

Article

Customizable vector acceleration in extreme-edge computing: a RISC-V software/hardware architecture study on VGG-16 implementation

Stefano Sordillo¹, Abdallah Cheikh¹, Antonio Mastrandrea¹, Francesco Menichelli¹, Mauro Olivieri^{1,*}

¹ DIET, Sapienza University of Rome

* Correspondence: mauro.olivieri@uniroma1.it

Abstract: Computing in the cloud-edge continuum, as opposed to cloud computing, relies on high performance processing on the extreme edge of the IoT hierarchy. Hardware acceleration is a mandatory solution to achieve the performance requirements, yet it can be tightly tied to particular computation kernels, even within the same application. Vector-oriented hardware acceleration has gained renewed interest to support AI applications like convolutional networks or classification algorithms. We present a comprehensive investigation of the performance and power efficiency achievable by configurable vector acceleration subsystems, obtaining evidence of both the high potential of the proposed microarchitecture and the advantage of hardware customization in total transparency to the software program.

Keywords: edge-computing, processors, hardware acceleration

1. Introduction

The cloud-edge continuum computing paradigm relies on the possibility of local processing in the edge of the IoT whenever it is convenient for reasons of energy efficiency, reliability, or data security. As a consequence, there is a gradual shift of artificial intelligence (AI) algorithm execution from the cloud down low power embedded IoT devices on the edge, to be used in real-time for example to take voice commands or extract image features, for biometric, security, or filtering purposes [5].

The resultant demand for very high processing speed on extreme edge computing devices turns into unprecedented design challenges, especially because of the usually limited energy budget. Therefore, the implementation of hardware acceleration on edge devices in the IoT hierarchy has become a major trend to reach the speed and energy efficiency requirements.

Vector computing acceleration was a major stream in high performance computing systems for decades and is gaining renewed interest in recent development in the supercomputing sector [22]. Yet, it is easy to note that the vector computing paradigm can also be applied to AI computing kernels that are run in embedded IoT devices on the edge. Nonetheless, the limited hardware budget usually available in edge devices makes it interesting to explore the possibility of configurable acceleration sub-systems to optimally exploit the available hardware resources according to the specific computation kernels being run during the application execution.

We implemented such exploration addressing the execution of the VGG-16 deep convolutional neural network inference, widely known for its image recognition performance as well as for the high computing power and storage demand. The VGG-16 execution is composed of consecutive layers having different computational characteristics. Therefore, it well represents a stress-test of the hardware micro-architecture with a time-variant workload profile. Our target micro-architecture is an open-source RISC-V [3] processor core supporting multi-threaded execution and featuring a highly customizable vector acceleration subsystem [23].

The contributions of this work to the reader interested in advanced embedded system design for IoT extreme-edge computing, are manifold:

- we report the quantitative evidence of the trade-offs in vector co-processor design and configuration targeting simple edge-computing soft-cores;
- we present details on the small custom RISC-V compliant instruction extension sufficient to support typical vector operations in a tiny soft-core;
- we present a complete yet very simple library of intrinsic functions to support application development, and we discuss the full detail of source code exploiting the co-processor instructions in each VGG-16 layer execution;
- we give insights into the open-source Klessydra processor core microarchitecture.

The rest of this article is organized as follows: Section 2 covers the related works on hardware acceleration for embedded computing on the IoT edge, including configurable solutions, Section 3 introduces the Klessydra T1 processor soft-core featuring configurable hardware acceleration subsystem, Section 4 describes the fundamental features of the VGG-16 application case and its implementation on Klessydra T1. Section 5 reports and discusses the results obtained for the different sub-parts of the chosen application cases, and Section 6 summarizes the outcomes of the work.

2. Related works

Several previous works reported the design of hardware accelerated microcontroller cores implemented in edge-computing silicon chips. In [6], a RISC-V processor with DSP hardware support is presented, targeting near-threshold voltage operation. The Diet-SODA design implements a similar approach by running its DSP accelerator in near-threshold regime [7]. In [8,9,10,11] application specific accelerators are reported, based on highly parallel operation and minimized off-chip data movements for energy efficiency.

All of the above works focus on silicon implementation of units tailored to specific computations. As opposed to this view, the proposed hardware architecture study is independent of technology assumptions, such as the supply voltage, and addresses any physical implementation, particularly soft-cores on commercial FPGA devices, in the view of exploiting application-driven configurability. Regarding FPGA-based implementations, in [12] the authors present a cluster of RISC-V cores connected to a tightly-coupled scratchpad memory and a special purpose engine dedicated to convolutions only. Thanks to FPGA implementation, the convolution engine can be configured at synthesis time to optimize the execution of each convolutional layers, yet exhibiting a severe performance degradation when executing layers it was not built to optimize.

A recently published work [13] presents a SIMD configurable CNN coprocessor connected to a 32-bit RV32IM RISC-V processor. Compared to the proposed Klessydra configuration that consumes almost the same amount of LUTs, the design in [13] performs significantly slower.

In [14] the authors present a coprocessor soft-core at the edge of IoT, designed to be energy efficient in executing CNN as well as other machine learning algorithms. In particular, they explore the potential impact of data parallelism on the energy efficiency due the increased memory bandwidth. In our study, memory traffic as well as the memory static power consumption are taken into account in energy estimations.

The works in [15][16] present a pipelined CNN coprocessor capable of accelerating convolutions based on the extremely high parallelism in the coprocessor, yet limited to convolutional computation kernels.

In [17] the authors present different coprocessor configurations integrated with a parallel cluster of RISC-V cores and evaluated which of the configurations is the fastest and most energy efficient. They introduce special co-processing cores dedicated to the standard instruction subset RV32M, without exploring more sophisticated co-processor operations.

In [18] the authors provide a DCNN accelerator for IoT. The accelerator itself is designed to work with 3x3 kernels, and being not configurable, in order to support larger kernels they use a technique called kernel decomposition, which in fact increases the waste in computational resources and decreases in the energy efficiency, similarly to the convolution engine in [12].

The coprocessor architecture proposed in this work is general purpose in nature, being based on vector operations, and can be tailored to support a given computation kernel in the most efficient way. Our work builds on the preliminary developments reported in [2,4] and complements the analysis presented in [23].

3. The Klessydra T1 customizable architecture

Hardware microarchitecture

Klessydra is a family of open-source, RISC-V compliant and PULPino [20] compatible cores, which includes basic processors (T0 sub-family), hardware accelerated processors (T1 sub-family), and fault-tolerant processors (F0 sub-family) [21]. A characteristic feature of all Klessydra cores is the hardware support for interleaved multi-threading on a single core [1].

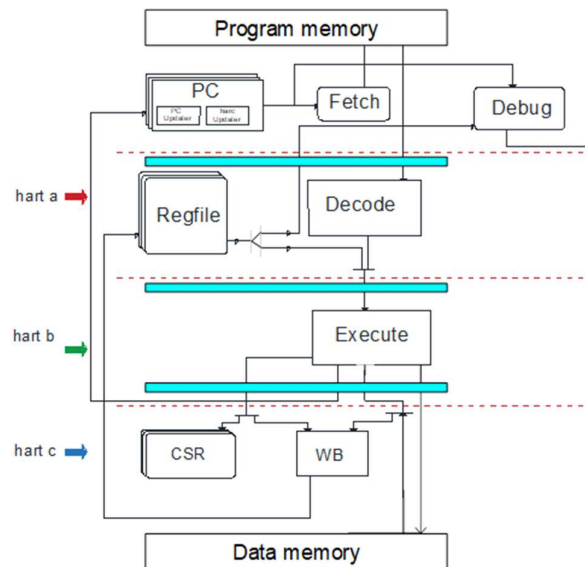


Figure 1. Klessydra T0 core microarchitecture

The hardware accelerated T1 cores are an extension of the basic T0 core, that is sketched in Figure 1. The T0 microarchitecture resembles a classic four-stage RISC pipeline, except for having multiple Program Counters to support multi-threading, and replicated register files and Control/Status Registers to keep the state of multiple threads. In each clock cycle a different Program Counter is used for instruction fetching, on a rotation basis. As a result, instructions belonging to different threads of execution are interleaved in the core pipeline, so that it is never possible that any two instructions in the pipeline manifest any register, structural or branch dependency. The only dependency between two threads can occur on explicit shared memory access, which is responsibility of the programmer. The supported number of interleaved threads is a parameter of the synthesizable RTL code of the core.

