# Sparq: A Custom RISC-V Vector Processor for Efficient Sub-Byte Quantized Inference

Théo Dupuis, Yoan Fournier, MohammadHossein AskariHemmat, Nizar El Zarif,
François Leduc-Primeau, Jean Pierre David, Yvon Savaria
*Department of Electrical Engineering, Polytechnique Montréal, Québec, Canada*

*Abstract*—Convolutional Neural Networks (CNNs) are used in a wide range of applications, with full-precision CNNs achieving high accuracy at the expense of portability. Recent progress in quantization techniques has demonstrated that sub-byte Quantized Neural Networks (QNNs) achieve comparable or superior accuracy while significantly reducing the computational cost and memory footprint. However, sub-byte computation on commodity hardware is sub-optimal due to the lack of support for such precision. In this paper, we introduce Sparq, a Sub-byte vector Processor designed for the AcceleRation of QNN inference. This processor is based on a modified version of Ara, an open-source 64-bit RISC-V "V" compliant processor. Sparq is implemented in GLOBAL FOUNDRIES 22FDX FD-SOI technology and extends the Instruction Set Architecture (ISA) by adding a new multiply-shift-accumulate instruction to improve sub-byte computation effciency. The floating-point unit is also removed to minimize area and power usage. To demonstrate Sparq performance, we implement an ultra-low-precision (1-bit to 4-bit) vectorized conv2d operation taking advantage of the dedicated hardware. We show that Sparq can significantly accelerate sub-byte computations with respectively 3.2 times, and 1.7 times acceleration over an optimized 16-bit 2D convolution for 2-bit and 4-bit quantization.

*Index Terms*—RISC-V, Vector ISA, Sub-byte, Convolution

## I. INTRODUCTION

Convolutional Neural Networks (CNNs) are used in a broad range of applications that include image processing, speech recognition, and natural language processing. As new CNN architectures and optimization methods allow for higher accuracy, the required amount of computations and on/off-chip memory accesses tends to increase. The high complexity and energy consumption pose a challenge for their deployment on resource-constrained devices such as low-power systems.

In the last few years, research has focused on minimizing computational cost and memory footprint while trying to preserve accuracy. This was achieved through the emergence of multiple new optimization techniques. Among them, reducing the bit-precision of weights and activations allows for a decrease in both the computational complexity and the memory footprint. In extreme cases, where weights and activations are quantized down to 1-bit or 2-bit, the resulting models typically experience an accuracy drop of less than 10% compared to their full-precision counterparts [1]–[5].

At the same time, new hardware solutions are proposed to address the issue of complexity through the emergence of vectorized Instruction Set Architectures (ISAs). Supporting these new instructions, vector processors can achieve high efficiency through parallelization of tensor computations while maintaining the flexibility of being fully programmable [6]–[8]. However, they are limited by the minimum granularity of their vector registers. Generally fixed at 8-bit, this limitation makes these architectures sub-optimal for sub-byte computing.

In this paper, we propose to modify Ara [6], a 64-bit RISC-V RVV1.0 compliant vector processor to add a new vectorized multiply-shift-accumulate instruction, `vmacsr`. In order to minimize area and power usage, we remove the vector floating-point unit (FPU). The main idea is to increase the performance of ultra-low-precision computations with a cost-effective instruction. To demonstrate the performance of Sparq, we implement several low-precision conv2d algorithms over a wide range of precisions, using the recently proposed ULPPACK [9] technique. We demonstrate that the new `vmacsr` instruction allows us to mitigate the constraints related to the ULPPACK technique and significantly improve performance. Using up to 2-bit quantization, we achieve a speedup of $3.2\times$ over an optimized 16-bit conv2d implementation and $1.7\times$ using up to 4-bit quantization. Lastly, we implement the processor in GLOBAL FOUNDRIES 22FDX FD-SOI technology and compare power and area usage between Sparq and Ara.

