# MARVEL: An End-to-End Framework for Generating Model-Class Aware Custom RISC-V Extensions for Lightweight AI

**Ajay Kumar M[1], Cian O'Mahoney[1], Pedro Kreutz Werle[1], Shreejith Shanker[2], Member, IEEE, Dimitrios S. Nikolopoulos[3], Fellow, IEEE, Bo Ji[3], Senior Member, IEEE, Hans Vandierendonck[4], Senior Member, IEEE, and Deepu John[1], Senior Member, IEEE**

[1]University College Dublin, Dublin 4, Ireland
[2]Trinity College Dublin, Dublin 4, Ireland
[3]Virginia Tech, Blacksburg, VA 24061, USA
[4]Queen's University Belfast, Belfast, UK

Corresponding author: Ajay Kumar M (email: ajay.kumarm@ucdconnect.ie).

**ABSTRACT** Deploying deep neural networks (DNNs) on resource-constrained IoT devices remains a challenging problem, often requiring hardware modifications tailored to individual AI models. Existing accelerator-generation tools, such as AMD's FINN, do not adequately address extreme resource limitations faced by IoT endpoints operating in bare-metal environments without an operating system (OS). To overcome these constraints, we propose MARVEL—an automated, end-to-end framework that generates custom RISC-V ISA extensions tailored to specific DNN model classes, with a primary focus on convolutional neural networks (CNNs). The proposed method profiles high-level DNN representations in Python and generates an ISA-extended RISC-V core with associated compiler tools for efficient deployment. The flow leverages (1) Apache TVM for translating high-level Python-based DNN models into optimized C code, (2) Synopsys ASIP Designer for identifying compute-intensive kernels, modeling, and generating a custom RISC-V and (3) Xilinx Vivado for FPGA implementation. Beyond a model-class specific RISC-V, our approach produces an optimized bare-metal C implementation, eliminating the need for an OS or extensive software dependencies. Unlike conventional deployment pipelines relying on TensorFlow/PyTorch runtimes, our solution enables seamless execution in highly resource-constrained environments. We evaluated the flow on popular DNN models such as LeNet-5*, MobileNetV1, ResNet50, VGG16, MobileNetV2 and DenseNet121 using the Synopsys trv32p3 RISC-V core as a baseline. Results show a $2\times$ speedup in inference and upto $2\times$ reduction in energy per inference at a 28.23% area overhead when implemented on an AMD Zynq UltraScale+ ZCU104 FPGA platform.

**INDEX TERMS** RISC-V, ISA extensions, TVM, Deep Neural Networks, FPGA, Hardware accelerators

## I. INTRODUCTION

The rapid proliferation of IoT devices has transformed the way we interact with physical systems, enabling efficient remote automation and monitoring. A key advancement in this space is Edge AI, which brings AI computation closer to the data source—typically at the network edge—to enhance responsiveness and reduce latency [1]. However, Edge AI solutions are inherently resource-constrained, driving research efforts toward optimizing hardware efficiency. While hardware acceleration has received significant attention, the methodology behind designing ISA extensions remains an equally important yet often underexplored aspect. Many existing works introduce ISA extensions based on assumed computational hotspots in DNNs such as matrix multiplications or rely on heuristic, domain-specific motivations without empirical profiling. For example, Xvpfloat [2] introduces a RISC-V ISA extension designed for dynamically variable and extended precision floating-point computations. The extension supports significant sizes up to 512 bits and includes features like specialized indexed loads/stores and hardware-assisted prefetching. However, this design is driven by the need for efficient hardware support in scientific computing, not by application profiling. Similarly, Convex [3] proposes a RISC-V ISA extension aimed at enhancing the performance

of convolution operations on microcontrollers. The extension includes bit concatenation, weight storage optimization, and activation value reuse, yet these are guided by high-level architectural goals rather than profiling data. Authors in [4] describe Confidential Virtual Machine Extension (CoVE) for RISC-V, aiming to provide hardware support for confidential computing. The proposed ISA and non-ISA extensions focus on security features necessary for trusted execution environments. Here too, the design decisions are driven by security requirements rather than application profiling.

Another often overlooked consideration is the ease of programming and software integration, which plays a crucial role in the practical adoption of these extensions. Despite advances in AI hardware, the lack of automated and user-friendly software solutions remains a significant barrier to widespread adoption. Existing solutions either focus on translating high-level domain-specific languages (DSLs) like PyTorch and TensorFlow into low-level code, or offloading the programming responsibility entirely to the user while focusing solely on hardware efficiency. However, very few solutions provide a complete, end-to-end framework that seamlessly integrates AI models with specialized hardware. For instance, BARVINN [5], a RISC-V-based DNN accelerator, processes ONNX inputs and generates executable command streams for a RISC-V controller. However, its software framework lacks flexibility, as it does not support residual connections, limiting compatibility with models like ResNet. RISC-VTF [6] extends the RISC-V instruction set to improve transformer model execution through custom matrix operations and activation functions, yet the absence of compiler support forces users to write low-level RISC-V assembly. Similarly, Eyeriss [7], designed as a co-processor to accelerate CNN workloads, requires a robust software stack to efficiently offload compute-intensive tasks. Other solutions, such as the LSTM accelerator by Kadetotad et al. [8], optimize speech recognition workloads but require substantial modifications in quantization, compression, and memory structure, making them unsuitable for off-the-shelf pretrained models. Another RISC-V-based processor [9] for IoT endpoints introduces DSP extensions for energy-efficient filtering, convolutions, and simple ML operations. While the compiler toolchain supports vectorization and hardware loops, it lacks an automated model conversion pipeline from high-level AI frameworks to embedded C, requiring manual coding efforts.

The challenge of limited software accessibility is also evident in model-specific instruction-based accelerators. XpulpNN [10] enhances a RISC-V core with SIMD operations for quantized neural networks (QNN), achieving significant speedups. However, details on its programming interface remain sparse. ZeroRiscy [11] accelerates CNN workloads but requires developers to write RISC-V assembly code to leverage custom ISA extensions. In another example, Aness et al. [12] optimize transformer-based models for constrained devices by porting ARM Keyword Spotting Transformer [13], reducing model size by 369× with only a 10% accuracy drop, yet this relies on manual C-code modifications without an automated pipeline.

A broader challenge arises when designing general-purpose instruction-based accelerators, where software-hardware co-design is crucial for usability. FlexACC [14] addresses model- accelerator incompatibilities by introducing ISA extensions for MLP, CNN, LSTM, GCN, and transformer models, achieving a 216× CNN speedup. However, its programming model remains low-level and requires C-like coding. Similarly, [15] explores an ISA extension that optimizes memory access patterns in AI applications, reducing power consumption and code size without significant area overhead. However, compiler updates are necessary for full utilization, and support for non-consecutive memory accesses remains limited. Efforts have also been made to distribute computation between IoT devices and edge servers to reduce latency and enhance privacy. For example, [16] presents a framework that integrates pruning and early exit strategies with FPGA-based offloading, achieving a 1.6× latency reduction and 3.9× power efficiency improvement. However, its complexity makes adaptation to different AI models challenging. Compute-in-memory architectures [17] introduce dedicated RISC-V ISA extensions for switching operational modes but lack an end-to-end flow that supports high-level AI model inputs as it expects the software program to be loaded to bootRAM from external flash with weight parameters already stored into multi-bank eMRAM. AI-RISC [18] represents a more integrated approach, co-designing hardware, ISA, and software to extend RISC-V for Edge AI acceleration. By embedding AI functional units into the RISC-V pipeline, it enables seamless execution of both AI and non-AI tasks. However, AI-RISC requires modifications to the TVM compiler to recognize its intrinsics, making adaptation to different models cumbersome. In contrast, our work eliminates the need for such modifications, providing a streamlined, adaptable, and automated solution for bridging high-level AI frameworks with edge hardware. By addressing the gap in software accessibility, we introduce a framework that seamlessly integrates AI models with resource-constrained accelerators, facilitating efficient and user-friendly deployment of Edge AI solutions.