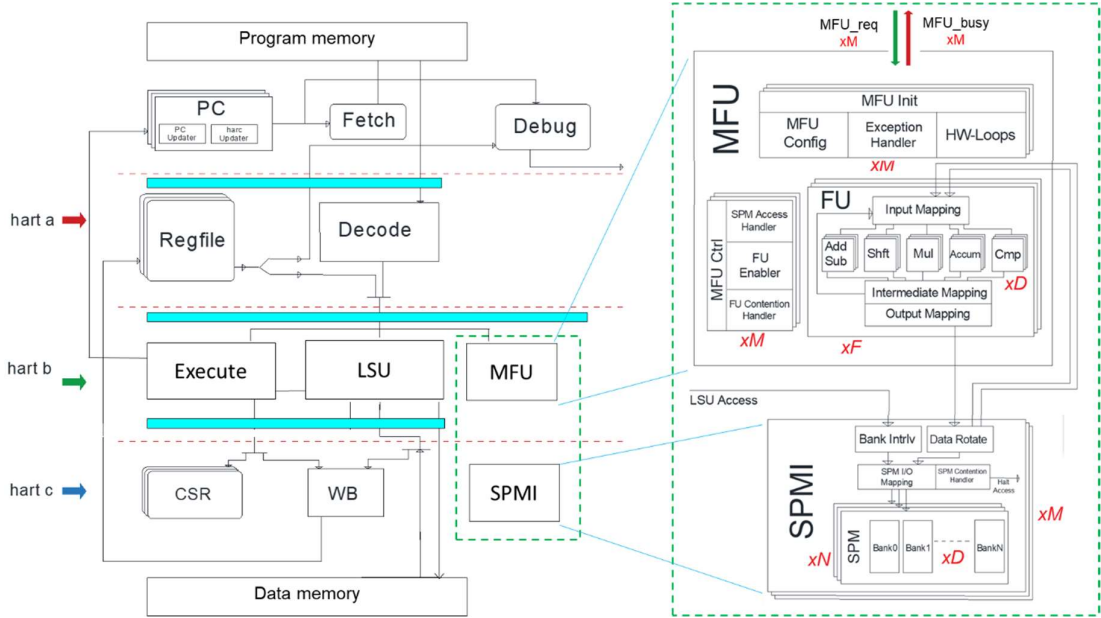


Figure 2. Klessydra T1 core microarchitecture

The T1 microarchitecture (Figure 2) is derived from the T0 by adding two execution units, namely the Load-Store Unit (LSU) and the Vector Co-processing Unit (VCU), the latter being internally comprised of Multi-Purpose Functional Units (MFU) and Scratch-Pad Memory Interface (SPMI).

At the instruction level, the T13 architecture supports the parallel execution of instructions of different types, belonging to the same hart. In fact, the LSU works in parallel with the other units when executing memory store instructions, that cannot cause a write-back conflict on the register file. The MFU is allowed to read operands from the register file but can only write its results to local scratchpad memories (SPMs), thus keeping the SPMs and the Registerfile decoupled and allowing parallel execution between instructions writing to each of these memories simultaneously. Data transfers to/from the data memory from/to the SPMs are managed by the LSU via dedicated instructions.

The MFUs execute vector arithmetic instructions, whose latency is proportional to the vector length. In an in-order IMT pipeline, a hart requesting access to the busy MFUs may result in stalling the whole pipeline, stalling other harts that may not need to access the MFU. To circumvent this, in the T13 architecture, the waiting hart executes a self-referencing jump until the MFU becomes free, avoiding unnecessary stalls of harts that are independent from the MFU being busy. Figure 3 demonstrates a cycle accurate diagram of the mechanism.

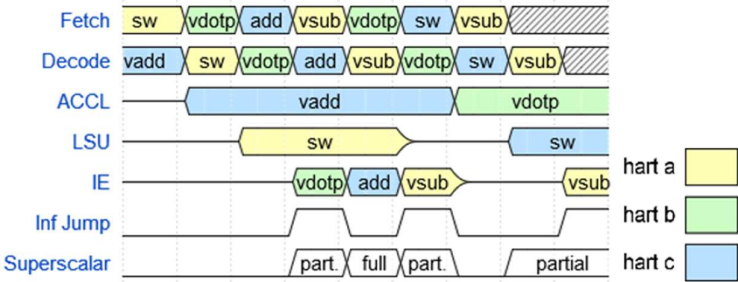


Figure 3. Hart interleaving and hart stall timing diagram

When deploying Klessydra T1 in a IoT edge device, one can configure the number of parallel lanes D in the MFU, the number of MFUs F, the SPM capacity, the number of SPMs N in each SPMI,

the number of SPMs M , as well as the way the MFUs and SPMs are shared between the harts. Representative configurations are the following:

- **Thread-Shared coprocessor:** All harts in the core share a single MFU/SPM subsystem. Harts in this scheme are required to execute an infinite jump when trying to access the MFU when it's busy. In this approach, instruction level parallelism is limited to occur only between coprocessor instructions writing to the SPM and non-coprocessor instructions writing to the main memory or registerfile. To mitigate the delays on a hart executing an infinite jump, the coprocessor here may exploit pure data level parallelism (DLP) acceleration, by multi-lane SIMD execution.
- **Thread-Dedicated coprocessor:** Each hart is appointed a full MFU/SPM subsystem, eliminating inter-hart coprocessor contention and allowing inter-coprocessor parallel execution. Stalls can only happen if the next instruction of the same hart that is using the MFU requests an MFU operation. DLP by multi-lane SIMD execution can still be exploited in this approach, but also thread level parallelism (TLP) by fully symmetric MIMD execution, allowing execution of multiple vector instructions in parallel.
- **Thread-Dedicated SPMs with a Shared MFU:** The harts here maintain a dedicated SPM address space, yet they share the functional units in the MFU. This scheme still allows inter-hart parallel execution of coprocessor instructions, provided they use different internal functional units of the MFU (e.g., adder, multiplier). Harts that request a busy internal unit in the MFU will be stalled, and their access will be serialized until the contended unit becomes free, while harts that request a free functional unit can work in parallel with the other active harts in the MFU. DLP by multi-lane SIMD execution can still be exploited in this approach, but also TLP by heterogeneous MIMD execution.

Error! Reference source not found. summarizes the design parameters and corresponding configurations, whose names will be used as references in reporting performance results.

Table 1. Summary of explored hardware configurations

M	F	D	Execution paradigm
1	1	1	SISD
1	1	2,4,8	SIMD
3	3	1	Symmetric MIMD
3	3	2,4,8	Symmetric MIMD + SIMD
3	1	1	Heterogenous MIMD
3	1	2,4,8	Heterogenous MIMD + SIMD

Programming paradigm

By default, a Klessydra core runs the maximum number of hardware threads (which is a synthesis parameter) allowed by the microarchitecture. The function *Klessydra_get_coreID()* can read the id number of the thread executing the function from the Mhartid CSR register, so this allows to distinguish threads and possibly have each thread to execute a different piece of program. Figure 4 shows a generic C program skeleton in which each of three threads executes its own instruction flow. The functions *sync_barrier_thread_registration()* and *sync_barrier()* allow implementing a synchronization barrier by based on inter-thread software interrupts, to synchronize thread execution at certain points of the program.

```

sync_barrier_thread_registration(); //Executed by all threads
if (Klessydra_get_coreID()==0){
    // thread_0 subroutine
}
if (Klessydra_get_coreID()==1){
    // thread_1 subroutine
}
if (Klessydra_get_coreID()==2){
    // thread_2 subroutine
}
sync_barrier(); //Executed by all threads

```

Figure 4. Code for multi-threaded execution on Klessydra-T13

Inter-thread data transfers may happen via shared global static variables allocated in the main data memory or, in the case of a shared coprocessor configuration, via shared SPM address space.

The custom instruction extension supported by the VCU and LSU is summarized in Table 2. The instructions supported by the coprocessor sub-system are exposed to the programmer in the form of very simple intrinsic functions, fully integrated in the RISC-V *gcc* compiler toolchain. The instructions implement vector operations without relying on a vector register file, but rather on a memory space mapped on the local SPMs, for sake of flexibility. The programmer can move vector data in any point of the SPM address space with no constraint except the total capacity of the SPMs, which in turn is a parameter of the microarchitecture design. The vector length applied by MFU operations is encoded in a user accessible custom control/status register (CSR) named MVSIZ.