## II. BACKGROUND AND RELATED WORK

### A. Sub-byte Convolutional Neural Networks

Allowing to address both computation complexity and memory footprint, quantization has aroused a real interest. In fact, quantizing properly can achieve considerable memory savings with limited accuracy drop. On Deep Neural Networks such as ResNeXt [10] and EfficientNet [11], only a low accuracy degradation was observed when quantizing the full-precision 32-bit weights and activations down to 8-bit [12]. In recent work, ultra-low-precision quantization ($\leq$ 8-bit), which we refer to as sub-byte quantization, has demonstrated that minimal accuracy degradation was achievable with new network structures and training strategies. In some cases sub-byte Quantized Neural Networks (QNNs) perform better than their full-precision counterparts [13], [14] as shown in Table I.

TABLE I
ACCURACY ON QUANTIZED RESNET18 USING LG-LSQ [14]

| Dataset | Model | Precision (W/A) | Top-1 | Top-5 |
|---|---|---|---|---|
| ImageNet | Resnet18 | LG-LSQ(3/3) | 70.31 | 89.55 |
| | | LG-LSQ(4/4) | 70.78 | 89.77 |
| | | FP32 | 69.76 | 89.08 |

This was made possible by the introduction of architectures and methods focused on ultra-low-precision. LSQ [13], SSG [14] and DSQ [15] propose to learn the quantization parameters for both weights and activations by minimizing the quantization loss of the network during training. SAWB [3] focuses on the weight quantization parameter using the weight distribution to estimate the optimal quantization scale. To address the issue of unbounded activation range after ReLU, PACT [3] proposes to train a clipping parameter to find the balance point between clipping and quantization error. Focusing on binary or ternary quantization, CNN architectures tailored for 1-bit or 2-bit such as XNORNet [2] or BinaryNet [4] have been introduced. While achieving very high efficiency and low memory footprint, their accuracy suffers a significant drop when compared to full-precision models. All these methods try to minimize the accuracy drop of sub-byte QNN by minimizing the error due to quantization on weights and activations.

### B. Sub-byte Computation on Commodity Hardware

CNNs rely heavily on the use of the *conv2d* operation, which can make up to 90% of the computation [16], [17]. By using ultra-low-precision quantization, or specific binary/ternary architectures, the computational cost can be drastically reduced. Typically, general-purpose processors cannot take advantage of sub-byte operands since their ISA is suited for byte instructions at best. The same applies to their corresponding vector extension ISA whose granularity is typically limited to vectorized 8-bit instructions. Thus, naive implementations of sub-byte algorithms are bounded by the performance of the 8-bit implementation.

Several acceleration algorithms have been introduced to counterbalance this issue. For example, bit-serial computation allows each operand bit to be processed in a serial manner. As a result, the operation between any arbitrary $N$-bit and $M$-bit precision operand can be computed. However, the computation complexity is defined by $\mathcal{O}(N \times M)$, meaning that this method is only suitable for ultra-low-precision, typically less than 3-bit. For this precision range, the performance over fp32 can be multiplied by a factor of 2 to 6 depending on the precision [18]. This paved the way for the development of specialized bit-serial hardware [19], [20].

Other techniques, such as ULPPACK, offer a less invasive alternative while achieving higher performance than standard precision implementation. Even though the expected speedup is lower than with bit-serial, typically $2\times$ to $4\times$ acceleration over fp32 [9], this method covers a wider precision range, typically 1-bit to 4-bit, by utilizing densely packed operands with commodity Single Instruction, Multiple Data (SIMD) architectures. Although no specialized hardware is required, this algorithm can benefit from specialized instruction such as a multiply-shift-accumulate.

### III. Software implementations

In this section, we provide an overview of different precision *conv2d* algorithms benchmarked on Sparq. Every implementation, except for the provided benchmarks, is handwritten
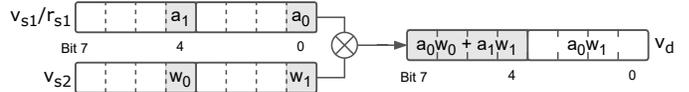


Fig. 1. Multiplication between 2 packed 8-bit vectors with 1-bit precision on weights and activation. The dot product of the two packed elements is computed using a single multiplication on the MSB. The result then needs to be shifted to retrieve the dot product.

using inline assembly, unrolled, and stored using a channel-first memory layout for the input, kernel, and output tensors.