This paper presents an automated workflow that converts any Python-based CNN model into a bare metal implementation for a RISC-V CPU enhanced with the proposed ISA extensions, optimizing both speed and energy efficiency. For evaluation, the design was implemented on a Xilinx ZCU104 FPGA and the performance is benchmarked on C implementations of pretrained DNNs. The latency required to perform model inference is compared between the baseline RISC-V core and the same core with added custom ISA extensions. The main contributions of this paper (and its comparison with some related works in Table 11) can be summarized as follows:

- An end-to-end flow showcasing how Python-based DNN models from high-level frameworks like PyTorch/TensorFlow are translated to C and profiled to develop a model class-specific custom RISC-V which enables bare metal programming.
- A RISC-V core with a custom ISA generated with the proposed MARVEL flow for lightweight CNNs which achieves an inference acceleration of up to **2**×.
- A tool flow based on ASIP Designer [19], to facilitate an accelerated modeling and development of the identified ISA extensions. The flow described in this paper tracks the design of these extensions from Python models to final hardware implementation.
- Benchmarking of DNN models to quantify the performance of the custom RISC-V core against the baseline. FPGA demonstration of the identified extensions with a reduction in energy per inference by up to 2× for models such as LeNet-5*[1], MobileNetV1, ResNet50, VGG16, MobileNetV2 and DenseNet121.

## II. Methodology

The trv32p3 processor core [20] from Synopsys is used as a starting point for the development of an Application Specific Instruction set Processor (ASIP). trv32p3 has a 32-bit wide data path and a three-stage pipeline. This core implements the RV32IM ISA that includes integer instructions and hardware support for integer multiplication, division and remainder. Additionally, on-chip debugging (OCD) for JTAG access and memory was integrated before compiling with ASIP designer to generate HDL files for FPGA implementation. The design was synthesized targeting an AMD ZCU104 development platform. The extended core was tested with quantised standard DNN models' C codes obtained through the proposed workflow as described further. The complete process, from converting a Python-based model to executing it on the extended RISC-V core implemented on an FPGA, is illustrated in Fig 1 and 2

### A. Workflow to generate C code

This section outlines the methodology for converting an AI model, implemented in Python using high-level deep learning frameworks such as TensorFlow, into C code suitable for compilation with Synopsys' ASIP Designer tools. The target architecture in this work is the RISC-V-based trv32p3 processor and its ISA extended variants. To enable this conversion, we leverage the TVM machine learning compiler [21], an open source framework designed to optimize and deploy deep learning models on a variety of hardware backends. Notably, there are other commercially available frameworks that also enable this translation [22], [23]. The translation process begins with the conversion of the high-level Python-based model into Relay, an intermediate representation (IR) that is hardware-agnostic and optimized for

computational graphs. Relay serves as an abstraction layer, enabling transformations and optimizations independent of the underlying hardware target. Once the model is represented in Relay IR, TVM applies a series of optimization passes to enhance computational efficiency before generating source code tailored for the specified backend. In this work, we configure TVM to target generic C, producing C-based computational kernels that can be compiled and executed on any processor that supports C compilation. This approach facilitates portability and enables seamless integration of AI workloads on RISC-V based embedded systems. The following subsections detail the step-by-step procedure for generating C code using TVM, from model parsing to final code generation. Additionally, an open-source repository has been made available on GitHub[2], containing example scripts and reference implementations for converting DNN models into C code suitable for RISC-V deployment.

### 1) Environment Setup
To begin, all required dependencies and software packages must be installed. The models used in this work are pre-trained keras deep learning architectures, which are subsequently fine-tuned for the target classification task. Five deep learning models—MobileNetV1, ResNet50, VGG16, MobileNetV2, and DenseNet121—are fine-tuned using the StanfordCars [24] and COCO [25] datasets. The input images are standardized to dimensions of 64×64×3 (height, width, channels), and the classification task is simplified to a binary problem: distinguishing between "Car" and "Not Car." In addition, a modified LeNet-5 model is employed to classify 28×28 grayscale images of handwritten digits into one of the ten-digit classes.

### 2) Model Training Using Transfer Learning
This step is optional and not part of the main workflow, though it is recommended. The proposed approach is designed for IoT-embedded systems with limited resources. For these systems, minimizing network size is crucial for compatibility. Therefore, in this step, we fine-tune the standard models to focus on two target classes. To adapt the pretrained models for the target application, transfer learning is employed. The models are re-trained for a limited number of epochs until the validation accuracy reaches an acceptable level. The dataset is split into an 80:20 ratio for training and validation. To accommodate the binary classification task, the final fully connected layer of the original network is removed and replaced with a dense layer containing two neurons, corresponding to the two output classes.

### 3) Model Size Optimization via Quantization
This step is optional but included in the proposed flow to further reduce the network size. Once the model is trained, it exists as a TensorFlow-based DNN tailored to the car recognition task. However, its memory footprint

---

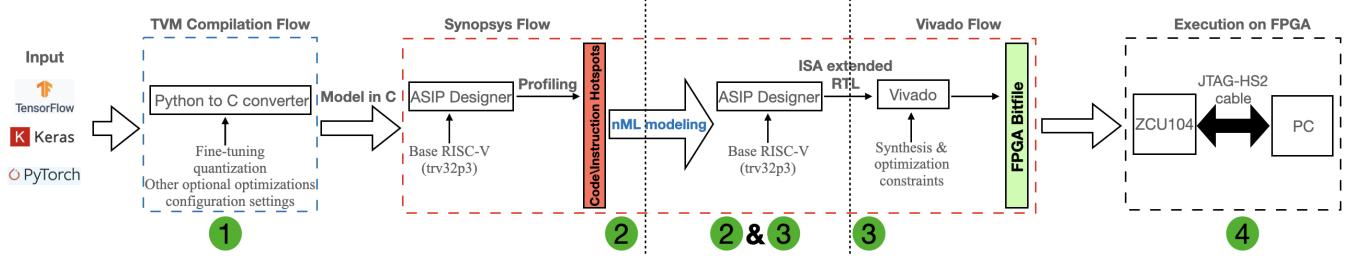[1]Hand-coded C implementation of LeNet-5-like model

FIGURE 1: MARVEL: An end-to-end framework to generate model-class specific custom RISC-V for DNN acceleration
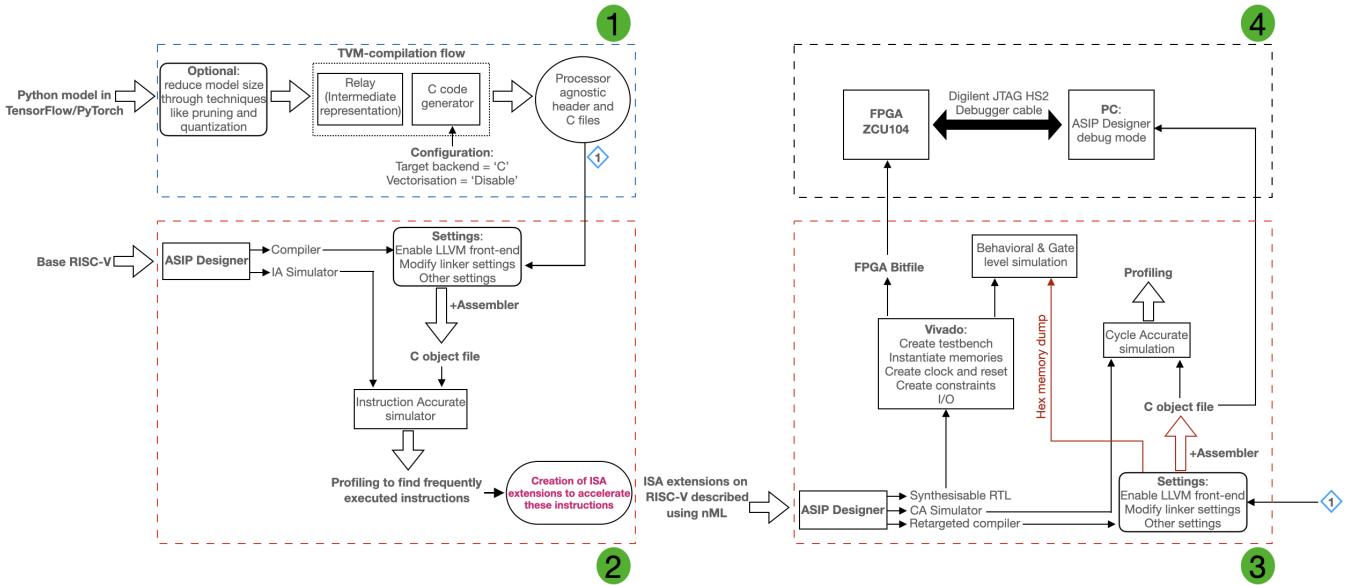


FIGURE 2: MARVEL: Detailed flow diagram

remains prohibitively large for deployment on embedded platforms with constrained memory resources. For instance, MobileNetV1, one of the models used in this work, contains 216,000 parameters, each represented as a 32-bit floating-point number in the native model. This translates to very high storage and computational requirements, making them non-viable for deployment on edge platforms. To mitigate this, post-training quantization is applied to reduce model size and improve inference efficiency. The TensorFlow Lite framework [26] is leveraged to convert the model's floating-point parameters into 8-bit integers, significantly reducing both memory consumption and computational overhead without substantial accuracy degradation [27].

