

**Towards EXtreme scale Technologies and Accelerators for euROhpc
hw/Sw Supercomputing Applications for exascale**



textarossa

**WP2 New accelerator designs exploiting mixed
precision**

D2.6 IP with data compression, part 1

Revised version

<http://textarossa.eu>



This project has received funding from the European Union's Horizon 2020 research and innovation programme, EuroHPC JU, grant agreement No 956831



textarossa

TEXTAROSSA Towards EXtreme scale Technologies and Accelerators for euROhpc hw/Sw Supercomputing Applications for exascale

Grant Agreement No.: 956831

Deliverable: D2.6 IP with data compression, part 1

Project Start Date: 01/04/2021

Duration: 36 months

Coordinator: AGENZIA NAZIONALE PER LE NUOVE TECNOLOGIE, L'ENERGIA E LO SVILUPPO
ECONOMICO SOSTENIBILE - ENEA, Italy.

Deliverable No	D2.6 (revised)
WP No:	WP2
WP Leader:	CINI-UNIFI
Due date:	M18 (September 30, 2022)
Delivery date:	13/10/2022, Revised 21/05/2023

Dissemination Level:

PU	Public	X
PP	Restricted to other programme participants (including the Commission Services)	
RE	Restricted to a group specified by the consortium (including the Commission Services)	
CO	Confidential, only for members of the consortium (including the Commission Services)	



This project has received funding from the European Union's Horizon 2020 research and innovation programme, EuroHPC JU, grant agreement No 956831



DOCUMENT SUMMARY INFORMATION

Project title:	Towards EXtreme scale Technologies and Accelerators for euROhpc hw/Sw Supercomputing Applications for exascale
Short project name:	TEXTAROSSA
Project No:	956831
Call Identifier:	H2020-JTI-EuroHPC-2019-1
Unit:	EuroHPC
Type of Action:	EuroHPC - Research and Innovation Action (RIA)
Start date of the project:	01/04/2021
Duration of the project:	36 months
Project website:	textarossa.eu

WP2 New accelerator designs exploiting mixed precision

Deliverable number:	D2.6					
Deliverable title:	IP with data compression, part 1					
Due date:	M18					
Actual submission date:	M27 (revised)					
Editor:	Sergio Saponara					
Authors:	S. Saponara, F. Rossi, F. Urbani					
Work package:	WP2					
Dissemination Level:	Public					
No. pages:	16 (revised 22)					
Authorized (date):	13/10/2022, revised 15/05/2023					
Responsible person:	Sergio Saponara					
Status:	Plan	Draft	Working	Final	Submitted	Approved

Revision history:

Version	Date	Author	Comment
0.1	2022-09-30	S. Saponara	Draft structure
0.2	2022-10-12	F. Rossi	First version completed
0.3	2022-10-12	D. Gregori	First Review
1.0	2022-10-13	F. Rossi	Final version completed
1.1	2023-05-10	F. Rossi	Revised version after mid-term review

Quality Control:

Checking process	Who	Date
Checked by internal reviewer	Daniele Gregori	October 12th, 2022
Revised, checked by internal reviewer	Daniele Gregori	May 15 th , 2023
Checked by WP Leader	Sergio Saponara	October 11th, 2022
Revised, checked by WP leader	Sergio Saponara	May 15 th , 2023
Checked by Project Coordinator	Massimo Celino	October 13th, 2022
Revised, Checked by Project Coordinator	Massimo Celino	May 15 th , 2023

COPYRIGHT

Copyright by the **TEXTAROSSA** consortium, 2021-2024

This document contains material, which is the copyright of TEXTAROSSA consortium members and the European Commission, and may not be reproduced or copied without permission, except as mandated by the European Commission Grant Agreement No. 956831 for reviewing and dissemination purposes.

ACKNOWLEDGEMENTS

This project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement no 956831. The JU receives support from the European Union's Horizon 2020 research and innovation programme and Italy, Germany, France, Spain, Poland.

Please see <http://textarossa.eu> for more information on the TEXTAROSSA project.