Table 2. RISC-V instruction set custom extension for Klessydra-T processors

<i>Assembly syntax – (r) denotes memory addressing via register r</i>	<i>Function declaration</i>	<i>Short description</i>
<code>kmemld (rd), (rs1), (rs2)</code>	<code>kmemld((void*) rd, (void*) rs1, (int) rs2);</code>	<i>load vector into scratchpad region</i>
<code>kmemstr (rd), (rs1), (rs2)</code>	<code>kmemstr((void*) rd, (void*) rs1, (int) rs2);</code>	<i>store vector into main memory</i>
<code>kaddv (rd), (rs1), (rs2)</code>	<code>kaddv((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>adds vectors in scratchpad region</i>
<code>ksubv (rd), (rs1), (rs2)</code>	<code>ksubv((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>subtract vectors in scratchpad region</i>
<code>kvmul (rd), (rs1), (rs2)</code>	<code>kvmul((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>multiply vectors in scratchpad region</i>
<code>kvred (rd), (rs1)</code>	<code>kvred((void*) rd, (void*) rs1);</code>	<i>reduce vector by addition</i>
<code>kdotp (rd), (rs1), (rs2)</code>	<code>kdotp((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>vector dot product into register</i>
<code>ksvaddsc (rd), (rs1), (rs2)</code>	<code>ksvaddsc((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>add vector + scalar into scratchpad</i>
<code>ksvaddrf (rd), (rs1), rs2</code>	<code>ksvaddrf((void*) rd, (void*) rs1, (int) rs2);</code>	<i>add vector + scalar into register</i>
<code>ksvmulsc (rd), (rs1), (rs2)</code>	<code>ksvmulsc((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>multiply vector + scalar into scratchpad</i>
<code>ksvmulrf (rd), (rs1), rs2</code>	<code>ksvmulrf((void*) rd, (void*) rs1, (int) rs2);</code>	<i>multiply vector + scalar into register</i>
<code>kdotpps (rd), (rs1), (rs2)</code>	<code>kdotpps((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>vector dot product and post scaling</i>
<code>ksrlv (rd), (rs1), rs2</code>	<code>ksrlv((void*) rd, (void*) rs1, (int) rs2);</code>	<i>vector logic shift within scratchpad</i>
<code>ksrav (rd), (rs1), rs2</code>	<code>ksrav((void*) rd, (void*) rs1, (int) rs2);</code>	<i>vector arithmetic shift within scratchpad</i>
<code>krelu (rd), (rs1)</code>	<code>krelu((void*) rd, (void*) rs1);</code>	<i>vector ReLu within scratchpad</i>
<code>ksvslt (rd), (rs1), (rs2)</code>	<code>ksvslt((void*) rd, (void*) rs1, (void*) rs2);</code>	<i>compare vectors and create mask vector</i>
<code>ksvslt (rd), (rs1), rs2</code>	<code>ksvslt((void*) rd, (void*) rs1, (int) rs2);</code>	<i>compare vector-scalar and create mask</i>
<code>kvcv (rd), (rs1)</code>	<code>ksrlv((void*) rd, (void*) rs1);</code>	<i>copy vector within scratchpad region</i>
<code>csr MVSZIE, rs1</code>	<code>mvsize((int) rs1);</code>	<i>vector length setting</i>
<code>csr MVTYPE, rs1</code>	<code>mvtype((int) rs1);</code>	<i>element width setting (8,16,32 bits)</i>
<code>csr MPSCLFAC, rs1</code>	<code>mpsclfac((int) rs1);</code>	<i>post scaling factor (kdotpps instruction)</i>

4. VGG-16 implementation on Klessydra T1

Implementation workflow

VGG-16 is a deep Convolutional Neural Network (CNN) used in computer vision for classification and detection tasks, consisting of 13 convolutional layers, 5 maxpooling layers, 2 fully-connected layers and one output/softmax layer. The original VGG-16 can label a 224x224 pixel RGB image to one class out of 1000, using approximately 554MB space for 32-bit floating-point weights and bias values.

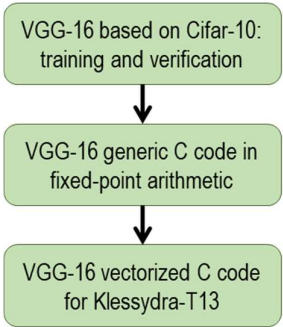


Figure 5. Workflow for the VGG-16 implementation

In the view of a realistic IoT edge embedded scenario, we implemented a VGG-16 derivation based on the widely known CIFAR-10 dataset, targeting 10 classes and 32x32 pixel RGB images and requiring 135 MB for weights and bias values. Table 3 reports the breakdown of the inference algorithm layers constituting the Cifar-10 VGG-16. The layers 19 to 21 do not compute operations on matrices, rather they implement dot-product operations between vectors of different sizes, similarly, layer 22 implements a Softmax function on a vector of length 10.

Table 3. Cifar-10 VGG-16 inference layers

Layer number	Computation type	Matrix size
1	<i>Convolution</i>	32x32
2	<i>Convolution</i>	32x32
3	<i>Max Pool</i>	16x16
4	<i>Convolution</i>	16x16
5	<i>Convolution</i>	16x16
6	<i>Max Pool</i>	8x8
7	<i>Convolution</i>	8x8
8	<i>Convolution</i>	8x8
9	<i>Convolution</i>	8x8
10	<i>Max Pool</i>	4x4
11	<i>Convolution</i>	4x4
12	<i>Convolution</i>	4x4
13	<i>Convolution</i>	4x4
14	<i>Max Pool</i>	2x2
15	<i>Convolution</i>	2x2
16	<i>Convolution</i>	2x2
17	<i>Convolution</i>	2x2
18	<i>Max Pool</i>	1x1
19	<i>Fully connected</i>	512x512
20	<i>Fully connected</i>	4096x4096
21	<i>Fully connected</i>	4096x4096
22	<i>Softmax</i>	10

Figure 5 illustrates the workflow to implement a Cifar-10 VGG-16 application on the Klessydra processor platform. Notably, since the target hardware platform supports fixed-point arithmetic, we based the implementation on fixed-point weights and data. We set the integer part to 11 bits and the fractional part to 21 bits, which gave us very good quality of output results, yet it is inessential to the performance results. Further algorithmic optimizations, such as quantization and compression techniques, are not in the scope of the present work. The learning and verification phase of the network in fixed point arithmetic was done via Matlab Deep Learning Toolbox. In order to be able to exploit the C language intrinsic functions of the Klessydra platform, the original Matlab code for VGG-16 was ported to C code. This generic C code implementation was used as the basis for the subsequent vectorization to exploit the hardware co-processor, and it was also used to run the same algorithm on the reference platforms used for performance comparison.

Generic fixed-point C code porting

The generic C code used for convolutional layers is reported in Figure 6. Image convolutions are implemented using the zero-padding technique: the feature map (FM) matrix is converted into a new matrix having two additional rows and columns of zeros on its borders, to avoid having filter elements without corresponding pixel values when the centroid element of the 3x3 kernel slides along the borders. As a general feature of the implementation, multiplications always need a pre-scaling and post-scaling operation in order to re-align the fixed-point representation of the result. The *convolution2D()* function performs the pre-scaling when creating the zero-padded matrix and also pre-scales the kernel values. The convolution is carried out by nested for loops, by which the Kernel map (KM) matrix slides across the input image with a stride of one element. The partial result of each multiplication is pre-scaled and added to the corresponding output pixel, completing the multiply and accumulate step. After the convolution is complete, a bias value is added to the output feature map, and the ReLU non-linear activation function is executed across all the matrix elements to conclude the convolutional layer.

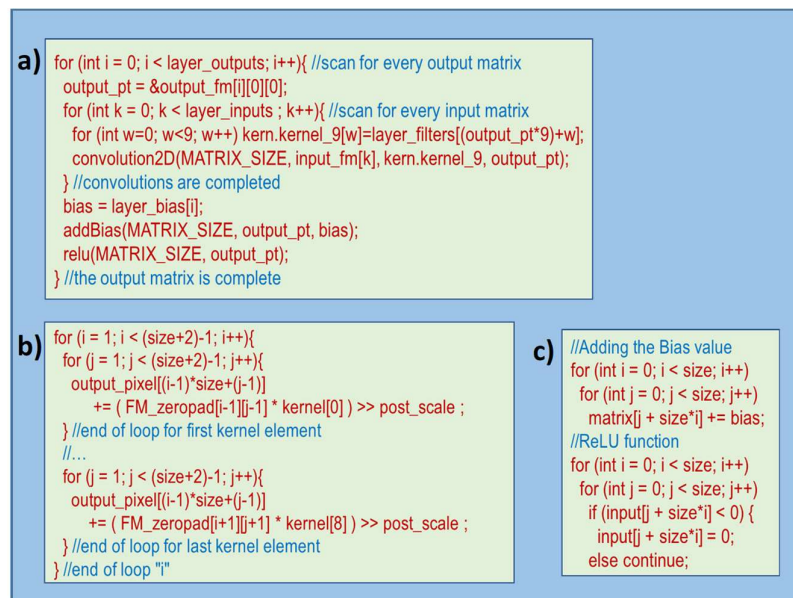


Figure 6. (a) Convolutional layer in generic C code; (b) Convolution2D function inner operations; (c) Bias addition and ReLU function inner operations

Figure 7 reports the C code adopted for Maxpool layers. The Maxpool layer halves the width and height of the FMs, sliding across them a 2x2 window, with a stride equal to two, filtering all the values except for the highest of the batch. In this way the most important features detected from the image are passed at the successive layers.