### 4) Conversion of Python Model to Bare-Metal C

Up to this stage, TensorFlow has been used to construct and train the DNN model. For deployment on the trv32p3 RISC-V processor, the model is first converted into C code [3] for compilation with ASIP Designer. This conversion is performed using TVM compiler as shown in Fig. 2. TVM

---

[3]This does not include the hand-coded LeNet-5* implementation, as it is already written in C.

translates high-level Python-based AI models into optimized C code, tailored for the specified target hardware. Since the RISC-V processor used in this work operates without an operating system (bare metal environment), we configured TVM to generate generic C code.

### 5) TVM-generated C Application

Once the model is successfully compiled with TVM, the generated C source files are stored in a structured directory. The key output files are:

- Header file defining structs for inputs and outputs → ./codegen/host/include/tvmgen_default.h
- C source files containing the compiled neural network weights and implementation → ./codegen/host/src/default_lib0.c , ./codegen/host/src/default_lib1.c

In addition to the compiled source files, sample input images are provided for validation purposes. No modifications are required for the generated C source files. To integrate the model into an application, a custom main function must be implemented to invoke the necessary TVM-generated functions for inference.

### B. ASIP Designer

The trv32p3 ISA is extended using ASIP Designer (version W-2024.12). ASIP Designer [19] is a retargetable compiler, simulation, and hardware generation environment which uses the special-purpose modeling languages nML and PDG (Primitives Definition and Generation). nML is a high-level language for modeling processor architectures. An nML description of a target processor can be converted into a hardware description with the use of the Synopsys Go® nml-to-hdl compiler [28]. The primitives used by an nML processor description are described in the separate PDG language. It is with this language that the specific processor functionality, I/O behaviour and control behaviour are described. The custom instructions described in the next section are added to the existing nML model for the trv32p3 processor, with the specific functionality defined in a separate PDG description.

### C. ISA extensions

Traditionally, ISA extensions are introduced to address known computational bottlenecks, often without a comprehensive analytical foundation. For example, LiteAIR5 [29] proposes ISA extensions targeting matrix-matrix multiplication, a well-known bottleneck in AI workloads. In this study, we took a more data-driven approach by profiling (step 2 in Fig. 1) C code generated through the proposed flow for MobileNetV1 on a baseline RISC-V architecture to systematically identify performance bottlenecks. Using Synopsys tools, we analyze execution patterns to pinpoint compute-intensive operations that could benefit from ISA extensions. Specifically, we enable instruction profiling in a instruction-accurate simulator, capturing execution counts for each instruction. We then sort and analyze the most cycle-intensive instructions, identifying key compute bottlenecks that can be addressed through ISA extensions. Additionally, it would be valuable to conduct a similar analysis using RISC-V's open-source toolchain and Spike simulator, as demonstrated in [30], to obtain detailed instruction coverage metrics for specific RISC-V extensions. Our profiling revealed a set of frequently executed instruction patterns that, when optimized, could significantly enhance performance. To validate the applicability of our proposed extensions, we extended our profiling to other models within the same class. The results confirmed that the identified patterns were not model-specific but rather class-specific, making them broadly applicable. Fig 3 shows the normalised count of the number of times these patterns are executed in the corresponding model, while Table 2 provides definitions for the legends used in this figure. These patterns were then optimized through the proposed ISA extensions utilizing the custom opcodes available in RISC-V as outlined in Table 3. The following subsections provide a detailed discussion of these extensions, with Table 1 defining the terminology used to describe different processor versions as they evolve with custom instructions.

TABLE 1: Terminology for different RISC-V variants

| Processor Version | Description |
|---|---|
| v0 | Baseline RISC-V processor (trv32p3) |
| v1 | `mac` extension enabled on v0 |
| v2 | `add2i` extension enabled on v1 |
| v3 | `fusedmac` extension enabled on v2 |
| v4 | Zero-overhead hardware loops (`zol`) extension enabled on v3 |

TABLE 2: Description of instruction patterns in Fig. 3

| Pattern | Metric | Description |
|---|---|---|
| `mac` | add_count | Number of times the `add` instruction occurs |
| | mul_count | Number of times the `mul` instruction occurs |
| | mul_add_count | Number of times the `mul` and `add` instructions occur consecutively |
| `add2i` | addi_count | Number of times the `addi` instruction occurs |
| | addi_addi_count | Number of times two `addi` instructions occur consecutively |
| `fusedmac` | add_count | Number of times the `add` instruction occurs |
| | addi_count | Number of times the `addi` instruction occurs |
| | mul_count | Number of times the `mul` instruction occurs |
| | fusedmac_count | Number of times the `mul`, `add`, and two `addi` instructions occur consecutively |

TABLE 3: Instruction opcode encoding

| inst[4:2] / inst[6:5] | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| 00 | LOAD | unused | custom-0 = fusedmac | unused | ALU | AUIPC | unused | unused |
| 01 | STORE | unused | custom-1 = add2i | unused | ALU/NOP/ MPY/DIV | LUI | unused | unused |
| 10 | unused | unused | unused | unused | unused | reserved | custom-2 = mac | zol(2/2) |
| 11 | BRANCH | JALR | reserved | JAL | SYSTEM | reserved = zol(1/2) | SWBRK | unused |

### 1) mac

The instruction "`mul rd, rs1, rs2`" performs a 32-bit multiplication, followed by the instruction "`add rd, rd, rs1`" which accumulates the result through addition. This pair of instructions can be replaced by a single custom "`mac rd, rs1, rs2`". This operation performs the same operation in half the number of clock cycles. This is the most common operation in all AI/ML workloads [31]. In this work, we fix the registers ("`rd = x20, rs1 = x21, rs2 = x22`"), which enables us to combine mac with another custom instruction in subsequent steps. Hardcoding the source/sink registers also reduces the hardware area overhead incurred by the processor version, where both these custom instructions co-exist. The trade-off is reduced flexibility in register allocation, though this had minimal impact in practice due to ample general-purpose registers. This extension uses the CUSTOM 2 opcode. Listing 1 shows the assembly syntax for this instruction, and Table 4 shows its decoding scheme.
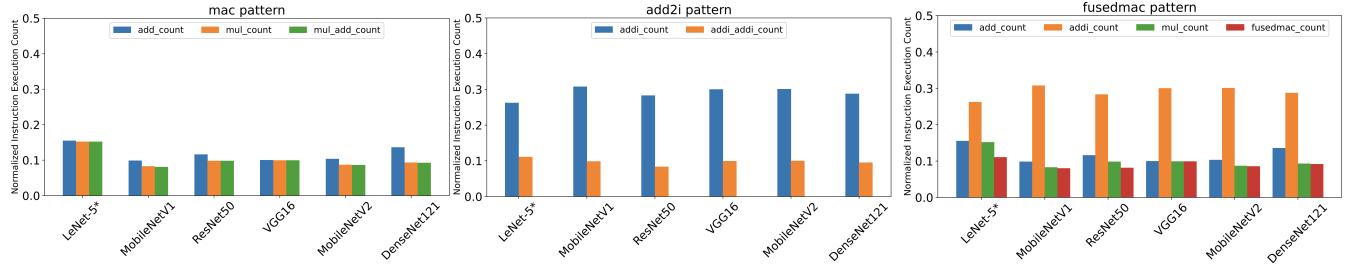
FIGURE 3: Frequently executed patterns identified through profiling of generated C code on baseline RISC-V

Listing 1: Assembly format of the custom mac instruction

```
1    ;rd = rd + rs1*rs2
2    mac
3    ;rd = x20, rs1 = x21, rs2 = x22
```