The partners in the project are AGENZIA NAZIONALE PER LE NUOVE TECNOLOGIE, L'ENERGIA E LO SVILUPPO ECONOMICO SOSTENIBILE (ENEA), FRAUNHOFER GESELLSCHAFT ZUR FOERDERUNG DER ANGEWANDTEN FORSCHUNG E.V. (FHG), CONSORZIO INTERUNIVERSITARIO NAZIONALE PER L'INFORMATICA (CINI), INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET AUTOMATIQUE (INRIA), BULL SAS (BULL), E4 COMPUTER ENGINEERING SPA (E4), BARCELONA SUPERCOMPUTING CENTER-CENTRO NACIONAL DE SUPERCOMPUTACION (BSC), INSTYTUT CHEMII BIOORGANICZNEJ POLSKIEJ AKADEMII NAUK (PSNC), ISTITUTO NAZIONALE DI FISICA NUCLEARE (INFN), CONSIGLIO NAZIONALE DELLE RICERCHE (CNR), IN QUATTRO SRL (in4). Linked third parties of CINI are POLITECNICO DI MILANO (CINI-POLIMI), Università di Torino (CINI-UNITO) and Università di Pisa (CINI-UNIPI); linked third party of INRIA is Université de Bordeaux; in-kind third party of ENEA is Consorzio CINECA (CINECA); in-kind third party of BSC is Universitat Politècnica de Catalunya (UPC).

The content of this document is the result of extensive discussions within the TEXTAROSSA © Consortium as a whole.

DISCLAIMER

The content of the publication herein is the sole responsibility of the publishers and it does not necessarily represent the views expressed by the European Commission or its services.

The information contained in this document is provided by the copyright holders "as is" and any express or implied warranties, including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose are disclaimed. In no event shall the members of the TEXTAROSSA collaboration, including the copyright holders, or the European Commission be liable for any direct, indirect, incidental, special, exemplary, or consequential damages (including, but not limited to, procurement of substitute goods or services; loss of use, data, or profits; or business interruption) however caused and on any theory of liability, whether in contract, strict liability, or tort (including negligence or otherwise) arising in any way out of the use of the information contained in this document, even if advised of the possibility of such damage.

Table of contents

List of acronyms	7
Executive summary	9
1. Introduction	10
2. Posit number and CppPosit library	11
3. Light PPU IP and interfacing with RISC-V	13
4. Implementation results and benchmarks	18
5. Conclusions and IP repository	21
6. References	22

List of Acronyms

Acronym	Definition
AI	Artificial Intelligence
ALU	Arithmetic Logic Unit
ASIC	Application Specific Integrated Circuit
CINI	Consorzio Interuniversitario Nazionale per l'Informatica
CPU	Central Processing Unit
DNN	Deep Neural Network
FP32	Floating Point 32 bit
FPGA	Field Programmable Gate Array
GTSRB	German Traffic Road Sign Benchmark
HDL	Hardware Description Language
HW	Hardware
HPC	High-Performance-Computing
INFN	Istituto Nazionale di Fisica Nucleare
IP	Intellectual Property
IPR	Intellectual Property Rights
ISA	Instruction Set Architecture
KPI	Key Performance Indicator
LUT	Look-up table
ML	Machine Learning
PMB	Project Management Board
PPU	Posit Processing Unit
RISC	Reduced Instruction Set Computer
SOC	System on Chip
SW	Software
UART	Universal Asynchronous Receiver Transmitter interface
VPU	Vector Processor Unit



Executive Summary

This document reports the activities done by Textarossa partner CINI (UNIPISA), with reference to preliminary HDL design, verification, and synthesis of accelerator macrocells in WP2 for IP with data compression.

The IP with data compression, that according to what foreseen in D2.1 is called Light PPU (Posit Processing Unit), has been implemented in FPGA technology.

Moreover, it has been integrated with a RISC-V open core like the Ariane RISC-V 64 bits to demonstrate the feasibility of integrating the Light PPU within a SoC based on RISC-V.

The IP with data compression is designed according to the specifications defined in D2.1 [1].

The IP can be accessed at:

https://bitbucket.org/federicorossifr/ppu_public/src/master/

The proposed work has been also submitted to the review of the scientific community and has been accepted on IEEE Transactions on Emerging Computing.

1. Introduction

This document D2.6 reports the activities done by Textarossa partner CINI (UNIPISA) in WP2.

D2.6 deals with the preliminary HDL design, using SystemVerilog, verification and synthesis of an IP for data compression that exploits the data compression capability of Posits.

D2.6 follows the revised D2.1 specifications [1].

The main innovation vs. what is already available in the state of art and in other EuroHPC projects like EPI, is in the hardware support of a new arithmetic format called Posit that, particularly for ML and DNN applications, has been proved in recent literature [2-4] to have a high compression effect:

Posit data format can ensure up to 4x compression for the same quality of floating point 32 (FP32), where quality is measured, as example, for a detection or classification task as accuracy (i.e. ratio between correct detection/classification and the total number of detection/classification done).