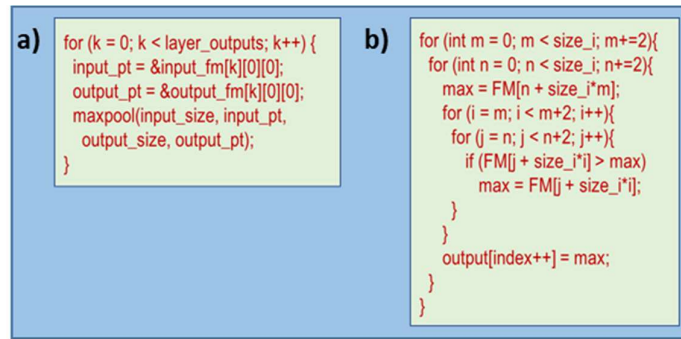


Figure 7. (a) Maxpool layer in generic C code; (b) Maxpool function inner operations

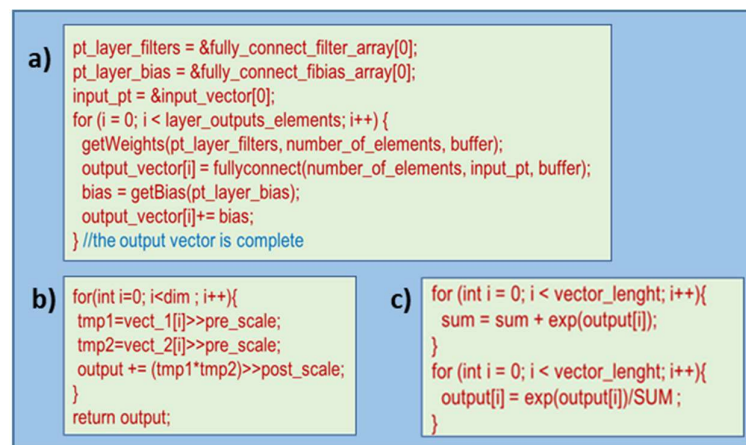


Figure 8. (a) Fully-connected layer in generic C code; (b) Fully-connect inner operations; (c) Softmax layer inner operations

The last three layers of the network are Fully Connected, corresponding to the code in Figure 8. The fully-connected layer is implemented by a dot-product, doing the pre-scaling of the inputs and post scaling of the results from every multiplication, needed for fixed point alignment. This is accomplished by the *fullyconnect()* function after putting the weights into local buffers and adding a bias to the output value. The results are passed through a Softmax layer, in which the network produces the classification of the image with a given probability.

Vectorized C code implementation

Program code vectorization targeting the Klessydra intrinsic function library is based on two types of intervention: data movement to efficiently exploit the scratchpad memory sub-system, and vector arithmetic operation exploiting the accelerator functional unit.

A loop of *kmemld()* functions transfer the FM and KMs operands into two SPMs, that we refer to as *spmA* and *spmB*, from the main memory. To implement zero-padding, when loading the feature maps into *spmA*, we first reset the SPM content to zero and then proceed with loading bursts of data from the FM rows, with exact offsets that grant the correctness of zero-padding. Figure 9(a) displays the code executed to set up the FM in *spmA*. The offsets added to the pointers passed to the *Kmemld()* function allow for the implementation zero-padding. The *ksrav()* function implements fixed-point pre-scaling by performing an arithmetic right shift operation of a vector. It requires a pointer to the vector, a pointer to store the resulting vector and a shift amount. Figure 9(b) similarly shows the

loading and pre-scaling of the 9-element KM into spmB and also the calling sequence of the *convolution2D()* function.

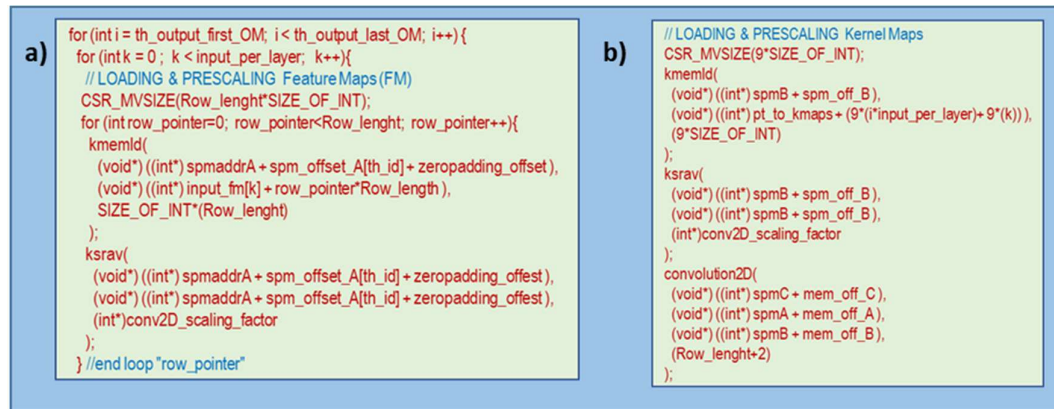


Figure 9. (a) Zero-padded, pre-scaled FM setup in SPM; (b) Pre-scaled KM collection in SPM and calling sequence of *convolution2D()*

The *Convolution2D()* function requires the addresses of the FM and KM first elements in spmA and spmB, an address pointing to a region in spmD for temporary value storage, and the address to store the output matrix in spmC. Figure 10 reports the internal operations, which are built upon knowing which vectors are to be multiplied as the kernel map slides across all the input map pixels. Taking into account which elements will be multiplied when the kernel completely slides across a row of the FM, and the fact that this process is replicated for every row, we can multiply the FM row values with the corresponding scalar from the KM, and update the output matrix (OM) row with a vector of partial results. This process is straightforward and allows to fully exploit the vector coprocessor capabilities by using matrix rows as vector operands.

```

CSR_MVSIZE(Row_size);
for(i=Zeropad_off; i Row_size-Zeropad_off; i++){
  k_element=0;
  for(FM_row_pointer=Zeropad_off; FM_row_pointer <= Zeropad_off; FM_row_pointer++){
    for(column_offset=0; column_offset < kernel_size; column_offset++){
      FM_offset = (i+FM_row_pointer)*Row_size+column_offset;
      ksvmulsc(SPM_D, (SPM_A+FM_offset), (SPM_B + k_element++));
      ksrav(SPM_D, SPM_D, scaling_factor);
      OM_offset = (Row_size*i)+Zeropad_off;
      kaddv((SPM_C+OM_offset), (SPM_C+OM_offset), SPM_D);
    }
  }
}

```

Figure 10. Convolution2D inner loops operations

Referring to Figure 10, after setting the vector length, the loop with index “i” scans the rows of the output matrix (OM); the *FM_row_pointer* loop and the *column_offset* loop iterate three times each to cover the necessary vector-scalar product required for the 3x3 kernel matrix. The *FM_offset* variable points to the proper FM row in spmA, from which the source vector is fetched. The *ksvmulsc()* function performs the scalar-vector multiplication between an FM row vector and a KM scalar, and the result is post-scaled by the *ksrav()* function for fixed-point alignment. The *kaddv()* function performs the vector addition, updating the OM row in spmC.

After the convolutions are done, the OM is updated with the addition of the bias value (Figure 11(a)). A *kmemld()* is required to have the single scalar value in the scratchpad memory, then the whole matrix is updated by *ksvaddsc_v2()*, which performs the vector plus scalar operation and includes a fourth parameter to adjust the vector length prior to doing the calculation.

