The operands pre-processing logic, which takes care of extending the mantissa and exponent parts to deal with both normalized and denormalized operands and of identifying special values such as NaNs, infinites, and zeros, is shared between all the functional units computing the actual operations. On the contrary, rounding logic is instantiated for each of the hardware-supported formats. In the example, result rounding is performed separately for *float32* and *bfloat16* results, which are multiplexed from the corresponding functional units.

## 3.2 Mixed-precision fixed-point hardware support

The mixed-precision fixed-point architecture [30] provides hardware support for fixed-point multiplication and division instructions of the RISC-V ISA. Any 32-bit fixed-point format is supported, with the only constraint that both operands and the result share the same fixed-point format.

The selection of such specific fixed-point format is not made at design time for the fixed-point hardware support, but it is encoded within the instruction opcodes of the fixed-point multiplication and division instructions, which also include the number of integer and fractional digits. In particular, the RISC-V ISA was extended with eight fixed-point multiplication and division instructions, each

**Figure 2: Overview of the changes, highlighted in red, applied to integer multiplier and divider functional units to support fixed-point operations in the mixed-precision fixed-point architecture [30]**
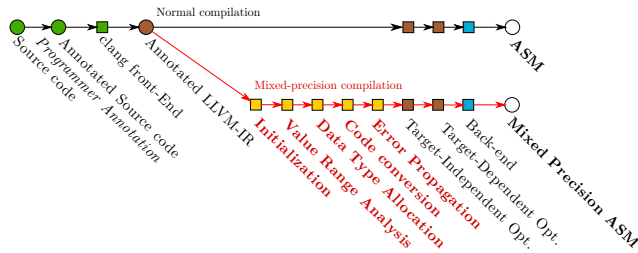


corresponding to a multiplication or division instruction from the standard RISC-V ISA M extension. The additional fixed-point instructions, coupled with the mixed-precision fixed-point hardware support, allow executing, at run time, multiplications and divisions instructions with any 32-bit fixed-point format while reusing the same hardware.

The fixed-point hardware support requires a limited number of additional FPGA resources compared to those required by the overall SoC implementing a CPU compliant to the RISC-V ISA integer (I) and multiplication/division (M) extensions. In particular, the experimental results highlighted an overhead in terms of resource utilization limited to 4%, compared to the baseline SoC packing a CPU supporting the lone RISC-V I and M extensions. On the energy-efficiency side, implementing fixed-point hardware support was shown to provide a 35% EDP improvement compared to the reference SoC also implementing an FPU, while maintaining a negligible accuracy loss.

Figure 2 depicts, highlighted in red, the changes applied to the baseline functional unit implementing integer multiplication and division operations in order to support the corresponding fixed-point operations. Logic meant to perform conversions between sign-magnitude and two's complement representations, to compute the sign of the result, and to manage signed and unsigned operations is instead not depicted, since it is not modified. In particular, the number of fractional bits of the operands and the result, encoded within the *imm7* 7-bit portion of the fixed-point instructions' opcode, is employed to shift the result of the integer multiplier and the divider operand of the integer divider, when executing fixed-point instructions, i.e., when the *isFixed* 1-bit flag is set to 1. On the contrary, when computing standard integer operations from the RISC-V ISA M extension, i.e., when the *isFixed* 1-bit flag is set to 0, the product, quotient, and remainder outputs are computed without applying any shift.

**Figure 3: TAFFO compilation flow performing floating to fixed point conversion [4]**



## 4 COMPILER-SUPPORTED PRECISION TUNING

In this section, we provide an overview of the goals, architecture and roadmap for the TAFFO compiler plugin set developed in TEXTAROSSA and APROPOS.

### 4.1 Precision Tuning Compiler: Architecture

Traditional compilers do not normally change the data types involved in the computation. The rationale is not only that compilers should not alter the semantics of the program (actually, most compilers can perform aggressive optimizations, although those are normally disabled by default), but that the compiler does not normally have information about the value ranges of variables, and thus cannot say much about the computed data as well. A precision tuning compiler, therefore, needs such information either from the programmer (by means of compiler hints expressed as annotations or pragmas) or from profiling (as in profile-guided optimization [18]). Our set of plugins for the LLVM compiler framework, TAFFO [4], takes the former approach, leveraging programmer annotations that express the value ranges of the input data.