This deliverable is organized as follows:

Section 2 discusses the Posit numbers and a SW library implementing them called CppPosit, focusing on their data compression properties vs Floating-point numbers.

Section 3 discusses the RISC-V Posit Instruction Set Architecture design and implementation to support thanks to the Light PPU; at minimal overhead, the translation from/to Posits and Floats, and from/to Posits and integers can ensure that the computation of a DNN can be done using conventional ALU and FPU, but the storage can be done more effectively in Posit format.

Section 4 discusses in Subsection 4.1 Hardware results when implementing the IP, called Light PPU, in FPGA technology according to the specifications defined in D2.1, and then verified and integrated with a Ariane RISC-V 64 bits processor core. In Subsection 4.2 functional benchmarks for DNN examples are shown.

Conclusions are drawn in Section 5, where also the IP repository is listed.

2. Posit numbers and CppPositLibrary

A posit number (see Fig. 2.1) 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 (N, ES).

The Posit 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. 2.1. Note that, if the regime field 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 zeros to the right of the exponent bits in the format.

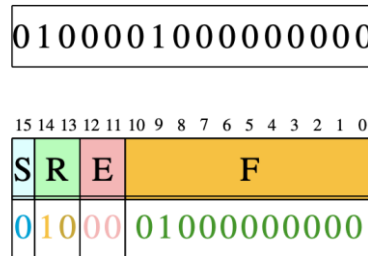


Figure 2.1: example of a Posit number

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

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

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

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

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

$$\text{used} := 2^{2^{ES}}$$

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:

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

The value F is the length of the fraction field.

Note that there will always be an implicit one in front of the fraction (i.e. 1.f1, f2, . . . , fF), without any subnormal number differently from IEEE binary32 numbers.

Software support for posits is enabled by our cppPosit library `cppPosit`, developed in Pisa and maintained by the Pisa authors, see here the public and open source part <https://github.com/eruffaldi/cppPosit>

The library uses templatization to define different posit configurations during compilation. The posit operations are put into four different levels L1-L4 with increasing computational complexity.

The main features of the cppPosit library are:

- Template-based posit definition and configuration
- Compile-time selection of different backends (software or hardware based)
- Use of operator overloading to hide backends implementation detail to end-user
- Ability to seamlessly target hardware accelerators or co-processors with operator overloading directly on the application code.

Thanks to this library we can transparently use the lightweight PPU in a neural network framework without changing the application code; indeed, we just to re-compile the application with a proper flag to enable the use of newly inserted RISC-V instructions.

3. Light PPU IP and interfacing to RISC-V

3.1 Light PPU IP design

The goal of this work is the design of an IP core for lightweight PPU (Posit Processing Unit) to be connected to a 64b RISC-V processor in the form of a co-processor with an extension of the Instruction Set Architecture.

Please note that the theory of Posit arithmetic, the structure of Posit numbers and their compression data benefit vs. classic integer and floating-point formats, particularly for DNN computation, have been already discussed by us in published works such as [2-4]. Therefore, the goal of this deliverable is on the design and verification of digital IPs supporting data compression for DNN thanks to the light support of Posits.

We focus on the compression abilities of posits by providing a co-processor with only conversions in mind, called light PPU, (as in Figure 3.1 below).

The difference between D2.6 and D2.2 is that D2.2 presents an AI accelerator IP giving full hardware support to posits number and hence using the IP in D2.2 the FPU of a RISC-V processor can be replaced by the Full PPU.

Instead, in D2.2 the idea is minimizing the circuit complexity overhead and hence the proposed IP supports only the data compression using Posits with Float to/from Posit translation, but operations have to be done still in the FPU.

We can convert binary32 floating point numbers to posit numbers with 16 and 8 total bits (and a variable number for the exponent since Posit 16,0 and Posit16,1 is considered).