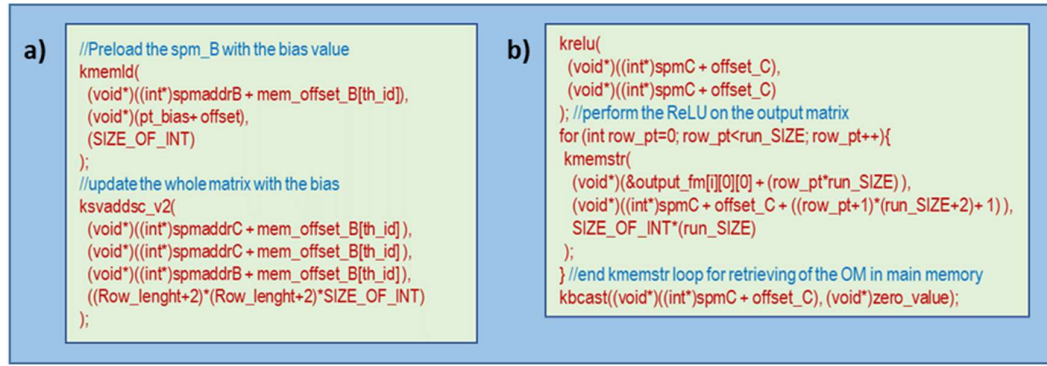


Figure 11. (a) Adding the bias to the Output; (b) Matrix and applying ReLu function

As the last part of the convolutional layers, the ReLU non-linear function is applied to all the OM pixels, which is stored back in main memory. The SPM region is cleared for the next iteration of the loop by broadcasting a zero value into the target memory region with *kbcst()* (Figure 11(b)).

The maxpooling layer is executed on the OM in main memory, through conventional scalar instructions, following the same implementation of the generic C code.

The fully-connected layer is comprised of a computation kernel based on dot products (Figure 12(a)). The source vector is moved into *spmA* as a single burst of data using the *kmemld()* function, and pre-scaled by *ksrav()*. A loop handles the properly transposed loading of the neurons parameters into *spmB*. The two vectors in the SPMs are processed by the dot-product function *kdotpps()*, which includes a post-scaling of the product before accumulation for fixed point alignment.

After the end of the loop, the vector of bias values is moved into *spmD* then added to the output vector of the layer. The result vector is processed by the *krelu()* function, and then it is stored back to the main memory. The *kbcst()* function clears the *spmC* memory space (Figure 12(b)).

The softmax layer is executed in main memory through conventional scalar instructions, with the same implementation of the generic C code.

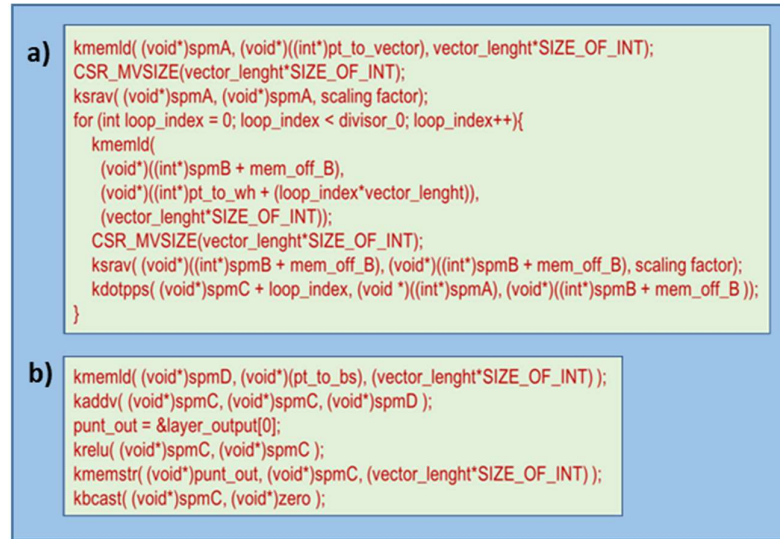


Figure 12. Fully-connected layer operations. (a) dot-product kernel; (b) Bias addition and ReLu.

The exact execution of the vectorized VGG-16 inference program running on Klessydra T13 cores was verified by comparing the full output produced by RTL simulation against the general purpose VGG-16 code running on an X86 server.

5. Performance and Power analysis

Setup

All Klessydra cores are compatible with the PULPino processor platform [20]. Yet, the original PULPino memory subsystem cannot support the execution of the full VGG-16 inference algorithm, which requires 255 MB storage for the constant data consisting of the neural network weights, and at least 1 MB memory space for global and local variables. Thus, we extended the PULPino memory sub-system to include 256 MB of addressable physical data memory, partitioned into a 1 cycle latency 1 MB RAM to be mapped on the FPGA BRAM, and a 6 cycle latency 255MB space mapped on an external flash memory device, connected via SPI interface. The program memory is 32 KB mapped in the FPGA BRAM.

The modified PULPino platform featuring Klessydra T13 processor cores was synthesized on a Kynntex7 FPGA board using the Vivado tool flow. Table 4 reports the hardware resource utilization and the maximum clock frequency results for all the processor configurations under analysis.

Table 4. Area and frequency summary of the Klessydra-T cores connected to 1MB Data Mem,

Configuration	Area Utilization					Top Freq. MHz
	FF	LUT	DSP	B-RAM	LUT-RAM	
SISD (M=1,F=1,D=1)	2482	7083	11	88	264	132.1
Pure SIMD (M=1,F=1,D=2)	2664	9010	15	88	264	127.0
Pure SIMD (M=1,F=1,D=4)	3510	11678	23	88	264	125.5
Pure SIMD (M=1,F=1,D=8)	4904	18531	39	88	264	112.6
Symmetric MIMD (M=3,F=3,D=1)	3509	10701	19	120	264	114.2
Symmetric MIMD+SIMD (M=3,F=3,D=2)	4659	16556	31	120	264	113.9
Symmetric MIMD+SIMD (M=3,F=3,D=4)	6746	27485	55	120	264	108.9
Symmetric MIMD+SIMD (M=3,F=3,D=8)	11253	52930	103	120	264	96.3
Heterogenous MIMD (M=3,F=1,D=1)	3025	10655	11	120	264	119.9
Heterogenous MIMD+SIMD (M=3,F=1,D=2)	3741	17161	15	120	264	115.7
Heterogenous MIMD+SIMD (M=3,F=1,D=4)	4767	25535	23	120	264	110.4
Heterogenous MIMD+SIMD (M=3,F=1,D=8)	7303	48066	39	120	264	91.5
No Accel T0 Core	1409	4079	7	72	176	194.6
RI5CY	1307	6351	6	72	0	65.1
Zeroriscy	1605	2834	1	72	0	77.2

The VGG-16 inference fixed-point code was also implemented on the following alternative computing systems, to accomplish a comprehensive comparative analysis:

- An FPGA board featuring a soft-processor comprised of the extended PULPino platform equipped with the DSP-accelerated RI5CY core, reaching 65 MHz clock frequency;
- An FPGA board featuring a soft-processor comprised of the extended PULPino platform equipped with a Zeroriscy core [19], reaching 77 MHz clock frequency;
- An STM32 single board computer featuring an 84 MHz ARM Cortex M4 core with DSP extension, 96 KB data memory;
- A Raspberry-PI 3b+ single board computer featuring a 1.4 GHz ARM Cortex A53 quad-core CPU, 16 KB L1 cache and 512 KB L2 cache, 1 GB LPDDR2 main memory;
- An x86 single board computer featuring a 3 GHz exa-core, 12-thread i7 CPU, 384 KB L1 cache, 1.5 MB L2 cache, 9 MB LLC, 8 GB DDR4 main memory.

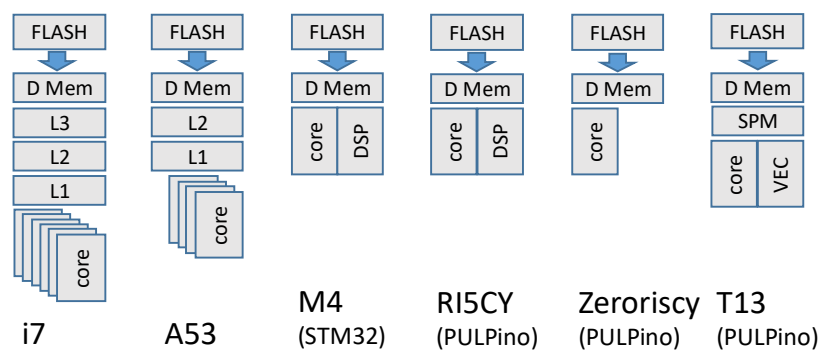


Figure 13. System architecture organization of the compared boards

The system architecture organization corresponding to the devices under comparison are sketched in Figure 13. The read-only storage space dedicated to the VGG-16 weights is hosted by an SPI-connected Flash memory expansion board in all the considered architectures, and the weights are preemptively loaded into the main RAM space for the inference algorithm execution.

Results

The first phase of performance analysis targeted the detailed account of the performance of each coprocessor hardware microarchitecture.