### A. Optimized Vector Conv2d

A wide range of benchmarks is provided with Ara[1], including a double-precision (DP) *conv2d* function. The DP implementation achieves high utilization of Ara's lanes by using the available optimized vector slides to maximize efficiency through high data reuse. We implemented int16 and fp32 *conv2d* based on the structure of the convolution example to serve as the baseline for our comparisons. The choice of a dedicated convolution algorithm over an image to column (*im2col*) operation followed by a GEneral Matrix Multiplication (GEMM) technique is motivated by the reduction of the memory footprint induced by the *im2col* operation. Thus, fewer on/off-chip memory accesses are required to compute the output. We benchmarked the int16 *conv2d* on Sparq and the fp32 version on Ara. Our implementations of int16 and fp32 achieve a lane utilization of respectively 93.8% and 93.6% at $1 \times 32 \times 512 \times 512$ input size.

### B. ULPPACK

ULPPACK [9] is a software-only technique aiming towards acceleration of ultra-low-precision computation via an effective packing operand scheme. As illustrated in Figure 1 with an 8-bit register example, by packing multiple low-precision operands in wider registers, a single multiplication instruction completes the dot product between multiple operands and returns the result on the 4 most significant bits. Formally, the result can be expressed as follows:

$$(a_0 + 2^4 a_1) \times (w_1 + 2^4 w_0)$$
$$= 2^8 a_1 w_0 + 2^4 (a_0 w_0 + a_1 w_1) + a_0 w_1$$

The dot product of packed vectors has to be extracted through a right logical shift (LSR) before accumulation. Local accumulations of the non-shifted result can be done to alleviate the number of shift operations, but are limited by the width of the low-precision result. Moreover, increasing the precision of weights and activations results in a reduction in the number of possible local accumulations due to overflow, implying a lower speedup. In the Figure 1 example, using 8-bit elements to compute the product of 2 packed operands limits the result on a 4-bit width. By using 1-bit precision on weights and activation, 8 local accumulations are possible without risking overflow. This method is an efficient way to compute vector-based algorithms since it computes the dot product between multiple low-precision operands with a unique higher precision

[1] See https://github.com/pulp-platform/ara

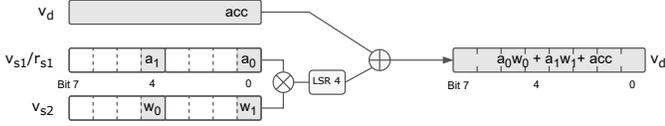multiplication, which significantly reduces the computational cost and resource usage.



Fig. 2. Diagram of the multiply-shift-accumulate operation on 8-bit packed registers with 1-bit precision for activations and weights. A shifter is inserted between the multiplication and the accumulation.

## IV. Proposed Architecture

### A. vmacsr Custom Instruction

To address the issue of local accumulation overflow described in Section III-B, we propose to implement a multiply-shift-accumulate operation as a custom vector instruction. As depicted in Figure 2, using the multiply-shift-accumulate operation, the low-precision multiplication is shifted and accumulated. This avoids the need for an intermediate register, a logical right shift, and an addition.
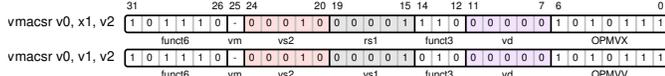


Fig. 3. vmacsr encoding example.