In order to provide a circuit design for our light-PPU we considered several key points to simplify the final logic design:

1. IEEE floating point values are encoded in a module and sign-like representation while posits are encoded using 2's complement representation. Therefore, when converting from IEEE floats we just ignore the sign and build the positive posit. Then we use the sign to apply the 2's complement to the result if negative.
2. Given a Posit<16,0> the size of the regime spans from a minimum of 2 to a maximum of 15 bits. As a result, the mantissa size spans from a minimum of 0 to a maximum of 12 bits. This means that, given a 23-bit mantissa IEEE Float, the 8 least significant bits of the float are set to 0. The same concepts hold for Posit<8,0>.
3. We can build the Posit<X,0> regime arithmetically shifting an appropriate value by the X least significant bits of the FP32 normalized exponent. For Posit<16,0> we shift the signed integer represented by 2^{15} or $(8000)_{16}$ in hexadecimal notation. For Posit <8,0> we shift the signed integer represented by 2^7 .
4. Decoding the regime is particularly interesting since we need to employ a *find first set* module (or *find first unset*) to evaluate the regime length. The output of the *find first set* module is the index i of the highest set bit (discarding the sign if present). For Posit<16,0> the regime length is actually computed as $l=14-i$.

Figures 3.1, 3.2 and 3.3 show simplified versions of the light-PPU architecture.

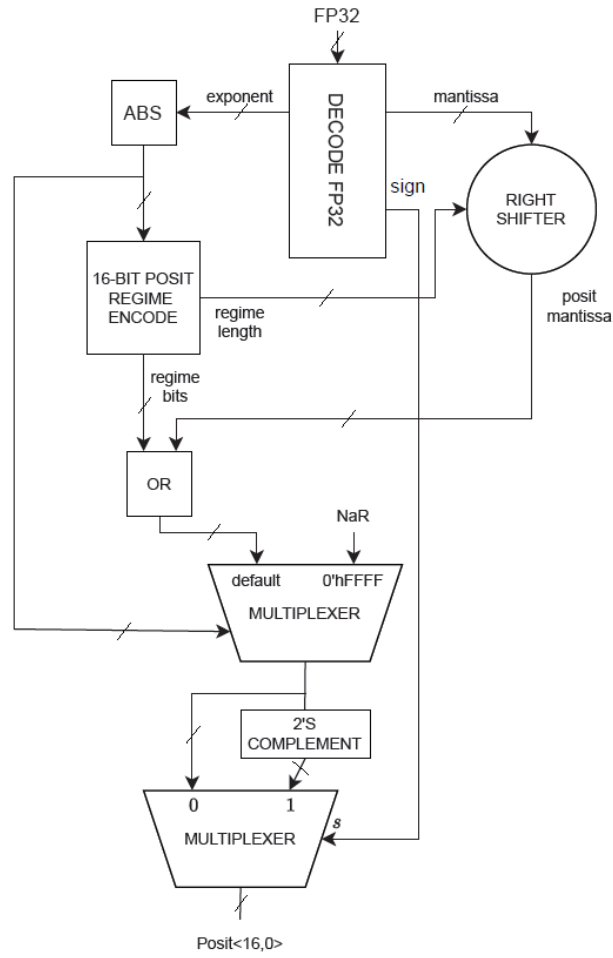


Figure 3.1: Logical circuit for the 32-bit floating point to posit16,0 converter

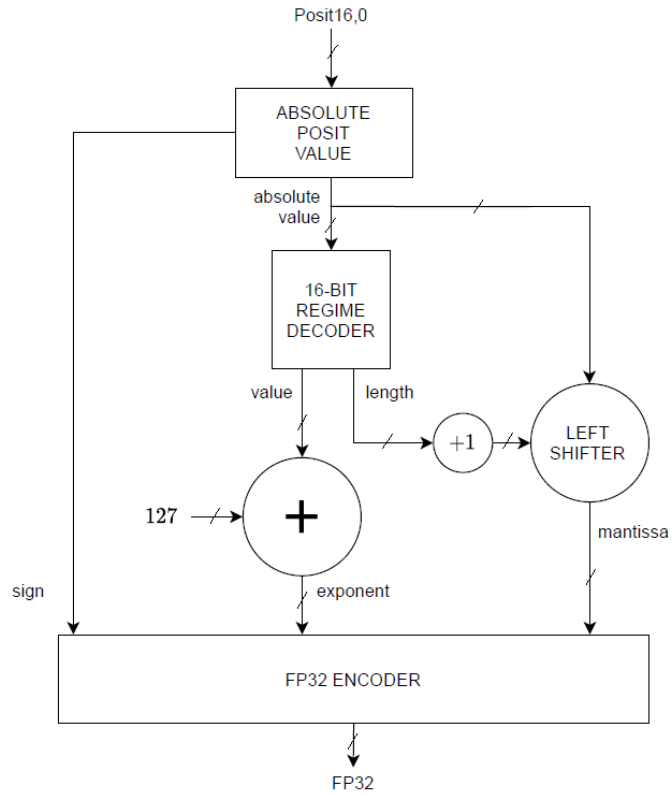


Figure 3.2: Logical circuit for the posit16,0 to 32-bit floating point converter.