TABLE 4: Instruction decoding: mac

| | 31 25 | 24 20 | 19 15 | 14 12 | 11 7 | 6 0 | |
|---|---|---|---|---|---|---|---|
| mac | 0100000 | 00000 | 00000 | 000 | 00000 | 1011011 | R-type |
| | funct7 | rs2 = x22 | rs1 = x21 | funct3 | rd = x20 | opcode | |

## 2) add2i

This instruction merges the execution of two consecutive "add immediate" operations targeting different registers. However, incorporating two destination registers into the instruction syntax leaves only 15 bits available for encoding the immediate values. A direct split, allocating 7 bits per immediate (assuming signed values), would result in a limited range of (-64, 63), reducing potential use cases. To determine an optimal bit allocation, we analyzed the usage patterns of the addi instruction within the inner convolution loops of the DNN models where this fused instruction would be most beneficial. As shown in Fig. 4, all observed patterns mostly utilize only unsigned immediate values. Moreover, the most frequent patterns involve a small immediate followed by a larger one. Based on this observation, we allocated 5 bits to one immediate (i1: 0–31) and 10 bits to the other (i2: 0–1023), covering 100% (for LeNet-5*), 86.03% (for MobileNetV1), 75.19% (for ResNet50), 66.89% (for VGG16), 71.39% (for MobileNetV2) and 95.13% (for DenseNet121) of the relevant use cases when weighted by cycle count. This extension uses the CUSTOM 1 opcode. Listing 2 shows the assembly syntax for this instruction, and Table 5 shows its decoding scheme.

Listing 2: Assembly format of custom add2i instruction

```
1    ;rs1 = rs1 + i1
2    ;rs2 = rs2 + i2
3    add2i rs1, rs2, i1, i2
```

TABLE 5: Instruction decoding: add2i

| | 31 20 | 19 15 | 14 12 | 11 7 | 6 0 | |
|---|---|---|---|---|---|---|
| add2i | i2[9:0] :: i1[4:3] | rs2 | i1[2:0] | rs1 | 0101011 | I-type |
| | | | funct3 | | opcode | |

## 3) fusedmac

This custom instruction integrates the functionality of the previous two instructions. While it may seem having this instruction is redundant, it is necessary to evaluate whether using this combined instruction in all instances where either of the individual instructions could be applied is advantageous. Specifically, this instruction utilizes two functional units (mac and add2i), leading to increased power usage. Therefore, in scenarios where only the add2i operation is needed, executing an unnecessary mac operation on fixed registers will result in wasted power. Furthermore, an analysis of the baseline code reveals multiple opportunities to utilize all these custom extensions independently. As shown in the third subfigure of Fig. 3, the fusedmac pattern - comprising of four instructions in the sequence addi, addi, mul, and add - constitutes nearly 10% of the total code. This substantial proportion suggests that this pattern recurs frequently enough to justify its own custom instruction and dedicated hardware acceleration. Moreover, the occurrence of this pattern increases with larger neural network models, where having a dedicated ISA extension like the fusedmac is justified. Therefore, despite apparent redundancy, including all custom instructions is justified to optimize power efficiency and flexibility. This extension uses the CUSTOM 0 opcode. Listing 3 shows the assembly syntax for this instruction and Table 6 shows its decoding scheme.

Listing 3: Assembly format of custom fusedmac instruction

```
1    ;x20 = x20 + x21*x22
2    ;rs1 = rs1 + i1
3    ;rs2 = rs2 + i2
4    fusedmac rs1, rs2, i1, i2
```

TABLE 6: Instruction decoding: fusedmac

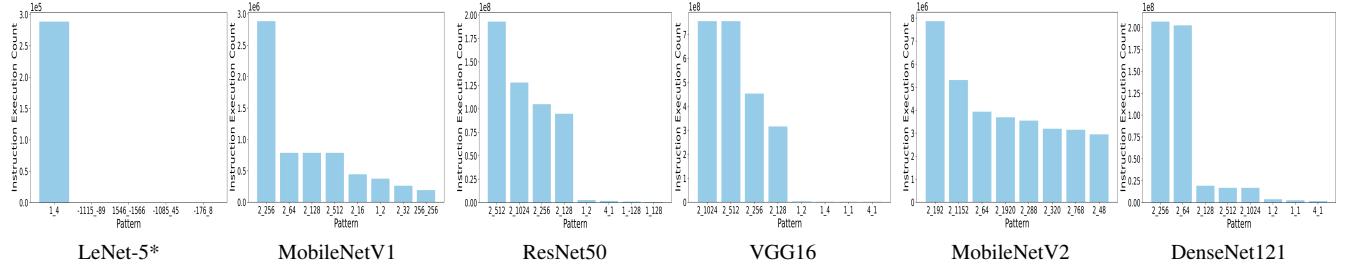| | 31 20 | 19 15 | 14 12 | 11 7 | 6 0 | |
|---|---|---|---|---|---|---|
| fusedmac | i2[9:0] :: i1[4:3] | rs2 | i1[2:0] | rs1 | 0001011 | I-type |
| | | | funct3 | | opcode | |

FIGURE 4: Instruction execution count for different consecutive `addi` immediate values. The x-axis represents the pattern `X_Y`, where `X` and `Y` denote the first and second immediate values in two consecutive `addi` instructions, respectively.

#### 4) zol

Profiling of C codes generated using the proposed flow on baseline RISC-V shows that the number of times `blt` instruction is executed is 923.2K in LeNet-5*, 23.6M in MobileNetV1, 1.89B in ResNet50, 6.88B in VGG16, 153.76M in MobileNetV2, and 1.43B in DenseNet121, highlighting the need for its optimization. Because of the way TVM generates code, lengths of convolutional `for` loops are known at compile time. Consequently, these loops can be fully unrolled, enabling the use of ISA extensions related to hardware loops to reduce cycle count. A hardware loop instruction essentially indicates in advance that the section of code following it is part of a loop, specifying the start and end addresses along with the iteration count. By leveraging hardware loops, the `blt` instruction is entirely eliminated from loop execution, along with the explicit increment of the loop counter, which is instead managed by the hardware itself. In this processor version implementation, the instructions `dlp`, `dlpi`, `zlp`, `set.zc`, `set.zs` and `set.ze` are introduced, based on Synopsys' reference processor design. This requires modifications to the Program Control Unit (PCU) and the addition of three new registers: `ZC` (Zero-overhead loop count), `ZS` (Start address), and `ZE` (End address). These registers are set up in the execution stage of the CPU pipeline. A total of five instructions are implemented to enable zero-overhead looping, with their decoding illustrated in Table 7. The hardware loop extensions utilize two opcodes: 11101, reserved for hardware loops, and 10111, which remains available if the implementation does not require instructions exceeding 32 bits.

Fig. 5 presents a comparison of the assembly code generated for the highlighted C code, comparing processor versions v0 and v4. The highlighted columns in the last two figures indicate the total number of clock cycles spent executing the corresponding instructions in the test case. Notably, the removal of the `blt` instruction, enabled by the hardware loop extension, leads to a significant reduction in execution cycles, contributing to improved efficiency.

#### D. Compilation

To ensure the successful compilation of the generated C code using ASIP Designer, the project file must be configured appropriately. This requires an understanding of the

TABLE 7: Instruction decoding: zol



Synopsys Chess® compiler [32], which consists of two key components:

1) **Front-End** – Responsible for parsing and interpreting the intent of the C source code.
2) **Back-End** – Responsible for translating the parsed code into the corresponding assembly instructions for the target architecture.