In this work, we mainly focus on a specific precision range following the condition $N + M \leq 7$-bit, where $N$ and $M$ are the operand's precision in bits. In this range, quantized models offer sufficient precision to achieve comparable or superior accuracy compared to full-precision models. We implement a *conv2d* with 16-bit packed register granularity to achieve high performance on this range. However, the use of multiply-shift-accumulate is applicable on higher precision operands and elements width, addressing the same overflow issue. To implement the new instruction, we modified Ara's dispatcher to specify the vmacsr funct6 encoding. As our instruction behaves very closely to vmacc, we use the free funct6 encoding following vmacc's funct6 [21]. We implement vmacsr in both OPMVV and OPMVX format allowing vector-vector and vector-scalar instructions to be used. The encoding of vmacsr is described in Figure 3. It should be noted that for this work, only the vector-scalar vmacsr instructions is used. We then modified Ara's SIMD multiplier to describe the expected behavior of the instruction, denoted as follows:

$$V_d \leftarrow V_d + ((V_{s1} \times V_{s2}) >> M)$$

where $M$ is the shifted value on the resulting product before accumulation. In this work, we only pack two operands per register to reduce the overhead related to the packing operation carried out at runtime. Thus, $M$ is hard-wired at half the granularity of the vector registers.

### B. 2D Convolution Algorithm

To our knowledge, the ULPPACK method has only been implemented on an ARM CPU using the ARM Neon intrinsic

**Algorithm 1** Proposed vector conv2d using vmacsr

| | |
|---|---|
| $H, C$ | : Input height, channels |
| $F_h, F_w$ | : Kernel height and width |
| $M$ | : Packed operand per element |
| $R, V$ | : Scalar and vector registers (initialized to zero) |

1: **for** $h \leftarrow 1$ to $H$ **do**
2:     $V_{F_h} \leftarrow 0$
3:     **for** $c \leftarrow 1$ to $\frac{C}{M}$ **do**
4:         $V_0 \leftarrow$ load one packed input row
5:         **for** $i \leftarrow 1$ to $F_w$ **do**
6:             $R_{[1:F_h]} \leftarrow$ load the $i^{th}$ packed kernel column
7:             **for** $j \leftarrow 1$ to $F_h$ **do**
8:                 $V_j \leftarrow$ vmacsr$(R_j, V_0, V_j)$
9:             $V_0 \leftarrow$ vslidedown$(V_0, 1)$
10:     **if** $h \geq F_h$ **then**
11:         $O[h] \leftarrow V_1$       ▷ store one output row
12:         **for** $j \leftarrow 1$ to $F_h - 1$ **do**
13:             $V_j \leftarrow V_{j+1}$

ISA extension [8], [9]. In addition, by rearranging the input and filters with an *im2col* operation, the convolution is performed with a GEMM. However, Ara's dedicated *conv2d* algorithm takes advantage of the vslidedown instruction provided by the RISC-V "V" ISA extension, which does not have any ARM equivalent. Thereby, the output stationary ULPPACK *conv2d* algorithm presented in Algorithm 1 is based on the *conv2d* described in Section III-A. Although the operand packing operation is not detailed in the algorithm, it is carried out at runtime. We use the ULPPACK P1 packing [9] scheme to compute the contribution of $M$ channels at each iteration, where $M$ is the number of operands packed per register. Following the example in Figure 1, $a_i$ and $w_j$ respectively represent activation and weight values from the channel $i$ and $j$.

Algorithm 1 presents a simplified ultra-low-precision *conv2d* implementation using the vmacsr instruction available on Sparq. For each loaded packed input row (line 4), the algorithm performs a vector multiply-shift-accumulate operation (line 8) to compute the partial result with the first packed kernel column. The input is then slided by one element to the left (line 9) to accommodate the next kernel column, and this operation is repeated until all the kernel columns and input channels have been processed. In order to complete and store the first output row (line 11), $F_h$ packed input rows must be processed (line 10). The partial results contained in the following vector registers (line 13) are moved down to be completed one by one with each new packed input row.