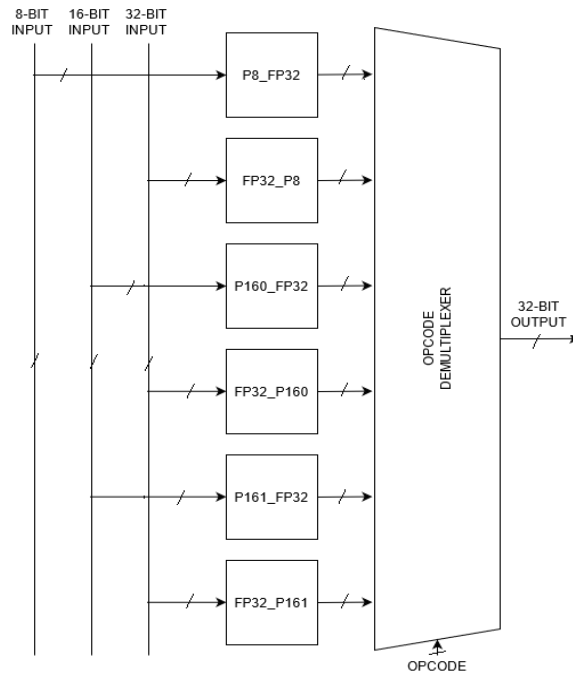


Figure 3.3: Overall architecture for the light PPU

3.2 Integration with RISC-V CPUs Ariane 64bits

The proposed co-processor, as shown in Figure 3.4, can be paired with a RISC-V core that already has a floating-point unit (e.g., the Ariane 64b RISC-V) without interrupting the existing pipeline.

On the other hand, we can use this unit to enable ALU computation of posit numbers with the posit-to-fixed conversion modules on a RISC-V core that does not support floating-point:

1. If the RISC-V processor embeds an FPU the light PPU can be used as a wrapper, providing a data compression by a factor up to 4, with little accuracy degradation. The cost of compression and decompression is the cost of converting a posit to a float and vice-versa.
2. If the RISC-V processor does not embed an FPU or we want to exploit only the ALU, the light PPU can function as a wrapper of fixed-point representation. Indeed, once we have converted between posit and fixed-point, the basic arithmetic operations can be computed just with the ALU. Also note that, for half of the posit domain, that is the [-1,1] range, the conversion between posit and fixed point is a simple left shift of 2 positions followed by a 0 padding on the most significant bits to reach desired fixed-point size.

We investigated the first use-case by outfitting a CVA6 core with our PPU co-processor and synthesizing it for a Xilinx Genesys 2 FPGA, resulting in a working RISC-V core capable of running a general-purpose Linux distribution.

The integration within the RISC-V core (Ariane CVA6 code, [7, 8]) can be done using the possibility to customize the instruction set.

The posit-based compression IP proposed in D2.6 can be integrated in addition to the Ariane integer ALU and in addition to the Ariane FPU.

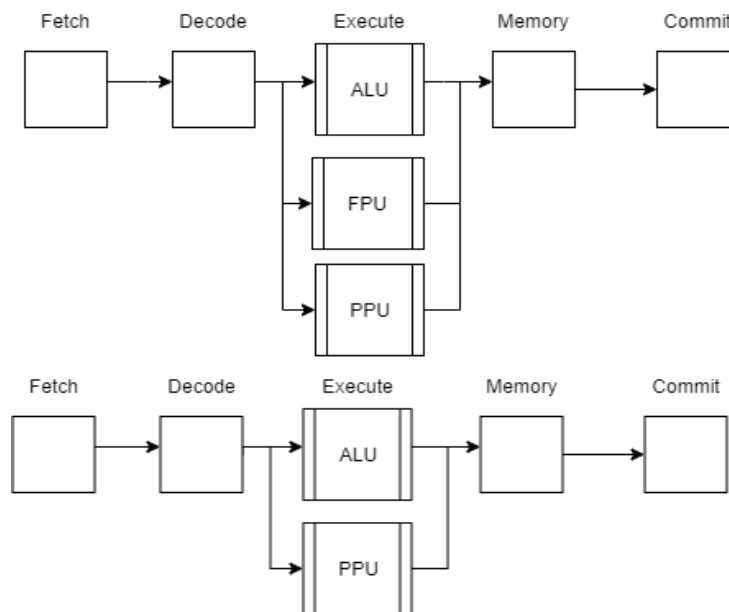


Figure 3.4: PPU possible integration modes within the RISC-V instruction set (with/without the FPU)

Table 3.1 Shows the implemented RISC-V instructions and opcodes for the new light PPU instructions. Note that we are leveraging the custom opcode space in the 6 least significant bits **0x0b**.