While the back-end of the Chess compiler is fixed, users can choose between two available front-ends:

- LLVM Front-End
- Chess Front-End

Additionally, ASIP Designer allows different front-ends to be used for different source files within the same project. The selection of the appropriate front-end depends on specific characteristics of the generated C code, as outlined below:

- Handling the `__attribute__` Construct: The generated C code makes use of the `__attribute__` construct, which is a standard feature in C but is not supported by the Chess front-end. However, it is fully supported by the LLVM front-end. Therefore, any source file that contains the `__attribute__` keyword must be compiled using LLVM.
- Optimizing Custom Instruction Utilization: A key drawback of the LLVM front-end is its inability to automatically identify and replace groups of assembly instructions with custom processor instructions. In contrast, the Chess front-end can efficiently infer and apply custom

```
for (int32_t ax0_ax1_fused_ax2_fused = 0; ax0_ax1_fused_ax2_fused < 64; ++ax0_ax1_fused_ax2_fused) {
  void* conv2d_nhwc_let = (&(global_workspace_29_var[24832]));
  for (int32_t ff = 0; ff < 64; ++ff) {
    ((int32_t*)conv2d_nhwc_let)[ff] = 0;
    for (int32_t rc = 0; rc < 64; ++rc) {
      ((int32_t*)conv2d_nhwc_let)[ff] = ((int32_t*)conv2d_nhwc_let)[ff] + (((int32_t)((int16_t*)pad_temp_let) \
      [((ax0_ax1_fused_ax2_fused * 64) + rc)]) * ((int32_t)((int16_t*)fused_constant_10_let)[((rc * 64) + ff)]));
    }
  }
  for (int32_t ax3_inner = 0; ax3_inner < 64; ++ax3_inner) {
    int32_t v_v = ((int32_t)((((((int64_t)((int32_t)conv2d_nhwc_let)[ax3_inner] \
    + ((int32_t)((int32_t*)fused_nn_conv2d_constant_10_let)[ax3_inner]) \
    * ((int64_t)((int32_t*)fused_nn_conv2d_add_constant_38_let)[ax3_inner])) \
    + ((int64_t)1) << (((int64_t)((int32_t*)fused_nn_conv2d_add_constant_32_let)[ax3_inner] + 31) - 1)))) \
    >> ((int64_t)((int32_t*)fused_nn_conv2d_add_constant_32_let)[ax3_inner] + 31)))] - 128;
    int32_t v_1 = (v) < (127) ? (v) : (127);
    T_cast[((ax0_ax1_fused_ax2_fused * 64) + ax3_inner)] = ((int8_t)((v_1) > (-128) ? (v_1) : (-128)));
  }
}
return 0;
}
```

```
  8192   00 00 08 b7   lui    x17,   0
524288   00 18 88 93   addi   x17,   x17,  1
524288   00 05 19 03   lh     x18,   0(x10)
524288   00 06 1a 03   lh     x20,   0(x12)
524288   03 2a 0a b3   mul    x21,   x20,  x18
524288   00 06 a9 83   lw     x19,   0(x13)
524288   01 3a 8b 33   add    x22,   x21,  x19
524288   01 66 a0 23   sw     x22,   0(x13)
524288   00 25 05 13   addi   x10,   x10,  2
524288   08 06 06 13   addi   x12,   x12,  128
1556480  fc 48 ce e3   blt    x17,   x4,   -36
  8192   00 18 08 13   addi   x16,   x16,  1
```

```
     2   00 00 05 37   lui      x10,   0
   128   01 63 00 77   dlp      x6,    44
   128   00 c3 83 df   zlp      x7,    12,   28
   128   00 00 00 33   nop
  8192   00 02 20 23   sw       x0,    0(x4)
   128   00 02 2a 03   lw       x20,   0(x4)
524288   00 07 1a 83   lh       x21,   0(x14)
524288   00 06 9b 03   lh       x22,   0(x13)
524288   20 06 a7 0b   fusedmac x14,   x13,  2,   128
524288   01 42 20 23   sw       x20,   0(x4)
  8192   00 42 02 13   addi     x4,    x4,   4
  8192   00 b6 86 b3   add      x13,   x13,  x11
  8192   f8 07 07 13   addi     x14,   x14,  -128
   128   00 c3 14 df   zlp      x6,    12,   164
   128   f0 02 02 13   addi     x4,    x4,   -256
```

|  |  |  |
|---|---|---|
| (a) | (b) | (c) |

FIGURE 5: Flow-generated C code of a convolution operation (a) in MobileNetV1 alongside its corresponding assembly in the baseline processor (v0) (b) and the fully extended processor (v4) (c). `blt` instruction in baseline is eliminated as a result of the ISA extension targeting hardware loops (`zol`)

instructions for performance optimizations. As a result, source files that could benefit from custom instruction optimizations should be compiled using the Chess front-end.

The Synopsys Chess compiler can be retargeted to a specific processor by defining "`chess_rewrite`" rules which specify how the custom instructions should be used. It should be noted that the actual PDG implementations of these instructions are opaque to the Chess compiler. The "`chess_rewrite`" rule applied to the "`mac`" instruction is as shown in listing 4:

Listing 4: Chess rewrite rule for mac instruction

```
chess_rewrite int mac_rule(int c, int a, int b)
{return c + a*b;} -> {return MAC(c,a,b);}
```

These rules allow the Chess compiler to identify patterns where the use of custom instructions may be more efficient than the baseline ISA instructions. We configured the Chess compiler so that a fully optimized C application is generated. This ensured compatibility with embedded RISC-V processors while leveraging the custom instructions of the target processor architecture for enhanced performance.

### E. Hardware Implementation

#### 1) RTL generation from ASIP Designer

The extensions obtained from Step 2 (in Fig. 1) of the flow are modeled in nML, with Fig. 6 presenting snippet of the `fusedmac` extension. Using the Go compiler, a synthesizable Verilog representation of the extended core is generated from this nML modeling of the obtained extensions. The corresponding hardware implementation for the `mac`, `add2i` and `fusedmac` extensions is depicted in Fig. 8. Additionally, for debugging purposes, On-Chip Debug (OCD) support is enabled by configuring the Go compiler to incorporate a JTAG interface within the processor module. Fig. 7 shows the key functional units within the extended trv32p3 processor.

The trv32p3 core implements a modified Harvard architecture. The memory interfaces that provide access to the data memory and program memory are defined in PDG code. We implement RAM blocks independent of the processor



FIGURE 6: nML model of `fusedmac` extension

module that are compatible with the interfaces defined in the PDG code. In this design, the processor data memory and program memory are implemented with the dedicated block RAM available within the ZCU104 FPGA. Logic specific to the target hardware must be added to provide an interconnection between the block RAM and the trv32p3 memory interfaces. In this case, the block RAM output register must be disabled. This reduces the read operation latency to a single clock cycle, as required by the trv32p3 interface.
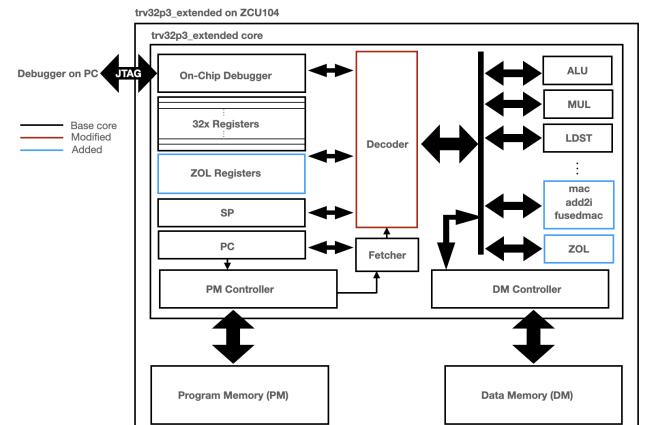


FIGURE 7: Custom RISC-V illustrating modifications on the baseline

#### 2) Implementation in Vivado

The hardware description generated by ASIP Designer is synthesized and implemented in Vivado (version 2024.2), as illustrated in Step 3 of Fig. 1. During this phase, a top-level
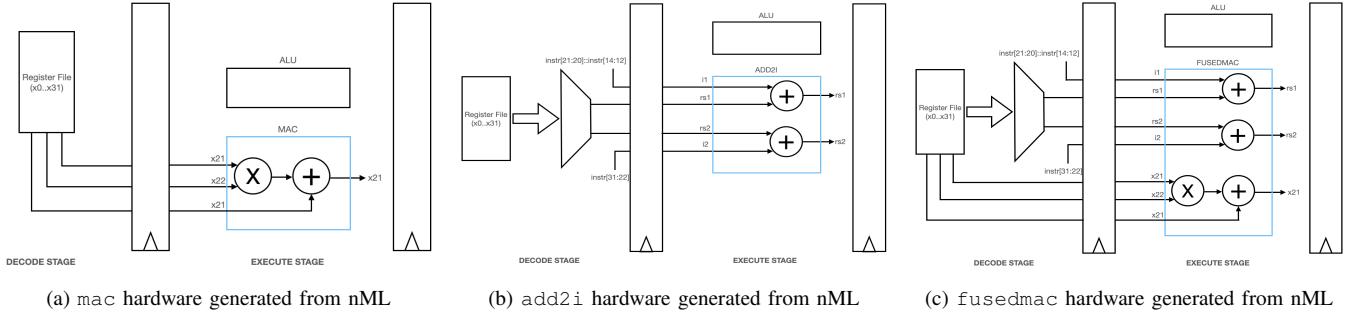
(a) `mac` hardware generated from nML    (b) `add2i` hardware generated from nML    (c) `fusedmac` hardware generated from nML