## V. Results and Performance Analysis

### A. Performance Analysis

We compare the performance of the different *conv2d* implementations over several precisions for activations and weights using RTL simulations with a 4-lane configuration for both Ara and Sparq. Our measured execution time includes both activations and weights packing done at runtime. However, the overhead induced by weights packing could be avoided by offline preprocessing. Figure 4 presents the performance in operations per cycle over the different *conv2d* implementations.
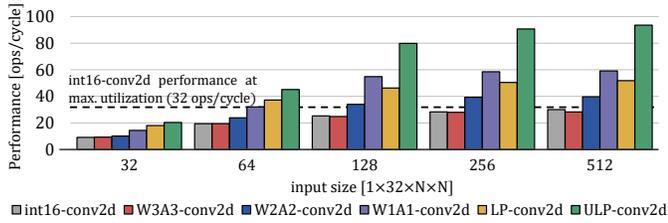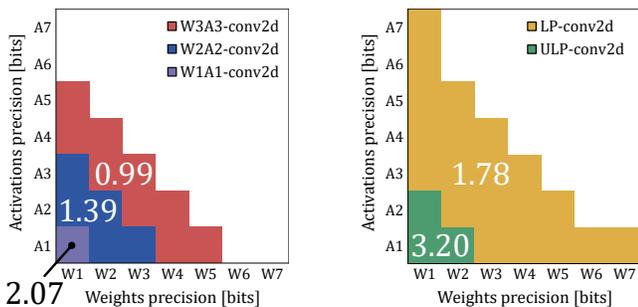
Fig. 4. Performance comparison between several *conv2d* implementation using a $7 \times 7$ kernel size. Ultra-Low-Precision (ULP) and Low-Precision (LP) *conv2d* take advantage of the `vmacsr` instruction on Sparq while W1A1, W2A2, W3A3 *conv2d* run on native RISC-V "V" ISA.

As expected, the performance of ULPPACK running on Ara shows an improvement over its 16-bit counterpart. Using $7 \times 7$ favors high data reuse, hence the small gap with the theoretical throughput. Moreover, using Sparq's `vmacsr` instruction removes the constraint of local accumulation, offering two major benefits:

- **Performance increase** resulting from the expected overall reduction in the number of instructions to compute the output matrix as presented in Figure 4.
- **Higher precision range**, as shown in Figure 5, is obtained without modifying the algorithm. The range is only limited by the 8-bit wide dot product result when using 16-bit packed registers for Low-Precision (LP) and the 4-bit result for the Ultra-Low-Precision (ULP).



(a) Native implementation  (b) `vmacsr` implementation

Fig. 5. Relative speedup over the 16-bit *conv2d* implementation on the overflow-free precision region. The kernel size is $7 \times 7$ over a $32 \times 256 \times 256$ input. (a) for the native implementation benchmarked on Ara (b) for the accelerated version taking advantage of the `vmacsr` instruction benchmarked on Sparq.

### B. Physical Implementation

Since the `vmacsr` circuit is located in the vector engine, we only implemented one lane of Ara and Sparq using GLOBALFOUNDRIES 22FDX FD-SOI technology. For synthesis and back-end, we used the SYNOPSYS DESIGN COMPILER S-2021.06-SP5 and CADENCE INNOVUS v21.15, respectively. The physical implementation results of these lanes are summarized in Table II. As anticipated, the Sparq lane exhibits significant improvements in terms of area (-43.3%) and power consumption (-58.8%) when compared to a standard Ara lane, primarily due to the FPU removal. Furthermore, the inclusion of the `vmacsr` instruction did not impact the typical corner frequency, indicating that it did not affect the critical path of the design. In fact, the removal of the FPU actually

yielded an increase (+8.7%) in the maximum lane clock speed. It is important to note that in Ara, the critical path is primarily located in the VLSU and SLDU, both of which are part of the interconnection between lanes [6]. These interconnections were not implemented as part of this work since the addition of `vmacsr` would have had a negligible impact on those modules. Figure 6 shows the physical layout of both designs.
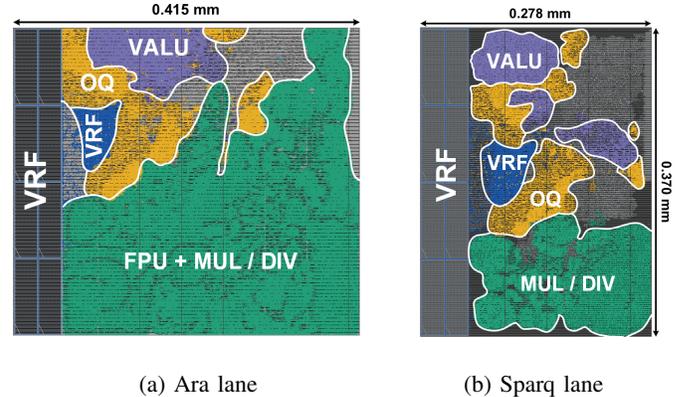