31	27	26	25	24	20	19	15	14	12	11	7	6	0	
1100000			00010			rs1		000		rd		0001011		FCVT.S.P8
1100000			00011			rs1		000		rd		0001011		FCVT.S.P16.0
1100000			00011			rs1		010		rd		0001011		FCVT.S.P16.1
1101000			00010			rs1		000		rd		0001011		FCVT.P8.S
1101000			00011			rs1		000		rd		0001011		FCVT.P16.0.S
1101000			00011			rs1		010		rd		0001011		FCVT.P16.1.S
1100000			00010			rs1		001		rd		0001011		FXCVT.H.P8
1100000			00011			rs1		001		rd		0001011		FXCVT.W.P16.0
1100000			00011			rs1		011		rd		0001011		FXCVT.L.P16.1
1101000			00010			rs1		001		rd		0001011		FXCVT.P8.H
1101000			00011			rs1		001		rd		0001011		FXCVT.P16.0.W
1101000			00011			rs1		011		rd		0001011		FXCVT.P16.1.L
1101000			00010			rs1		001		rd		0001011		FXCVT.P8.H
1101000			00011			rs1		001		rd		0001011		FXCVT.P16.0.W
1101000			00011			rs1		011		rd		0001011		FXCVT.P16.1.L
1100000			00010			rs1		100		rd		0001011		FCVT.P8.P16.0
1100000			00011			rs1		100		rd		0001011		FCVT.P16.0.P8
1101000			00011			rs1		111		rd		0001011		FCVT.P16.1.P16.0
1101000			00010			rs1		101		rd		0001011		FCVT.P16.1.P8
1100000			00011			rs1		110		rd		0001011		FCVT.P8.P16.1
1101000			00011			rs1		101		rd		0001011		FCVT.P16.0.P16.1

Table 3.1: ISA extension for the light PPU with conversions between 16-bit/8-bit posits and FP32/Fixed

4. Implementation results and benchmarks

4.1 Hardware implementation results

We chose the Xilinx Genesys 2 board for the hardware implementation (equipped with a Kintex 7 XC7K325T-2FFG900C FPGA component).

We chose this board to reduce the work required to construct our PPU inside a RISC-V core. We did, in fact, use the ARIANE RISC-V core that was originally built for this board.

The resulting design was then implemented on the same board.

This approach of using the same prototyping board also facilitates a fair comparison of the RISC-V without the Light PPU and the RISC-V with the Light PPU.

The Genesys 2 board can be connected to a host controller for configuration and diagnosis via an UART interface.

For the PPU component, we performed power, circuit complexity, and propagation delay (worst case combinatorial propagation delay of the PPU) reports:

1. Look-up table (LUT) utilization: 747/203800 (0.36%) LUTs used.
2. Component latency: 6.332ns (worst propagation delay).

Finally, the new instruction set architecture was merged into the Ariane RISC-V core and synthesized for the Xilinx Genesys 2.

The following key performance indicators (KPIs) were obtained:

1. Clock frequency: 125MHz (the same of the original Ariane design on the same board, so adding the Light PPU does not change the maximum achievable frequency)
2. Total power on FPGA components (Kintex 7): 2.056W of which the contribution of the Light PPU is less than 1%.
3. Look-up table (LUT) utilization: 63805/203800 (31.54%) LUTs used, of which only 0.36% is due to the added Light PPU.

To be noted that the target of Textarossa is the implementation of the IP on FPGA, not on a ASIC, and hence the KPI of keeping the same frequency of the original Ariane design in the range of hundred of MHz for FPGA is a good result. An ASIC implementation can be part of future evolutions of this IP design but is out of scope for Textarossa.

Moreover, the work on Textarossa in D2.2 and D2.6 refers to single line posit processing units. A vectorization of the Posit operators could be evaluated as a future extension of this work to become compliant with vectorized RISC-V instructions thus creating a Vector Processing Unit (VPU).