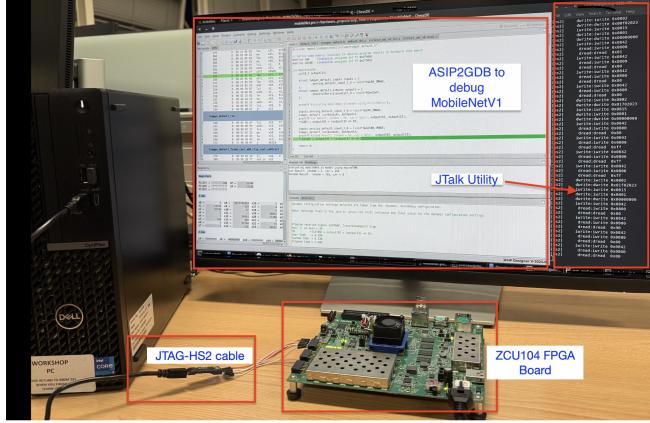FIGURE 8: Hardware diagrams generated from nML for `fusedmac`, `mac`, and `add2i` instructions.



FIGURE 9: Experimental setup showing MobileNetV1 deployed and debugged on RISC-V through JTAG

testbench is provided, defining the clock, reset, and JTAG interface signals. Additionally, synthesis and implementation constraints are specified to map the design pins to the corresponding physical pins on the ZCU104 FPGA.

### 3) On-Chip Debugging
The JTalk® cable driver is used to facilitate the on-chip debugging of applications running on the target hardware (step 4 in Fig. 1). This work uses the Diligent® JTAG HS2 cable. Synopsys ASIP2GDB debugging utility is used to communicate with the FPGA connected to the host PC running JTalk utility as shown in Fig. 9

## III. Results
The performance improvements and resource utilisations associated with the extended RISC-V core discussed in Section (C) are evaluated in this section. The performance is benchmarked with C implementations of CNN models obtained through the flow as described in Section (A). The performance improvement and memory utilisation (shown in Table 10) associated with each of the five processor model variations are compared.

TABLE 8: FPGA utilisation of all processor variants

| Processor | LUT | MUX | Registers | DSP | Power |
|---|---|---|---|---|---|
| **v0: Baseline** | 4,492 | 905 | 1,923 | 4 | 830 mW |
| **v1: v0 + mac** | 5,463 | 904 | 1,927 | 7 | 852 mW |
| **v2: v1 + add2i** | 6,409 | 912 | 1,946 | 7 | 850 mW |
| **v3: v2 + fusedmac** | 5,845 | 910 | 1,938 | 7 | 847 mW |
| **v4: v3 + hardware loops** | 6,207 | 910 | 2,268 | 7 | 849 mW |
| **Overhead:** | 1,715 | 5 | 345 | 3 | 19 mW |
| | (38.17%) | (0.5%) | (17.94%) | (75%) | (2.28%) |

### A. Implementation
The resource utilisation associated with each processor model variation is shown in both Fig. 10 and Table 8. The extended core (v4) comes with a LUT overhead of **38.17%** when compared to the baseline core. However, an overhead of only **0.5%** for multiplexers and **17.94%** for registers is associated with the fully extended core. The majority of the overhead comes from the addition of the `mac` and `add2i` instructions, with the `fusedmac` and `zol` instructions being more resource efficient in terms of LUT usage.

The power utilisation is estimated in the post-implemented design with typical process corner and default switching activity. The extended core (v4) is power efficient with an estimated power overhead of only **2.28%**.

### B. Performance Evaluation
Benchmarking was conducted on six neural network models: MobileNetV1, ResNet50, VGG16, MobileNetV2, and DenseNet121, generated using the proposed flow, along with a hand-coded C implementation of an LeNet-5-like CNN (denoted as LeNet-5* in this paper). The architecture of LeNet-5* model is detailed in Table 9. The inclusion of a hand-coded model demonstrates that the proposed extensions can be applied to both tool-generated and manually written C code, ensuring independence from specific tool artifacts. As long as the required patterns exist within the provided C code, the extensions remain applicable.

Simulations were performed using ASIP Designer as well as a hardware testbench implemented in Xillinx Vivado. The results from both simulation environments were identical and are presented in Fig. 11. The reported cycle count

(a) Look-Up Tables

(b) Multiplexers
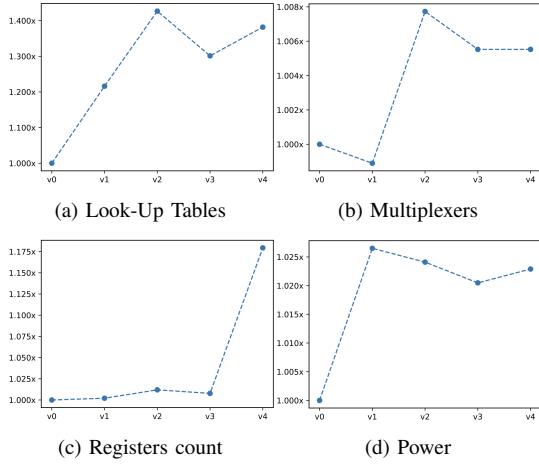
(c) Registers count

(d) Power

FIGURE 10: FPGA resource utilisation and power consumption with the addition of each successive extension, shown as a proportion of base core utilisation.

TABLE 9: LeNet-5* model structure

|  | First Layer | Second Layer | Third Layer |
|---|---|---|---|
| **Layer Type:** | Convolution | Convolution | MLP |
| **Input Size:** | 28×28 | 12×12×12 | 32×4×4 |
| **Output Size:** | 12×12×12 | 32×4×4 | 10×1 |
| **Activation:** | ReLU | ReLU | Softmax |
| **Kernel Size:** | 6×6 | 6×6 | N/A |
| **Stride:** | 2 | 2 | N/A |
| **Num. Filters:** | 12 | 32 | N/A |

represents the average number of cycles required for two CNN inference operations.

Energy efficiency is estimated for each model on each RISC-V variant as per (1)

$$E = P \times \left( \frac{C}{f} \right) \tag{1}$$

where E is the energy per inference, P is the estimated design power consumption, C is the number of clock cycles per inference and f is the processor clock frequency. In this design, the processor clock frequency is 100 MHz. While Synopsys trv cores target 80 MHz for prototyping by default, we evaluated 100 MHz to validate timing of processor version v4. This frequency was chosen as it meets timing without violations and requires no RTL modifications. The power utilisation is estimated by Xilinx Vivado based on the post-implementation design. Fig. 12 shows how the extensions reduce the energy consumed per CNN inference over different models.

## IV. Conclusion

This paper introduced MARVEL, an automated end-to-end framework that generates custom RISC-V ISA extensions specifically optimized for convolutional neural networks (CNNs) targeting resource-constrained IoT devices. MARVEL bridges high-level Python-based AI models and low-level bare-metal C implementations, enabling efficient deployment without operating system dependencies. The
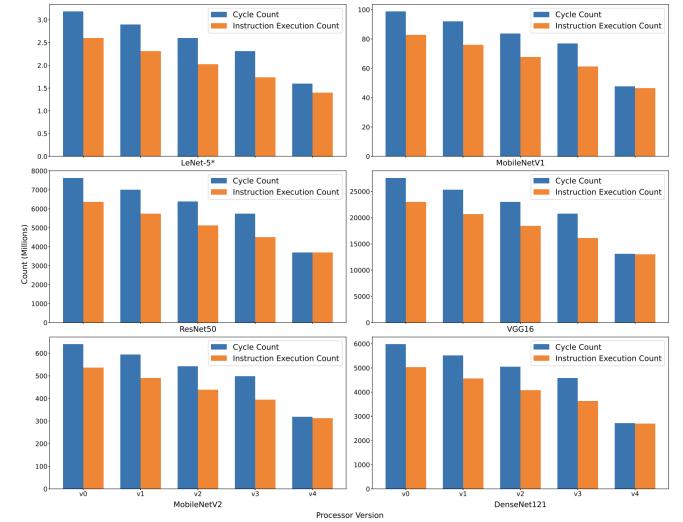