(a) Ara lane  (b) Sparq lane

Fig. 6. Ara and Sparq lanes placed and routed designs. ■ is the vector register file, ■ is the operands queue, ■ is the vector fixed point multiplication and division units with or without the FPU, ■ is the vector ALU.

TABLE II
PHYSICAL IMPLEMENTATION OF ARA AND SPARQ LANES

|  | Ara Lane | Sparq Lane |
|---|---|---|
| Number of Lanes | 4 | 4 |
| VRF Size [KiB] | 16 | 16 |
| Lane Cell Area [mm$^2$] | 0.120 | 0.068 |
| Lane Core Frequency [GHz] | 1.346 | 1.464 |
| Lane Power [mW] | 159.2 | 65.6 |

At typical corner (TT/0.8V/25°C)

### VI. CONCLUSION

This paper introduced the extension of the RISC-V "V" ISA with a multiply-shift-accumulate custom instruction, allowing better performance of sub-byte computations. We presented Sparq, a FPU-free RISC-V vector processor supporting the multiply-shift-accumulate instruction. To demonstrate its improved performance, we implemented a *conv2d* algorithm for 1 to 4-bit precision operands. In this range, quantized models offer sufficient precision to achieve comparable or superior levels of accuracy compared to full-precision models. We showed that, with `vmacsr`, we achieved up to $1.7\times$ and $3.2\times$ speedup over an optimized 16-bit *conv2d* depending on the precision. Lastly, the implementation reports of the modified version in GF22 nm showed that removing the FPU significantly reduced the area and power usage, while the newly added instruction had no impact on the critical path.

In future work, we plan to improve the flexibility of `vmacsr` using a runtime configurable shifter, as well as testing the proposed algorithm on an FPGA emulation of Sparq.

### VII. ACKNOWLEDGMENTS

## REFERENCES

[1] Z. Liu, B. Wu, W. Luo, X. Yang, W. Liu, and K.-T. Cheng, *Bi-Real Net: Enhancing the Performance of 1-Bit CNNs with Improved Representational Capability and Advanced Training Algorithm: 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XV*, 09 2018, pp. 747–763.

[2] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," vol. 9908, 10 2016, pp. 525–542.

[3] J. Choi, S. Venkataramani, V. V. Srinivasan, K. Gopalakrishnan, Z. Wang, and P. Chuang, "Accurate and efficient 2-bit quantized neural networks," in *Proceedings of Machine Learning and Systems*, A. Talwalkar, V. Smith, and M. Zaharia, Eds., vol. 1, 2019, pp. 348–359. [Online]. Available: https://proceedings.mlsys.org/paper/2019/file/006f52e9102a8d3be2fe5614f42ba989-Paper.pdf

[4] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1," 2016. [Online]. Available: https://arxiv.org/abs/1602.02830

[5] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," *NIPS*, vol. 28, 11 2015.

[6] M. Cavalcante, F. Schuiki, F. Zaruba, M. Schaffner, and L. Benini, "Ara: A 1-GHz+ scalable and energy-efficient RISC-v vector processor with multiprecision floating-point support in 22-nm FD-SOI," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 28, no. 2, pp. 530–543, feb 2020. [Online]. Available: https://doi.org/10.1109%2Ftvlsi.2019.2950087

[7] M. Perotti, M. Cavalcante, N. Wistoff, R. Andri, L. Cavigelli, and L. Benini, "A "new ara" for vector computing: An open source highly efficient RISC-v v 1.0 vector processor design," in *2022 IEEE 33rd International Conference on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE, jul 2022. [Online]. Available: https://doi.org/10.1109%2Fasap54787.2022.00017