We validated the new ISA assuring the coherency between the two versions of cppPosit, compiled with or without the light PPU support. In the former case, the ISA instructions were emulated in the Spike simulator and the operation results were compared to the latter case in which posit operations were implemented completely in cppPosit.

We developed a test-bench using the same HDL language used for the light PPU unit to test the functionality of the newly introduced components. We tried every combination for the inputs and tested the outputs of conversion against our cppPosit software library. This means that we have also verified the correct handling of all the corner cases, such as the redundant negative zero and the redundant representations for the Infinity and Not-A-Number in IEEE 32-bit floats.

4.2 DNN Benchmarks

To assess the variation of accuracy with the compressed formats, we tested two well-known datasets often used in literature for benchmarking [2, 5, 6 10] (the MNIST dataset and the GTSRB datasets on the networks shown in Figure 4.1, pre-trained using 32-bit floats). We then tested the model using the posit compression to assess accuracy variations.

We used the following neural network (LeNet-5) as benchmark base for weight compression:

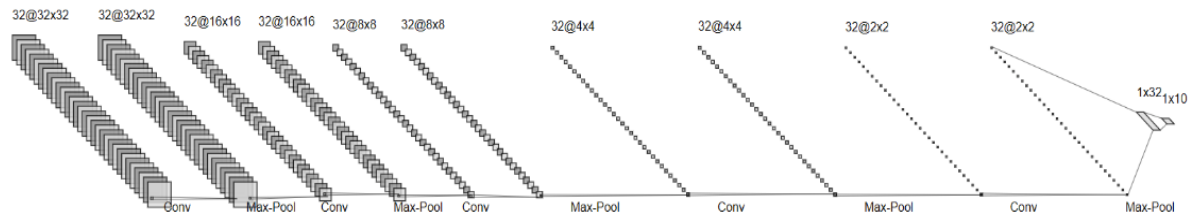


Figure 4.1: example neural network used for test of the light PPU data compression IP

We also give earlier accuracy results on this neural network in Figure 4.1 utilizing both posit and binary32 integers with varying workloads for completeness. This result was reached by using posit compression for the weights and computing using binary32 (fp32) format. We measured the relative speedup to a software-based posit compression.

As a result, we investigated the overall system compression times with the weights of a tiny LeNet-5 neural network, yielding the result displayed in Table 4.1.

LeNet		
	MNIST	GTSRB
FP32	98.83%	91.8%
posit(16, 1)	98.83%	91.8%
posit(16, 0)	98.50%	90.5%
posit(8, 0)	98.34%	90.4%

Table 4.1: Accuracy performance of Posit vs FP32 for different benchmark data sets (MNIST [5] and the German Traffic Road Sign Benchmark [6])

To assess the performance of the customized core we ran the same benchmarks used in the simulation phase on the ARIANE RISC-V core, equipped with the OpenPiton 12 Linux distribution (based on the ARIANE Linux 4.2).

As before, timing performance was measured using C++ internal software *chrono* directives. This approach involves the conversion between posits and floats at each operation (e.g. sum or multiplication).

As a result, for each operation, we are performing two more instructions for type conversion. Table 4.2 and 4.3 summarizes the results obtained with the evaluation of this trade-off. Note that the value obtained with IEEE FP32 is totally independent from the presence of the light PPU.

	w/ PPU (s)	wo/ PPU (s)	Speedup
posit 8,0	5.4	58.87	10.90
posit(16, 0)	11.6	64.54	5.56

Table 4.2: Real HW (FPGA) timing performance on a 10-layer convolutional neural network with and without the synthesized hardware PPU light support for the cppPosit library.

	Time (s)	DNN size (bytes)	Compression
IEEE FP32	2,1	224894	-
posit(16, 0)	11.6	112874	1.99
posit(8, 0)	5.4	56864	3.95

Table 4.3: Tradeoff between processing time and compression factor of Posit vs FP32 for DNN

We may instead take an entire posit network and convert it beforehand to IEEE FP32, in order to exploit compression as much as possible without slowing down actual image processing. This use case is relevant if we think about resource constrained environments where volatile memory is scarce (e.g. embedded or automotive systems) or when frequent transfer of network models are done (e.g. smartphones with recognition software). Moreover, these systems often use an adaptive approach where, depending on the surrounding environment (e.g. snow, night-time, off-road etc.), different machine learning models need to be loaded. This means that having multiple models on volatile storage can highly benefit from compression even when they are not actually loaded into main memory. In this case, we need to decompress the network into IEEE FP32 format only one time at the beginning of execution.