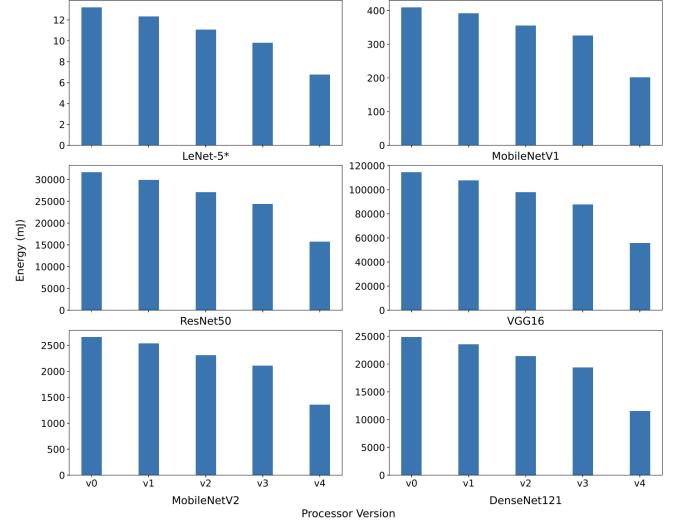FIGURE 11: Cycle and Instruction count for inferences across DNN models on five RISC-V variants



FIGURE 12: Energy/Inference on RISC-V variants

proposed framework demonstrated up to 2× improvement in inference speed and energy efficiency with minimal hardware overhead (28.23% area increase) compared to a baseline RISC-V core for popular DNN models. Future work includes refining power estimation accuracy, exploring additional RISC-V baselines, and extending support for diverse deep learning model classes with support for additional quantization levels, thereby enhancing edge AI hardware-software co-design.

TABLE 10: Comparison of Data and Program Memory Usage Across Processor Versions

| Processor | LeNet-5* | | MobileNetV1 | | ResNet50 | | VGG16 | | MobileNetV2 | | DenseNet121 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DM (kB) | PM (kB) | DM (kB) | PM (kB) | DM (MB) | PM (kB) | DM (MB) | PM (kB) | DM (MB) | PM (kB) | DM (MB) | PM (kB) |
| v0: Baseline | 60.83 | 1.47 | 622.04 | 24.19 | 43.65 | 39.50 | 15.76 | 13.18 | 4.19 | 40.86 | 14 | 105.14 |
| v1: v0 + mac | 31.53 | 1.46 | 592.74 | 24.12 | 43.62 | 39.28 | 15.73 | 13.21 | 4.19 | 40.70 | 14 | 104.39 |
| v2: v1 + add2i | 31.53 | 1.44 | 592.74 | 23.08 | 43.62 | 37.94 | 15.73 | 12.82 | 4.19 | 39.15 | 14 | 100.96 |
| v3: v2 + fusedmac | 31.53 | 1.43 | 592.74 | 22.97 | 43.62 | 37.73 | 15.73 | 12.77 | 4.19 | 38.97 | 14 | 100.49 |
| v4: v3 + hardware loops | 31.48 | 1.32 | 592.74 | 21.87 | 43.62 | 35.50 | 15.73 | 12.31 | 4.19 | 36.72 | 14 | 102.51 |
| Total Memory Saved (%) | 48.24 | 10.20 | 4.71 | 9.59 | 0.06 | 10.12 | 0.19 | 6.60 | 0 | 10.13 | 0 | 2.50 |

TABLE 11: Comparison of Related Works on ISA Extensions and HW-SW Co-Design for Edge AI

| Work | Input Format | Modifications to Input | ISA Extensions | HW-SW Co-Design | Automated End-to-End Flow |
|---|---|---|---|---|---|
| BARVINN [5] | ONNX | Required[1] | Generic | No | No |
| FlexACC [14] | C | Required[2] | Model class-specific | No | No |
| RISC-VTF [6] | Assembly | N/A | Transformer-specific | No | No |
| XPulpNN [10] | C[3] | Required | Generic | No | No |
| PULP [9] | C | Required | DSP-specific | No | No |
| AI-RISC [18] | DSL | Required | Generic | Modifies TVM-based flow | Yes |
| LiteAIR5 [29] | DSL | Required | Generic | Modifies TVM-based flow | Yes |
| **Ours** | DSL | No | Model class-specific | Yes | Yes |

[1] Requires the input to not have any residual connections. [2] Software mapping needed to map FlexACC primitives in C. [3] Manuscript description is ambiguous. DSL = Domain Specific Language (PyTorch/TensorFlow, etc.)

## REFERENCES

[1] R. Singh and S. S. Gill, "Edge ai: A survey," *Internet of Things and Cyber-Physical Systems*, vol. 3, pp. 71–92, Feb. 2023.

[2] E. Guthmuller et al., "Xvpfloat: Risc-v isa extension for variable extended precision floating point computation," *IEEE Transactions on Computers*, vol. 73, no. 7, pp. 1683–1697, 2024.

[3] H. Liu et al., "Convex: A risc-v instruction set extension scheme for accelerating convolution operations on mcus," in *Proceedings of the 2024 8th International Conference on High Performance Compilation, Computing and Communications*, ser. HP3C '24. Association for Computing Machinery, 2024, p. 133–138.

[4] R. Sahita et al., "Cove: Towards confidential computing on risc-v platforms," in *Proceedings of the 20th ACM International Conference on Computing Frontiers*, ser. CF '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 315–321.

[5] M. Askarihemmat et al., "Barvinn: Arbitrary precision dnn accelerator controlled by a risc-v cpu," in *Proceedings of the 28th Asia and South Pacific Design Automation Conference*, ser. ASPDAC '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 483–489.

[6] Q. Jiao, W. Hu, F. Liu, and Y. Dong, "Risc-vtf: Risc-v based extended instruction set for transformer," in *2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 2021, pp. 1565–1570.

[7] Y.-H. Chen et al., "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, 2017.

[8] D. Kadetotad et al., "An 8.93 tops/w lstm recurrent neural network accelerator featuring hierarchical coarse-grain sparsity for on-device speech recognition," *IEEE Journal of Solid-State Circuits*, vol. 55, no. 7, pp. 1877–1887, 2020.

[9] M. Gautschi et al., "Near-threshold risc-v core with dsp extensions for scalable iot endpoint devices," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 25, no. 10, pp. 2700–2713, 2017.

[10] A. Garofalo, G. Tagliavini, F. Conti, D. Rossi, and L. Benini, "Xpulpnn: Accelerating quantized neural networks on risc-v processors through isa extensions," in *2020 Design, Automation and Test in Europe Conference and Exhibition 2020*, 2020, pp. 186–191.

[11] X. Yu et al., "Cnn specific isa extensions based on risc-v processors," in *2022 5th International Conference on Circuits, Systems and Simulation (ICCSS)*, 2022, pp. 116–120.

[12] A. Al-Qawlaq, M. Ajay Kumar, and D. John, "Kwt-tiny: Risc-v accelerated, embedded keyword spotting transformer," in *2024 IEEE 37th International System-on-Chip Conference (SOCC)*, 2024, pp. 1–6.

[13] A. Berg et al., "Keyword transformer: A self-attention model for keyword spotting," in *Proceedings of Interspeech 2021*, August 2021.

[14] E.-Y. Yang et al., "Flexacc: A programmable accelerator with application-specific isa for flexible deep neural network inference," in *2021 IEEE 32nd International Conference on Application-specific Systems, Architectures and Processors (ASAP)*, 2021, pp. 266–273.

[15] A. Kumar M, V. Kumar, D. John, and S. Shanker, "Implementation and analysis of custom instructions on risc-v for edge-ai applications," in *Proceedings of the 14th International Symposium on Highly Efficient Accelerators and Reconfigurable Technologies*, ser. HEART '24. ACM, 2024, p. 126–129.

[16] G. Korol and A. C. S. Beck, " IoT–Edge Splitting With Pruned Early-Exit CNNs for Adaptive Inference ," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, no. 01, pp. 1–0, Mar. 5555.

[17] H. Sang et al., "An 2.31uj/inference ultra-low power always-on event-driven ai-iot soc with switchable nvsram compute-in-memory macro," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 71, no. 5, pp. 2534–2538, 2024.

[18] V. Verma, "Ai-risc: Scalable risc-v processor for iot edge ai applications," *University of Virginia, PHD Thesis*, 2022.