[8] N. Stephens, S. Biles, M. Boettcher, J. Eapen, M. Eyole, G. Gabrielli, M. Horsnell, G. Magklis, A. Martinez, N. Premillieu, A. Reid, A. Rico, and P. Walker, "The ARM scalable vector extension," *IEEE Micro*, vol. 37, no. 2, pp. 26–39, mar 2017. [Online]. Available: https://doi.org/10.1109%2Fmm.2017.35

[9] J. Won, J. Si, S. Son, T. J. Ham, and J. W. Lee, "Ulppack: Fast sub-8-bit matrix multiply on commodity simd hardware," in *Proceedings of Machine Learning and Systems*, D. Marculescu, Y. Chi, and C. Wu, Eds., vol. 4, 2022, pp. 52–63. [Online]. Available: https://proceedings.mlsys.org/paper/2022/file/14bfa6bb14875e45bba028a21ed38046-Paper.pdf

[10] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5987–5995.

[11] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 09–15 Jun 2019, pp. 6105–6114. [Online]. Available: https://proceedings.mlr.press/v97/tan19a.html

[12] M. Nagel, M. Fournarakis, R. A. Amjad, Y. Bondarenko, M. van Baalen, and T. Blankevoort, "A white paper on neural network quantization," *CoRR*, vol. abs/2106.08295, 2021. [Online]. Available: https://arxiv.org/abs/2106.08295

[13] S. K. Esser, J. L. McKinstry, D. Bablani, R. Appuswamy, and D. S. Modha, "Learned step size quantization," in *International Conference on Learning Representations*, 2020. [Online]. Available: https://openreview.net/forum?id=rkgO66VKDS

[14] S.-T. Lin, Z. Li, Y.-H. Cheng, H.-W. Kuo, C.-C. Lu, and K.-T. Tang, "Lg-lsq: Learned gradient linear symmetric quantization," 2022. [Online]. Available: https://arxiv.org/abs/2202.09009

[15] R. Gong, X. Liu, S. Jiang, T. Li, P. Hu, J. Lin, F. Yu, and J. Yan, "Differentiable soft quantization: Bridging full-precision and low-bit neural networks," 10 2019, pp. 4851–4860.

[16] J. Cong and B. Xiao, "Minimizing computation in convolutional neural networks," in *Artificial Neural Networks and Machine Learning – ICANN 2014*, S. Wermter, C. Weber, W. Duch, T. Honkela, P. Koprinkova-Hristova, S. Magg, G. Palm, and A. E. P. Villa, Eds. Cham: Springer International Publishing, 2014, pp. 281–290.

[17] T. Abtahi, A. Kulkarni, and T. Mohsenin, "Accelerating convolutional neural network with fft on tiny cores," in *2017 IEEE International Symposium on Circuits and Systems (ISCAS)*, 2017, pp. 1–4.

[18] M. Cowan, T. Moreau, T. Chen, J. Bornholt, and L. Ceze, "Automatic generation of high-performance quantized machine learning kernels," 02 2020, pp. 305–316.

[19] P. Judd, J. Albericio, T. Hetherington, T. M. Aamodt, and A. Moshovos, "Stripes: Bit-serial deep neural network computing," in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2016, pp. 1–12.

[20] M. AskariHemmat, T. Dupuis, Y. Fournier, N. E. Zarif, M. Cavalcante, M. Perotti, F. Gurkaynak, L. Benini, F. Leduc-Primeau, Y. Savaria, and J.-P. David, "Quark: An integer risc-v vector processor for sub-byte quantized dnn inference," 2023. [Online]. Available: https://arxiv.org/abs/2302.05996

[21] K. Asanović and A. Waterman, "The risc-v vector extension, v1.0," 2021. [Online]. Available: https://github.com/riscv/riscv-v-spec/releases/download/v1.0/riscv-v-spec-1.0.pdf