This will lead to a much slower start but a faster computation time (that is the same as IEEE FP32 in Table 4.3). Table 4.4 summarizes the results of this evaluation.

	Time (s)	DNN size (bytes)	Compression
IEEE FP32		224894	
posit(16, 0)	51.5s	112874	1.99
posit(8, 0)	49.8s	56864	3.95

Table 4.4: Trade-off between compression time of the network in Figure 4.1 and compression factor Processing time was obtained from the real hardware implementation

5. Conclusions and IP repository

In this work we dealt with the preliminary HDL design, using SystemVerilog, verification and synthesis of an IP for data compression exploiting posit arithmetic.

The key novelty is the hardware support for a new arithmetic format called Posit, which has been shown to have a strong compression effect, notably for ML and DNN applications: 4x compression for the same quality.

We designed the data compression IP using FPGA technology targeting the Xilinx Genesys 2 FPGA platform. Furthermore, we tested the IP design against a golden-model software library for posit arithmetic.

The data compression IP has been implemented in different Xilinx FPGA devices and it has been designed according to the specifications defined in D2.1, aiming to demonstrate its platform independence.

Future activities will expand the implementation to other platforms (e.g. the ALVEO U280 platform also selected by other partners of the project such as INFN).

Finally, we integrated the data compression IP within the ARIANE 64-bits 6-stage RISC-V core.

The IP can be accessed: https://bitbucket.org/federicorossifr/ppu_public/src/master/

The proposed work has been also submitted to the review of the scientific community and has been accepted on

Marco Cococcioni; Federico Rossi; Emanuele Ruffaldi; Sergio Saponara, “A Lightweight Posit Processing Unit for RISC-V Processors in Deep Neural Network Applications”, IEEE Transactions on Emerging Topics in Computing, 2022, vol. 10, n. 4

6. References

- [1] D21, “Consolidated specs of accelerators IPs”, Textarossa project, May 2022, Revised May 2023
- [2] Marco Cococcioni, Federico Rossi, Emanuele Ruffaldi, & Sergio Saponara, Benoit De Dinechin, “Novel Arithmetics in Deep Neural Networks Signal Processing for Autonomous Driving: Challenges and Opportunities”, IEEE Signal Processing Magazine, Volume: 38, Issue: 1, 2021
- [3] Marco Cococcioni, Federico Rossi, Emanuele Ruffaldi, Sergio Saponara, A Novel Posit-based Fast Approximation of ELU Activation Function for Deep Neural Networks, 2020 IEEE International Conference on Smart Computing (SMARTCOMP)
- [4] Marco Cococcioni, Federico Rossi, Emanuele Ruffaldi, & Sergio Saponara. (2021). A Lightweight Posit Processing Unit for RISC-V Processors in Deep Neural Network Applications. IEEE Trans on Emerging topics in Computing, <https://zenodo.org/record/7128760#.Y0ZfpC8QNbU>
- [5] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition." Proceedings of the IEEE, 86(11):2278-2324, November 1998.
- [6] J. Stallkamp, M. Schlipsing, J. Salmen, C. Igel, Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition, Neural Networks, Available online 20 February 2012, ISSN 0893-6080, 10.1016/j.neunet.2012.02.016.
(<http://www.sciencedirect.com/science/article/pii/S0893608012000457>) Keywords: Traffic sign recognition; Machine learning; Convolutional neural networks; Benchmarking
- [7] F. Zaruba, and L. Benini, “The Cost of Application-Class Processing: Energy and Performance Analysis of a Linux-ready 1.7GHz 64bit RISC-V Core in 22nm FDSOI Technology”, arXiv e-prints, 2019.
- [8] <https://github.com/openhwgroup/cva6>.
- [9] Marco Cococcioni; Federico Rossi; Emanuele Ruffaldi; Sergio Saponara, A Lightweight Posit processing Unit for RISC-V Processors in Deep Neural Network Applications, IEEE Transactions on Emerging Topics in Computing. 2022 | Volume: 10, Issue: 4
- [10] https://benchmark.ini.rub.de/gtsdb_dataset.html#:~:text=The%20German%20Traffic%20Sign%20Detection%20Benchmark%20is%20a%20single-image,it%20is%20introduced%20on%20the%20IEEE%20International%20Joint