Table 5. Absolute execution time [s]. Best performing coprocessor configurations are highlighted for each layer.

Layer	SISD (M=1, F=1, D=1)	Pure SIMD (M=1, F=1, D=2)	Pure SIMD (M=1, F=1, D=4)	Pure SIMD (M=1, F=1, D=8)	Symm MIMD (M=3, F=3, D=1)	Symm MIMD +SIMD (M=3, F=3, D=2)	Symm MIMD +SIMD (M=3, F=3, D=4)	Symm MIMD +SIMD (M=3, F=3, D=8)	Heter. MIMD (M=3, F=1, D=1)	Heter. MIMD +SIMD (M=3, F=1, D=2)	Heter. MIMD +SIMD (M=3, F=1, D=4)	Heter. MIMD +SIMD (M=3, F=1, D=8)	Non Accel T0 core
1	0.057	0.037	0.029	0.023	0.027	0.022	0.022	0.021	0.034	0.023	0.021	0.021	0.196
2	1.121	0.707	0.538	0.431	0.498	0.396	0.375	0.362	0.630	0.426	0.374	0.361	2.650
3	0.005	0.005	0.005	0.005	0.005	0.005	0.006	0.006	0.006	0.005	0.005	0.006	0.003
4	0.626	0.419	0.325	0.284	0.267	0.246	0.211	0.244	0.356	0.254	0.221	0.236	1.857
5	1.251	0.837	0.649	0.566	0.532	0.490	0.426	0.485	0.709	0.506	0.440	0.469	4.844
6	0.002	0.002	0.003	0.003	0.002	0.003	0.003	0.003	0.004	0.003	0.003	0.004	0.002
7	0.769	0.535	0.476	0.434	0.355	0.298	0.327	0.315	0.448	0.332	0.322	0.332	4.152
8	1.535	1.070	0.950	0.867	0.708	0.596	0.652	0.640	0.895	0.663	0.642	0.663	8.361
9	1.535	1.070	0.951	0.867	0.708	0.580	0.652	0.640	0.895	0.663	0.642	0.663	8.383
10	0.001	0.001	0.001	0.001	0.001	0.001	0.002	0.002	0.002	0.001	0.002	0.002	0.001
11	1.040	0.839	0.821	0.818	0.455	0.499	0.468	0.534	0.660	0.544	0.524	0.589	11.619
12	2.080	1.694	1.636	1.635	0.909	0.997	0.954	1.086	1.319	1.086	1.047	1.176	23.237
13	2.080	1.678	1.641	1.635	0.909	0.997	0.954	1.086	1.319	1.086	1.047	1.176	23.237
14	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001
15	0.892	0.811	0.851	0.854	0.453	0.427	0.485	0.535	0.606	0.501	0.529	0.594	20.639
16	0.892	0.823	0.859	0.858	0.453	0.427	0.485	0.535	0.617	0.511	0.539	0.600	41.278
17	0.892	0.811	0.859	0.854	0.453	0.427	0.485	0.535	0.617	0.511	0.540	0.606	41.278
18	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.001	0.001	0.000	0.000	0.001	0.000
19	0.108	0.093	0.091	0.087	0.092	0.090	0.104	0.123	0.133	0.110	0.111	0.121	0.386
20	0.931	0.806	0.796	0.763	0.725	0.894	0.911	0.969	1.053	0.870	0.873	0.956	1.655
21	0.002	0.002	0.002	0.002	0.004	0.004	0.003	0.005	0.002	0.006	0.006	0.007	0.004
22	0.007	0.007	0.007	0.007	0.007	0.008	0.008	0.010	0.009	0.008	0.008	0.009	0.001

Table 5 shows the breakdown of the execution time obtained by all the explored T1 coprocessor configurations and by the non-accelerated T0 core, for each VGG-16 layer. The results give evidence to the fact that the performance of the coprocessor hardware configurations varies with the algorithm layer it executes. The Symmetrical MIMD configuration (D=1) results to be the best performing for layers 3, 6, 10-14, 18, 20 and 22, while the Symmetrical MIMD configuration (D=2) results to be the optimal choice for layers 7-9, 15-17. Notably, the Maxpool and Softmax layers exhibit worse execution time with respect to the non-accelerated core, because in the present software implementation, they are executed as scalar computation in the core, and so the data transfer to/from the SPMs constitutes an overhead with no corresponding computation acceleration. Nonetheless, the relative impact of those layers on the overall execution time is minor.

Figure 14 presents the total execution time speed-up obtained by each coprocessor configuration over the non-accelerated T0 core. The diagram also includes the speed-up obtained assuming to use the optimal configuration for each layer, giving evidence of performance advantage.

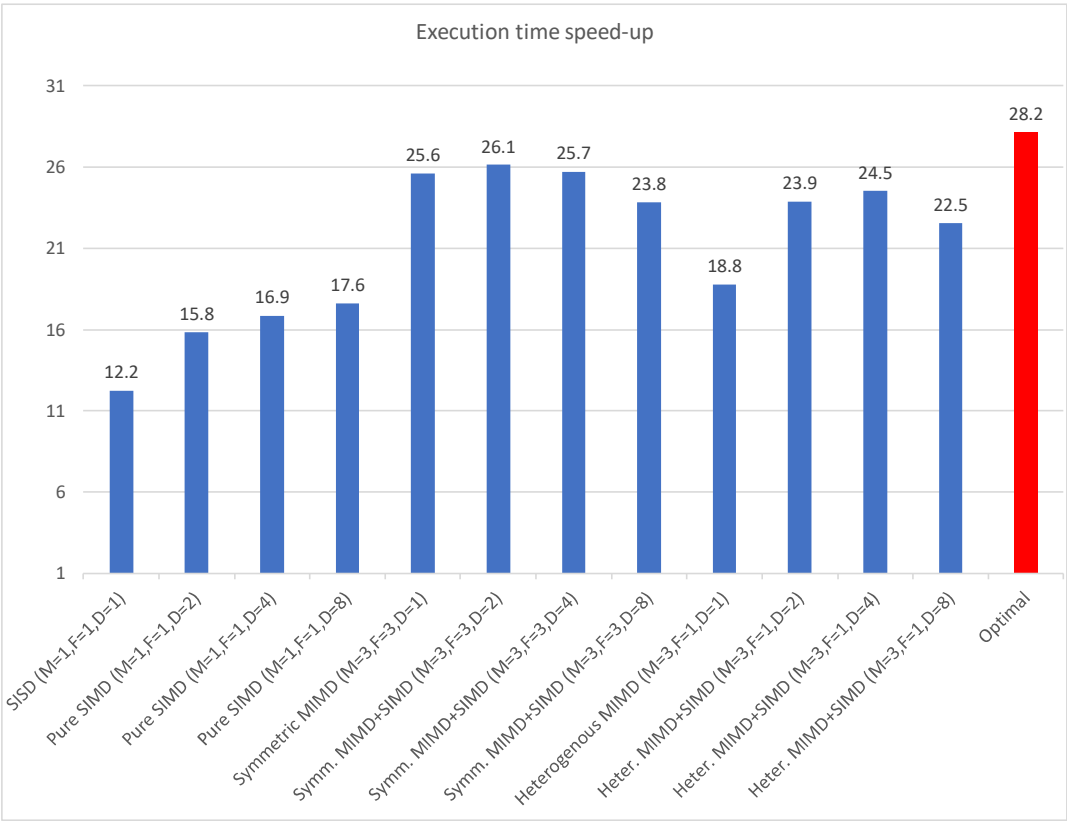


Figure 14. Total execution time speed-up over non-accelerated core obtained by each coprocessor configuration, along with the speed-up obtained by using the optimal configuration for each layer

Table 6 shows the breakdown of the total energy consumed by all the explored T1 coprocessor configurations and by the non-accelerated T0 core, for each VGG-16 layer. Again, it is evident that the optimal coprocessor configuration for energy efficiency is not unique for all the layers, yet it depends on the layer being executed. Optimal energy efficiency unlike absolute performance swings between Pure SIMD and Symmetrical MIMD configurations. Similarly to the execution time analysis, for pure scalar computation layers the energy consumption worsens in the vector-accelerated microarchitecture, due to the SPM data transfer overhead. Yet, the overall impact of those layers on the total energy is minor.