[19] Synopsys ASIP Designer. [Online]. Available: https://www.synopsys.com/dw/ipdir.php?ds=asip-designer

[20] trv32p3 processor core. [Online]. Available: https://www.synopsys.com/designware-ip/processor-solutions/asips-tools/asip-models.html#trv

[21] T. Chen et al., "Tvm: an automated end-to-end optimizing compiler for deep learning," in *Proceedings of the 13th USENIX Conference on Operating Systems Design and Implementation*, ser. OSDI'18. USA: USENIX Association, 2018, p. 579–594.

[22] EDGECORTIX, "Mera compiler and software framework." [Online]. Available: https://www.edgecortix.com/en/products/mera

[23] "Intermediate representation execution envioonment." [Online]. Available: https://github.com/iree-org/iree?tab=readme-ov-file

[24] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, "3d object representations for fine-grained categorization," in *2013 IEEE International Conference on Computer Vision Workshops*, 2013, pp. 554–561.

[25] COCO dataset. [Online]. Available: https://cocodataset.org/

[26] LiteRT overview. [Online]. Available: https://ai.google.dev/edge/litert

[27] T. Liang et al., "Pruning and quantization for deep neural network acceleration: A survey," *Neurocomputing*, vol. 461, pp. 370–403, 2021.

[28] Synopsys Go HDL Generator Manual. [Online]. Available: https://spdocs.synopsys.com/dow_retrieve/latest/sg/asip_designer/go-manual.pdf

[29] Y. Gao *et al.*, "Liteair5: A system-level framework for the design and modeling of ai-extended risc-v cores," in *2023 IEEE 36th International System-on-Chip Conference (SOCC)*, 2023, pp. 1–6.

[30] M. Ali *et al.*, "Rv-proviler: Evaluating risc-v isa for application-specific requirements," in *2024 IEEE Nordic Circuits and Systems Conference (NorCAS)*, 2024, pp. 1–7.

[31] V. Verma, T. Tracy II, and M. R. Stan, "Extrem-edge—extensions to risc-v for energy-efficient ml inference at the edge of iot," *Sustainable Computing: Informatics and Systems*, vol. 35, p. 100742, 2022.

[32] Synopsys Chess Compiler Manual. [Online]. Available: https://spdocs.synopsys.com/dow_retrieve/latest/sg/asip_designer/chess_user-manual.pdf

**Ajay Kumar M** received the B.Tech degree in Electronics and Communication Engineering from National Institute of Technology Karnataka, India, in 2017. During 2017-2022, he worked with Texas Instruments India and Intel India. He is currently a PhD student with the School of Electrical and Electronic Engineering, University College Dublin, Ireland. His research interests include Embedded Systems, Edge-AI, Digital hardware design, FPGA, and DNN accelerators.

**Cian O'Mahoney** received a B.Sc. degree in electronic and electrical engineering from University College Dublin, Dublin, Ireland, in 2023 and a M.E. in electronic and computer engineering from University College Dublin, Dublin, Ireland, in 2024. He is currently working as an FPGA Engineer in Susquehanna International Group LLP, Dublin, Ireland.

**Pedro Kreutz Werle** received the B.Sc. and M.Eng. degrees in electronic and computer engineering from University College Dublin, Dublin, Ireland, in 2024 and 2025, respectively. In 2024/2025 he worked at the Intelligent Systems Lab at the University College Dublin as a Graduate Researcher and join as a Research Assistant. His research interests include embedded AI and energy-efficient architectures.

**Shreejith Shanker** is an Assistant Professor in Reconfigurable Computing at the School of Engineering, Trinity College Dublin, Ireland. He obtained his B.Tech. degree in Electronics and Communication Engineering from the University of Kerala, India and his Ph.D. degree in Computer Science and Engineering from Nanyang Technological University (NTU), Singapore, in 2006 and 2016, respectively. From 2006 to 2008, he worked as an FPGA design and development engineer. From 2008 to 2011, he worked as a scientist in Digital Systems Group, Vikram Sarbhai Space Centre, Trivandrum, under the Indian Space Research Organisation (ISRO). From 2015 to 2016, he was a Research Fellow at the School of Computer Science and Engineering, NTU, Singapore. From 2017 to 2018, he was a Teaching Fellow at the School of Engineering, University of Warwick, UK. From 2018 to 2019, he was a Research Fellow at TUM CREATE Ltd, Singapore. Since 2019, he has been an Assistant Professor at the School of Engineering, Trinity College Dublin. He has also served as a member of the organising or technical committee for several IEEE/ACM conferences, such as FPL, DATE, FPT, ASAP, and HEART. He is also a reviewer for IEEE/ACM journals and conferences. His research interests include reconfigurable and adaptive accelerators, distributed embedded systems, and in-network computing.

**Dimitrios S. Nikolopoulos** is the John W. Hancock Professor of Engineering at Virginia Tech, with appointments in the Department of Computer Science and the Department of Electrical and Computer Engineering. He received a Ph.D. in Computer Engineering from the University of Patras, Greece, in 2000. His research interests include high-performance computing, parallel and distributed systems, computer architecture, and systems software. He is an IEEE Fellow.

**Bo Ji** received his Ph.D. degree in Electrical and Computer Engineering from The Ohio State University, Columbus, OH, USA, in 2012. Dr. Ji is an Associate Professor of Computer Science and a College of Engineering Faculty Fellow at Virginia Tech. Prior to joining Virginia Tech, he was an Associate Professor in the Department of Computer and Information Sciences at Temple University, where he was an Assistant Professor from July 2014 to June 2020. He was also a Senior Member of Technical Staff at AT&T Labs, San Ramon, CA, from January 2013 to June 2014. His research interests include the multidisciplinary intersections of Computing and Networking Systems, Artificial Intelligence and Machine Learning, Security and Privacy, and Extended Reality. He has been the general co-chair of IEEE/IFIP WiOpt 2021 and the technical program co-chair of IEEE/IFIP WiOpt 2026, ACM MobiHoc 2023, and ITC 2021, and he has also served on the editorial boards of multiple IEEE and ACM journals (IEEE/ACM Transactions on Networking, ACM SIGMETRICS Performance Evaluation Review, IEEE Transactions on Network Science and Engineering, IEEE Internet of Things Journal, and IEEE Open Journal of the Communications Society). Dr. Ji is a senior member of the IEEE and the ACM and a member of the AAAS. He was a recipient of the National Science Foundation (NSF) CAREER Award in 2017, the NSF CISE Research Initiation Initiative Award in 2017, the IEEE INFOCOM 2019 Best Paper Award, the IEEE/IFIP WiOpt 2022 Best Student Paper Award, the IEEE TNSE Excellent Editor Award in 2021, 2022, and 2024, and the Dean's Faculty Fellow Award from the College of Engineering at Virginia Tech in 2023.

**Hans Vandierendonck** received the Ph.D. and M.Sc. degrees from Ghent University, Ghent, Belgium, in 2004 and 1999, respectively. He is currently a Professor of high-performance and data-intensive computing with the School of Electronics, Electrical Engineering and Computer Science, and a Senior Member of ACM and IEEE. His research aims to discover novel solutions to data-intensive parallel systems and algorithms.

**Deepu John** is an Assistant Professor at the School of Electrical and Electronics Engineering, University College Dublin, Ireland. He received his Ph.D. in Electrical Engineering from the National University of Singapore in 2014. From 2014 to early 2017, he worked as a postdoctoral researcher at the Bio-Electronics Lab, National University of Singapore. Prior to this, he served as a senior engineer at Sanyo Semiconductors, Japan. He has received several awards, including the 2024 IEEE Transactions on Biomedical Circuits and Systems Best Paper Award, the Institution of Engineers Singapore Prestigious En-

gineering Achievement Award, Best Design Award at the Asian Solid-State Circuit Conference (2013), and IEEE Young Professionals, Region 10 individual award. He has served as an Associate Editor for IEEE Internet of Things Magazine and Guest Editor for IEEE Transactions on Circuits and Systems-I and IEEE Open Journal of Circuits and Systems. Currently, he serves as an Associate Editor for IEEE Transactions on Biomedical Circuits & Systems, Wiley International Journal of Circuit Theory & Applications and as a Senior Associate Editor for IEEE Transactions on Circuits & Systems-II. His research interests include IoT/Wearable sensing, Biomedical Circuits & Systems, and Edge Computing. He is a Senior Member of the IEEE.