Figure 3 shows the TAFFO compiler pipeline, compared with the standard compilation flow performed by the LLVM. Beyond some Initialization steps, the main activities that TAFFO performs are the Value Range Analysis, which propagates the information contained in the programmer annotations through the program data flow, computing the ranges for all intermediate values, and the Data Type Allocation, which selects the optimal allocation taking into account not only the beneficial effect on performance given by the reduced precision, but also the cost of converting data between different types. This process results in a clustering of operations, so that only a limited amount of type conversions are performed. The Code Generation and Error Propagation steps finally perform the conversion of the actual code, applying the decisions taken in the previous step, and check that the selected transformation does not catastrophically affect the computation error (in which case the transformation is undone).

TAFFO was first developed as the result of research ideas proposed in the FETHPC ANTAREX project [26], and currently developed under the umbrella of both the EuroHPC TEXTAROSSA and the MSCA-ITN APROPOS projects.

### 4.2 Precision Tuning for Heterogeneous Parallel HPC Architectures in TEXTAROSSA

When targeting parallel architectures, precision tuning tools need to tackle additional challenges that are not present in conventional single-threaded tasks. In particular, the tool needs to reliably detect each parallel region and the sets of variables shared between parallel execution threads. If this is not done, the transformed code might not be correct or appropriately transformed to mixed-precision. This task is even more troublesome for languages which do not support parallel programming paradigms without an auxiliary support library. This category includes languages of particular interest to HPC architectures such as C or C++, and languages used for compiler development like LLVM-IR.

To address these challenges, in the context of the TEXTAROSSA project we plan to integrate the TAFFO precision tuning plugins for the LLVM compiler framework with parallel-oriented languages supported by the same compiler. The choice of TAFFO is supported by its integrated architecture with the LLVM compiler framework, contrary e.g. to Daisy [10], which obviates the need of a specialized parser and code-generator as in a source-to-source compiler. Additionally, TAFFO is up to date with recent LLVM versions, contrary to Precimonious [23], which requires a severely outdated version of LLVM. The languages we plan to target encompass a large variety of parallel HPC architectures. In particular, we support OpenMP for CPU-based multiprocessing architectures, and we plan to support OpenCL and CUDA for the GPU-based SIMD paradigm and for GP-GPUs. In the future, we also envision additional extensions to support OmpSs, and the recently-proposed Posit numeric representation [14].

In order to add support for OpenMP-aware optimizations to TAFFO, we modify the Initializer and Conversion passes. The modifications allow the detection of specific OpenMP pragmas and of the outlined functions inserted by the *clang* frontend. This allows to mantain code correctness and to propagate the contextual information needed by precision tuning inside of parallel blocks. In Initializer, the program is searched for instances of call sites of the OpenMP fork function. At each call site, such function is temporarily deleted and replaced by a local trampoline, whose body simply calls the OpenMP outlined function This allows TAFFO's existing code to handle OpenMP programs without additional modifications. Additionally, we detect OpenMP's loop initialization library function to improve the loop trip count analysis already provided by LLVM.

A similar approach will be used to implement support for OpenCL and CUDA, with the additional complication that it is necessary to subject both host code and kernel code to optimization. In particular, the data types in the signature of the kernel functions must be kept coherent between host and device. Furthermore, it is necessary to detect where buffers are created in the host code in order to propagate annotations from the host code to the kernel code.

### 4.3 Precision Tuning for Embedded Systems in APROPOS

The APROPOS project leverages novel micro-controller architectures that expose mixed-precision or trans-precision arithmetic units, and develops compiler-based techniques to achieve the best

performance and energy efficiency within the application constraints on accuracy. To this end, we need to extend the TAFFO plugins set to support multiple data types, both floating and fixed point. Depending on the target hardware, it may be possible to achieve a co-design scenario, where the compiler can analyze the application based on the developer hints providing the quality of service requirements in terms of expected maximum relative error, determine the optimal data type selection for the various regions, and then, based on the architectural options available, apply the appropriate transformation to the code to generate the best mix of data types, as well as a configuration file for the generation or selection of the actual hardware platform.

To this end, APROPOS will need to extend the TAFFO framework to perform the error estimation and the data type selection not only for fixed point operations, but also for floating point ones, to enable the support of mixed-precision floating point units such as those developed in TEXTAROSSA and described in Section 3.1. In APROPOS, the TAFFO pipeline will be extended to operate first the partitioning of operation in two sets – those that can be performed in fixed point arithmetics, and those that need to be performed using floating point. Then, TAFFO will need to select the fixed-point width using the LuIS methodology [5], and finally to select the floating point width. New metrics combining sufficient precision and limited computational effort will be needed to perform this second step. Finally, the back-end of TAFFO will be extended to support the RISC-V instruction set architecture, as well as its relevant extensions.