Table 6. Total energy consumption [J]. Best performing coprocessor configurations are highlighted for each layer

Layer	SISD (M=1, F=1, D=1)	Pure SIMD (M=1, F=1, D=2)	Pure SIMD (M=1, F=1, D=4)	Pure SIMD (M=1, F=1, D=8)	Symm MIMD (M=3, F=3, D=1)	Symm MIMD +SIMD (M=3, F=3, D=2)	Symm MIMD +SIMD (M=3, F=3, D=4)	Symm MIMD +SIMD (M=3, F=3, D=8)	Heter. MIMD (M=3, F=1, D=1)	Heter. MIMD +SIMD (M=3, F=1, D=2)	Heter. MIMD +SIMD (M=3, F=1, D=4)	Heter. MIMD +SIMD (M=3, F=1, D=8)	Non Accel T0 core
1	9.0E-03	5.5E-03	4.4E-03	4.7E-03	7.8E-03	5.2E-03	5.0E-03	6.6E-03	8.3E-03	5.4E-03	5.1E-03	6.8E-03	6.3E-02
2	1.8E-01	1.1E-01	8.3E-02	8.8E-02	1.4E-01	9.2E-02	8.7E-02	1.1E-01	1.5E-01	9.9E-02	9.0E-02	1.2E-01	8.6E-01
3	5.3E-04	5.6E-04	6.0E-04	7.3E-04	1.0E-03	1.0E-03	1.1E-03	1.6E-03	1.0E-03	1.0E-03	1.1E-03	1.6E-03	5.6E-04
4	8.9E-02	5.9E-02	4.7E-02	5.5E-02	6.8E-02	5.4E-02	4.6E-02	7.0E-02	7.4E-02	5.5E-02	5.0E-02	7.0E-02	5.4E-01
5	1.8E-01	1.2E-01	9.4E-02	1.1E-01	1.4E-01	1.1E-01	9.3E-02	1.4E-01	1.5E-01	1.1E-01	9.9E-02	1.4E-01	1.4E+00
6	2.7E-04	2.9E-04	3.1E-04	3.8E-04	5.2E-04	5.3E-04	5.7E-04	8.4E-04	6.4E-04	6.5E-04	6.9E-04	9.9E-04	2.9E-04
7	1.0E-01	7.6E-02	6.9E-02	8.8E-02	8.8E-02	6.6E-02	7.3E-02	9.3E-02	9.1E-02	7.4E-02	7.5E-02	1.0E-01	1.0E+00
8	2.0E-01	1.5E-01	1.4E-01	1.8E-01	1.8E-01	1.3E-01	1.5E-01	1.9E-01	1.8E-01	1.5E-01	1.5E-01	2.1E-01	2.1E+00
9	2.0E-01	1.5E-01	1.4E-01	1.8E-01	1.8E-01	1.3E-01	1.5E-01	1.9E-01	1.8E-01	1.5E-01	1.5E-01	2.1E-01	2.1E+00
10	1.5E-04	1.6E-04	1.7E-04	2.1E-04	2.9E-04	3.0E-04	3.2E-04	4.7E-04	2.9E-04	3.0E-04	3.2E-04	4.6E-04	1.6E-04
11	1.2E-01	1.0E-01	1.0E-01	1.3E-01	1.0E-01	1.0E-01	9.2E-02	1.4E-01	1.2E-01	1.1E-01	1.1E-01	1.5E-01	2.9E+00
12	2.5E-01	2.1E-01	2.0E-01	2.6E-01	2.0E-01	2.0E-01	1.9E-01	2.8E-01	2.3E-01	2.2E-01	2.1E-01	3.1E-01	5.9E+00
13	2.5E-01	2.1E-01	2.0E-01	2.6E-01	2.0E-01	2.0E-01	1.9E-01	2.8E-01	2.3E-01	2.2E-01	2.1E-01	3.1E-01	5.9E+00
14	9.8E-05	1.0E-04	1.1E-04	1.4E-04	1.9E-04	2.0E-04	2.0E-04	3.0E-04	1.9E-04	2.0E-04	2.1E-04	3.0E-04	1.0E-04
15	1.0E-01	9.7E-02	1.0E-01	1.3E-01	9.9E-02	8.4E-02	9.3E-02	1.3E-01	1.1E-01	9.7E-02	1.0E-01	1.5E-01	4.0E+00
16	1.0E-01	9.9E-02	1.0E-01	1.3E-01	9.9E-02	8.4E-02	9.3E-02	1.3E-01	1.1E-01	9.9E-02	1.1E-01	1.5E-01	8.1E+00
17	1.0E-01	9.7E-02	1.0E-01	1.3E-01	9.9E-02	8.4E-02	9.3E-02	1.3E-01	1.1E-01	9.9E-02	1.1E-01	1.5E-01	8.1E+00
18	4.9E-05	5.1E-05	5.5E-05	6.8E-05	9.4E-05	9.5E-05	1.0E-04	1.5E-04	9.2E-05	9.4E-05	1.0E-04	1.4E-04	5.2E-05
19	1.1E-02	1.0E-02	9.8E-03	1.3E-02	1.9E-02	1.7E-02	1.9E-02	3.1E-02	2.2E-02	2.1E-02	2.1E-02	3.1E-02	1.1E-01
20	9.5E-02	8.8E-02	8.6E-02	1.1E-01	1.5E-01	1.7E-01	1.7E-01	2.5E-01	1.7E-01	1.6E-01	1.7E-01	2.5E-01	4.9E-01
21	2.5E-04	2.4E-04	2.6E-04	3.4E-04	8.7E-04	8.1E-04	4.9E-04	1.2E-03	2.7E-04	1.1E-03	1.1E-03	1.7E-03	1.2E-03
22	1.8E-04	1.8E-04	1.8E-04	1.8E-04	2.1E-04	2.1E-04	2.0E-04	2.0E-04	2.0E-04	2.0E-04	2.0E-04	2.0E-04	4.1E-05

Figure 15 gives significance of the total energy saving obtained by each coprocessor configuration over the non-accelerated T0 core. The energy saving is expressed as the fraction of the energy consumed in the accelerated core over the energy consumed in the non-accelerated core, obtaining energy consumption between 6.4% and 4% of the non-accelerated core (energy saving between 93.6% and 96%). The diagram also includes the energy reduction obtained assuming to use the optimal configuration for each layer.

The outcome of Tables 5 and 6 is that dynamically changing the coprocessor microarchitecture, by updating the FPGA bitstream when a new algorithm layer is to be computed, allows an IoT device system to always use the optimal hardware scheme to achieve the desired goal of computation speed or power efficiency. Software controlled bitstream updating is available in several commercial FPGA devices.

The second phase of performance analysis aimed at comparing the efficiency of the proposed soft-processor architecture with the alternative hardware architecture solutions for the execution of the same application. In this analysis, the proposed solution consisted of the extended PULPino

platform equipped with the Klessydra T13 core + optimal vector coprocessor for each layer being executed.

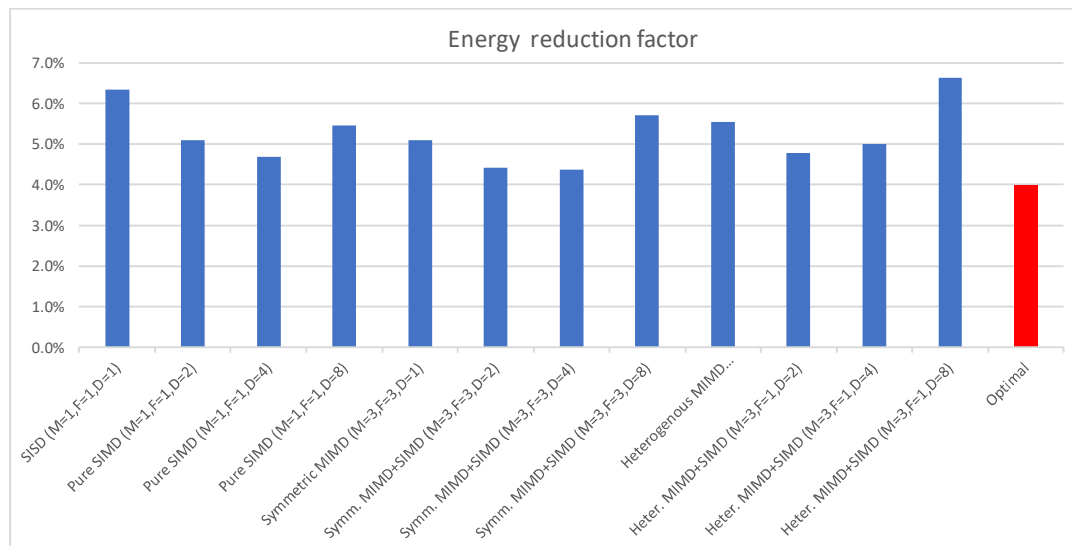


Figure 15. Energy reduction factor with respect to non-accelerated core (lower is better) obtained by each coprocessor configuration, along with the energy obtained by using the optimal configuration for each layer

Table 7 summarizes the performance comparison results, expressed as expressed as total execution time, total energy consumption, and average energy consumed per algorithmic operation. Algorithmic operations are the data multiplications and additions that are inherent to the algorithm being computed, and do not depend on the actual software implementation. The absolute execution time obviously favors high-end computing devices, yet the results give evidence of the effectiveness of the Klessydra T1 customizable vector coprocessor sub-system with respect to other single-core PULPino soft-processor FPGA implementations. Also, the energy efficiency results show the potential advantage of a Klessydra T1 vector-accelerated soft-processor FPGA implementation, with respect to general purpose single-board computers.

Table 7. Performance comparison with alternative solutions

Processor	Time [s]	Energy [J]	Energy per op [pJ/op]
Core i7 PC board	0.08	2.90	21
Cortex A53 Raspberry Pi 3	0.89	2.32	17
Cortex M4 STM32	117.78	7.77	55
RISCY PULPino on FPGA	315.91	40.06	285
Zeroriscy PULPino on FPGA	360.56	38.90	277
Klessydra-T1 PULPino on FPGA	6.88	1.74	12

6. Conclusion

The validation of the VGG-16 output data produced by Klessydra processors against VGG-16 inference demonstrate the suitability of the Klessydra open-source infrastructure for the implementation of FPGA based configurable RISC-V soft-cores equipped with hardware acceleration for vector computing. The detailed porting of the target application routines has been documented in this work.

Performance results show the effectiveness of the Klessydra microarchitecture scheme, built upon interleaved multi-threading and vector coprocessor hardware acceleration, with respect to other FPGA-based single-core solutions. Looking at energy efficiency, the Klessydra FPGA soft-core solution shows superior performance with respect to commercial single-board computers that may be used as IoT extreme-edge devices.

The results of the performance analysis conducted on the Klessydra T1 vector coprocessor schemes demonstrate the dependency of the optimal hardware configuration on the algorithm layer being executed. This evidence opens the way to the development of software configurable accelerators and further to the implementation of self-adapting coprocessor microarchitectures in IoT extreme-edge nodes.

Supplementary Materials: The Klessydra processor core family and accelerators are openly available online at <https://www.github.com/klessydra>

References

1. Cheikh, A., Cerutti, G., Mastrandrea, A., Menichelli, F. and Olivieri, M., 2017, September. The microarchitecture of a multi-threaded RISC-V compliant processing core family for IoT end-nodes. In International Conference on Applications in Electronics Pervading Industry, Environment and Society (pp. 89-97). Springer, Cham.
2. Olivieri, M., Cheikh, A., Cerutti, G., Mastrandrea, A., and Menichelli, F., "Investigation on the optimal pipeline organization in RISC-V multi-threaded soft processor cores", In 2017 New Generation of CAS (NGCAS), pp. 45-48. IEEE, 2017.
3. RISC-V Instruction Set specifications. [Online] "<https://riscv.org/specifications/>"
4. [P5] Cheikh, A., Sordillo, S., Mastrandrea, A., Menichelli, F. and Olivieri, M., 2019, September. Efficient Mathematical Accelerator Design Coupled with an Interleaved Multi-threading RISC-V Microprocessor. In International Conference on Applications in Electronics Pervading Industry, Environment and Society (pp. 529-539). Springer, Cham.
5. Samie, F.; Bauer, L.; Henkel, J. "From Cloud Down to Things: An Overview of Machine Learning in Internet of Things". IEEE Internet Things J. 2019, 4662, 1.
6. Gautschi, M., Schiavone, P., Traber, A., Loi, I., Pullini, A., Rossi, D., Flamand, E., Gürkaynak, F., Benini, L., "Near-threshold RISC-V core with DSP extensions for scalable IoT endpoint devices." IEEE Trans. on Very Large Scale Integration (VLSI) Systems 25, no. 10 (2017): 2700-2713.
7. Seo, S., Dreslinski, R.G., Woh, M., Chakrabarti, C., Mahlke, S. and Mudge, T., 2010, August. Diet SODA: A power-efficient processor for digital cameras. In Proceedings of the 16th ACM/IEEE international symposium on Low power electronics and design (pp. 79-84).
8. Chen, Y.H., Krishna, T., Emer, J.S. and Sze, V., 2016. Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks. IEEE journal of solid-state circuits, 52(1), pp.127-138.
9. Moini, S., Alizadeh, B., Emad, M. and Ebrahimpour, R., 2017. A resource-limited hardware accelerator for convolutional neural networks in embedded vision applications. IEEE Transactions on Circuits and Systems II: Express Briefs, 64(10), pp.1217-1221.
10. Du, L., Du, Y., Li, Y., Su, J., Kuan, Y.C., Liu, C.C. and Chang, M.C.F., 2017. A reconfigurable streaming deep convolutional neural network accelerator for Internet of Things. IEEE Transactions on Circuits and Systems I: Regular Papers, 65(1), pp.198-208.
11. Conti, Francesco, and Luca Benini. "A ultra-low-energy convolution engine for fast brain-inspired vision in multicore clusters." In 2015 Design, Automation & Test in Europe Conference & Exhibition (DATE), pp. 683-688. IEEE, 2015.
12. Meloni, P., Deriu, G., Conti, F., Loi, I., Raffo, L. and Benini, L., 2016, May. Curbing the roofline: a scalable and flexible architecture for CNNs on FPGA. In Proceedings of the ACM International Conference on Computing Frontiers (pp. 376-383).

13. Wu, N., Jiang, T., Zhang, L., Zhou, F. and Ge, F., 2020. A Reconfigurable Convolutional Neural Network-Accelerated Coprocessor Based on RISC-V Instruction Set. *Electronics*, 9(6), p.1005.
14. Watanabe, D., Yano, Y., Izumi, S., Kawaguchi, H., Takeuchi, K., Hiramoto, T., Iwai, S., Murakata, M. and Yoshimoto, M., 2020. An Architectural Study for Inference Coprocessor Core at the Edge in IoT Sensing. In 2020 2nd IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS) (pp. 305-309). IEEE.
15. Wu, Y., Wang, J.J., Qian, K., Liu, Y., Guo, R., Hu, S.G., Yu, Q., Chen, T.P., Liu, Y. and Rong, L., 2020. An energy-efficient deep convolutional neural networks coprocessor for multi-object detection. *Microelectronics Journal*, p.104737.
16. Chang, M.C., Pan, Z.G. and Chen, J.L., 2017, October. Hardware accelerator for boosting convolution computation in image classification applications. In 2017 IEEE 6th Global Conference on Consumer Electronics (GCCE) (pp. 1-2). IEEE.
17. Lima, P., Vieira, C., Reis, J., Almeida, A., Silveira, J., Goerl, R. and Marcon, C., 2020, March. Optimizing RISC-V ISA Usage by Sharing Coprocessors on MPSoC. In 2020 IEEE Latin-American Test Symposium (LATS) (pp. 1-5). IEEE.
18. Du, L., Du, Y., Li, Y., Su, J., Kuan, Y.C., Liu, C.C. and Chang, M.C.F., 2017. A reconfigurable streaming deep convolutional neural network accelerator for Internet of Things. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 65(1), pp.198-208.
19. Schiavone P.D., Conti F., Rossi D., Gautschi M., Pullini A., Flamand E., Benini L., Slow and steady wins the race? a comparison of ultra-low-power risc-v cores for internet-of-things applications. In 2017 27th International Symposium on Power and Timing Modeling, Optimization and Simulation (PATMOS) 2017 Sep 25 (pp. 1-8). IEEE.
20. Traber A, Gautschi M., PULPino: datasheet. ETH Zurich, University of Bologna. 2017 Jun 9.
21. Blasi L, Vigli F, Cheikh A, Mastrandrea A, Menichelli F, Olivieri M. A RISC-V Fault-Tolerant Microcontroller Core Architecture Based on a Hardware Thread Full/Partial Protection and a Thread-Controlled Watch-Dog Timer. In International Conference on Applications in Electronics Pervading Industry, Environment and Society 2019 Sep 11 (pp. 505-511). Springer, Cham.
22. European Processor Initiative (EPI), EU H2020 research and innovation programme GA No 826647, [Online] "<https://www.european-processor-initiative.eu/project/epi/>".
23. A. Cheikh, S. Sordillo, A. Mastrandrea, F. Menichelli, G. Scotti and M. Olivieri, "Klessydra-T: Designing Vector Coprocessors for Multi-Threaded Edge-Computing Cores," in *IEEE Micro*, doi: 10.1109/MM.2021.3050962.