## 5 CONCLUSIONS

In this paper, we have presented a roadmap towards an effective co-design methodology for mixed precision computing, supporting a range of different platform options for both HPC and Embedded Systems scenarios. During the next two years, in the context of the EuroHPC TEXTAROSSA and MSCA-ITN APROPOS projects, we will work towards the effective implementation of this vision in terms of both RISC-V-based hardware platforms and extensions to the TAFFO open source precision tuning tool set.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Giovanni Agosta et al. 2021. TEXTAROSSA: Towards EXtreme scale Technologies and Accelerators for euROhpc hw/Sw Supercomputing Applications for exascale. In *2021 24th Euromicro Conference on Digital System Design (DSD)*. 286–294. https://doi.org/10.1109/DSD53832.2021.00051

[2] Giovanni Agosta et al. 2022. Towards EXtreme scale technologies and accelerators for euROhpc hw/Sw supercomputing applications for exascale: The TEXTAROSSA approach. *Microprocessors and Microsystems* 95 (2022), 104679. https://doi.org/10.1016/j.micpro.2022.104679

[3] Marc Baboulin et al. 2009. Accelerating scientific computations with mixed precision algorithms. *Computer Physics Communications* 180, 12 (2009), 2526–2533.

[4] Daniele Cattaneo, Michele Chiari, Giovanni Agosta, and Stefano Cherubin. 2022. TAFFO: The compiler-based precision tuner. *SoftwareX* 20 (2022), 101238. https://doi.org/10.1016/j.softx.2022.101238

[5] Daniele Cattaneo, Michele Chiari, Nicola Fossati, Stefano Cherubin, and Giovanni Agosta. 2021. Architecture-aware Precision Tuning with Multiple Number Representation Systems. In *2021 58th ACM/IEEE Design Automation Conference (DAC)*. 673–678. https://doi.org/10.1109/DAC18074.2021.9586303

[6] Stefano Cherubin and Giovanni Agosta. 2020. Tools for Reduced Precision Computation: a Survey. *Comput. Surveys* 53, 2 (Apr 2020), 35 pages. https://doi.org/10.1145/3381039

[7] Stefano Cherubin, Daniele Cattaneo, Michele Chiari, and Giovanni Agosta. 2020. Dynamic Precision Autotuning with TAFFO. *ACM Trans. Archit. Code Optim.* 17, 2, Article 10 (May 2020), 26 pages. https://doi.org/10.1145/3388785

[8] Stefano Cherubin, Daniele Cattaneo, Michele Chiari, Antonio Di Bello, and Giovanni Agosta. 2019. TAFFO: Tuning Assistant for Floating to Fixed point Optimization. *IEEE Embedded Systems Letters* (2019). https://doi.org/10.1109/LES.2019.2913774

[9] Eva Darulova, Babak Falsafi, Andreas Gerstlauer, and Phillip Stanley-Marbell. 2021. Approximate Systems (Dagstuhl Seminar 21302). *Dagstuhl Reports* 11, 6 (2021), 147–163. https://doi.org/10.4230/DagRep.11.6.147

[10] Eva Darulova, Einar Horn, and Saksham Sharma. 2018. Sound Mixed-precision Optimization with Rewriting. In *Proceedings of the 9th ACM/IEEE International Conference on Cyber-Physical Systems* (Porto, Portugal) *(ICCPS '18)*. 208–219. https://doi.org/10.1109/ICCPS.2018.00028

[11] Eva Darulova and Viktor Kuncak. 2017. Towards a Compiler for Reals. *ACM Trans. Program. Lang. Syst.* 39, 2, Article 8 (March 2017), 28 pages.

[12] Florent de Dinechin and Bogdan Pasca. 2011. Designing Custom Arithmetic Data Paths with FloPoCo. *IEEE Design & Test of Computers* 28, 4 (July 2011), 18–27.

[13] Ronald G Dreslinski et al. 2010. Near-threshold computing: Reclaiming moore's law through energy efficient integrated circuits. *Proc. IEEE* 98, 2 (2010), 253–266.

[14] J L Gustafson and I T Yonemoto. 2017. Beating floating point at its own game: Posit arithmetic. *Supercomputing frontiers and innovations* 4, 2 (2017), 71–86.

[15] H. Kaul et al. 2012. A 1.45GHz 52-to-162GFLOPS/W variable-precision floating-point fused multiply-add unit with certainty tracking in 32nm CMOS. In *2012 IEEE International Solid-State Circuits Conference*. 182–184.

[16] Davy Koene. 2021. *Implementation and Evaluation of Packed-SIMD Instructions for a RISC-V Processor*. Master's thesis. TU Delft. https://repository.tudelft.nl/islandora/object/uuid%3Ac4162ff8-9419-4434-852d-c1c3297df808

[17] Michael O. Lam, Tristan Vanderbruggen, Harshitha Menon, and Markus Schordan. 2019. Tool Integration for Source-Level Mixed Precision. In *2019 IEEE/ACM 3rd International Workshop on Software Correctness for HPC Applications (Correctness)*. 27–35. https://doi.org/10.1109/Correctness49594.2019.00009

[18] Eduard Mehofer, R Gupta, and Y Zhang. 2002. *The Compiler Design Handbook: Optimizations and Machine Code Generation*. CRC Press, Chapter Profile guided code optimizations.

[19] Aleksandr Ometov and Jari Nurmi. 2022. Towards approximate computing for achieving energy vs. accuracy trade-offs in *2022 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 632–635.

[20] OpenRISC. 2022. mor1kx - an OpenRISC processor IP core. https://github.com/openrisc/mor1kx. [Online; accessed 27-May-2022].

[21] Federico Reghenzani, Giuseppe Massari, and William Fornaciari. 2020. Timing Predictability in High-Performance Computing With Probabilistic Real-Time. *IEEE Access* 8 (2020), 208566–208582. https://doi.org/10.1109/ACCESS.2020.3038559

[22] RISC-V Foundation. 2019. RISC-V "P" Extension Proposal, version 0.9.11-draft20211209. https://github.com/riscv/riscv-p-spec/raw/5a12c90b2c206c501a4489eb79e5d4d46afa1014/P-ext-proposal.pdf

[23] Cindy Rubio-González et al. 2013. Precimonious: Tuning Assistant for Floating-point Precision. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis* (Denver, Colorado) *(SC '13)*. Article 27, 12 pages. https://doi.org/10.1145/2503210.2503296

[24] Giovanni Scotti and Davide Zoni. 2019. A fresh view on the microarchitectural design of FPGA-based RISC CPUs in the IoT Era. *Journal of Low Power Electronics and Applications* 9 (02 2019), 19. https://doi.org/10.3390/jlpea9010009

[25] Stelios Sidiroglou-Douskos, Sasa Misailovic, Henry Hoffmann, and Martin Rinard. 2011. Managing performance vs. accuracy trade-offs with loop perforation. In *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering*. 124–134.

[26] Cristina Silvano et al. 2016. AutoTuning and Adaptivity appRoach for Energy efficient eXascale HPC systems: the ANTAREX Approach. In *Proc. of the 2016 Conf. on Design, Automation & Test in Europe* (Dresden, Germany). 708–713.

[27] G. Tagliavini, S. Mach, D. Rossi, A. Marongiu, and L. Benini. 2018. A transprecision floating-point platform for ultra-low power computing. In *2018 Design, Automation Test in Europe Conference Exhibition (DATE)*. 1051–1056. https://doi.org/10.23919/DATE.2018.8342167

[28] Florian Zaruba and Luca Benini. 2019. The Cost of Application-Class Processing: Energy and Performance Analysis of a Linux-Ready 1.7-GHz 64-Bit RISC-V Core in 22-nm FDSOI Technology. *IEEE Transactions on Very Large Scale Integration*

*(VLSI) Systems* 27, 11 (2019), 2629–2640.

[29] Davide Zoni, José Flich, and William Fornaciari. 2016. CUTBUF: Buffer Management and Router Design for Traffic Mixing in VNET-Based NoCs. *IEEE Transactions on Parallel and Distributed Systems* 27, 6 (2016), 1603–1616. https://doi.org/10.1109/TPDS.2015.2468716

[30] Davide Zoni and Andrea Galimberti. 2022. Cost-effective fixed-point hardware support for RISC-V embedded systems. *Journal of Systems Architecture* 126 (2022), 102476. https://doi.org/10.1016/j.sysarc.2022.102476

[31] Davide Zoni, Andrea Galimberti, and William Fornaciari. 2020. Flexible and Scalable FPGA-Oriented Design of Multipliers for Large Binary Polynomials. *IEEE Access* 8 (2020), 75809–75821. https://doi.org/10.1109/ACCESS.2020.2989423

[32] Davide Zoni, Andrea Galimberti, and William Fornaciari. 2021. An FPU design template to optimize the accuracy-efficiency-area trade-off. *Sustainable Computing: Informatics and Systems* 29 (2021), 100450. https://doi.org/10.1016/j.suscom.2020.100450