

IndexMAC: A Custom RISC-V Vector Instruction to Accelerate Structured-Sparse Matrix Multiplications

V. Titopoulos, K. Alexandridis, C. Peltekis
Electrical and Computer Engineering
Democritus University of Thrace, Greece

C. Nicopoulos
Electrical and Computer Engineering
University of Cyprus, Cyprus

G. Dimitrakopoulos
Electrical and Computer Engineering
Democritus University of Thrace, Greece

Abstract—Structured sparsity has been proposed as an efficient way to prune the complexity of modern Machine Learning (ML) applications and to simplify the handling of sparse data in hardware. The acceleration of ML models – for both training and inference – relies primarily on equivalent matrix multiplications that can be executed efficiently on vector processors or custom matrix engines. The goal of this work is to incorporate the simplicity of structured sparsity into vector execution, thereby accelerating the corresponding matrix multiplications. Toward this objective, a new vector index-multiply-accumulate instruction is proposed, which enables the implementation of low-cost indirect reads from the vector register file. This reduces unnecessary memory traffic and increases data locality. The proposed new instruction was integrated in a decoupled RISC-V vector processor with negligible hardware cost. Extensive evaluation demonstrates significant speedups of $1.80\times$ – $2.14\times$, as compared to state-of-the-art vectorized kernels, when executing layers of varying sparsity from state-of-the-art Convolutional Neural Networks (CNNs).

Index Terms—Structured sparsity, Matrix multiplication, Vector processor, Machine learning accelerator.

I. INTRODUCTION

The computation of Machine Learning (ML) models relies primarily on equivalent matrix multiplications that can be efficiently executed by vector processors [1], or specialized matrix engines built on top of systolic arrays [2]. To reduce memory storage and computation cost, the weights of ML models are often pruned, thereby leading to sparse models [3]. The derived zero weights are not stored and the corresponding computation is skipped.

The achieved sparsity can either be *unstructured* [4], or *structured* [5], [6]. In unstructured sparsity, there is no constraint on the locations of the zeros, as shown in Fig. 1(a). In this case, together with the non-zero elements, multiple indexes are also required to identify the original position of each non-zero element.

On the contrary, in structured sparsity, there is an upper limit on the number of non-zero elements that may be present within a block of consecutive elements. For instance, in Fig. 1(b), for every 4 elements in each row, there are up to two non-zero elements. Such block-based structure simplifies both the indexing required to identify the position of each non-zero element inside each block, and the hardware needed to operate on such sparse data. In most practical applications [5], [7], blocks are small and $N:M$ sparsity patterns of 1:2, 1:4 or 2:4 are supported, where each block of M elements may contain up to N non-zero elements.

This work was supported by a research grant from Codasip, a provider of customizable RISC-V IP and Codasip Studio design toolset, and its University Program to Democritus Univ. of Thrace for “RISCV vector processor design”.

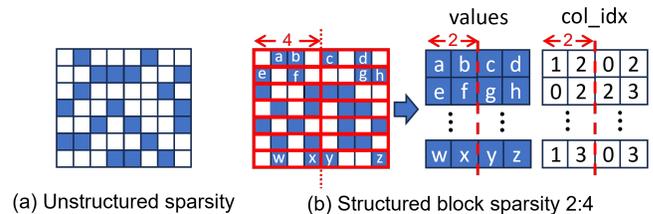


Fig. 1. Example of (a) unstructured sparsity; and (b) structured block sparsity of 2:4 (i.e., up to 2 non-zero elements in every 4 consecutive elements) and their respective representation. Each blue square represents a non-zero element.

Structured sparsity promises high performance and low storage overhead. However, in certain cases, it may also lead to less accurate ML models, due to the constrained sparsification process [3], [8]. If sparsification involves model retraining, the possible accuracy loss is ameliorated by allowing the model to adapt to the removal of certain weights.

In this work, our goal is to incorporate the simplicity of structured sparsity into vector execution at the minimum hardware cost. Prior work on vector processors exploits unstructured sparsity to achieve significant speedups. Various sparse vector matrix multiplication algorithms have been developed [9], [10], and recently extended [11], [12], with the goal being to improve performance by efficiently handling unpredictable sparsity patterns. Besides such software-only approaches, VIA [13] adds a scratchpad memory and new custom vector instructions to deal with the complexity of unstructured sparse data.

Aiming to avoid the substantial hardware overhead incurred in exploiting *unstructured* sparsity, this work makes the case that *structured* sparsity, which appears frequently in state-of-the-art CNN applications [5], can be exploited very effectively – and with negligible hardware cost – in vector processors to achieve high-performance sparse \times dense matrix ($A \times B$) multiplications.

To achieve this goal we follow a three-fold approach. First, we adopt a vectorized kernel of the row-based matrix multiplication algorithm that yields good performance under both dense and sparse workloads [11], [14], [15]. Other widely-applicable vectorized algorithms for sparse data, such as SPA [9] and its variants [11], [12], are avoided, since they target highly sparse matrices and do not perform as well with structured sparse data. Second, tiles of matrix B , which is treated as dense, are loaded in the vector register file to increase locality. Pre-loaded tiles are replaced only when they have been fully used for matrix multiplication and

are no longer needed. The third and most important step is the introduction of a new vector index-multiply-accumulate instruction that enables the implementation of low-cost indirect reads from the portion of the vector register file where tiles of matrix B reside. This new instruction exploits the structured-sparse format of matrix A to reduce effectively the number of instructions executed per iteration.

Overall, the contributions of this work can be summarized as follows:

- A state-of-the-art vectorized sparse matrix multiplication [11], [15] is re-organized to enable the handling of structured sparse data in long-vector ISAs, such as RISC-V. Structured sparsity of one matrix allows us to keep tiles of the other dense matrix in the vector register file, thus improving data locality and simplifying loop unrolling [16], [17].
- The reformed vectorized matrix multiplication is accelerated with a new index-multiply-accumulate (`vindexmac`) instruction, which enables the implementation of low-cost indirect reads from the vector register file. This eliminates unnecessary memory traffic often encountered in sparse matrix multiplication algorithms and reduces the number of instructions per iteration.
- The proposed vector instruction was integrated in a high-end decoupled RISC-V vector processor at negligible hardware cost. Extensive evaluation in Gem5 [18], [19] of the vector processor attached to a superscalar out-of-order core demonstrates significant speedups of $1.80\times$ – $2.14\times$, as compared to a state-of-the-art vectorized kernel, when executing state-of-the-art Convolutional Neural Networks (CNNs) pruned for structured sparsity.

The rest of the paper is organized as follows: Section II presents the formulation of a state-of-the-art vectorized algorithm for sparse matrix multiplications with structured-sparse data. Section III introduces the proposed vector index-multiply-accumulate instruction and its hardware implementation. Experimental results are presented in Section IV and conclusions are drawn in Section V.

II. VECTOR SPARSE \times DENSE MATRIX MULTIPLICATION

Vectorized matrix multiplications with sparse data can be implemented with many approaches [9], [11], [12]. The row-wise approach [14], [15], [20], also known as Gustavson’s algorithm [21], has been shown to be highly effective in computing the matrix product $A \times B$, and it is a better fit to the targeted structured sparsity context. Other approaches [9], [11], [12] target extremely high sparseness and are not as effective with structured sparsity that exhibits modest sparsity.

Matrix A is sparse and assumed to follow a structured-sparsity template, and, without loss of generality, matrix B is considered to be dense. Algorithmically, all the non-zero elements in a single row of matrix A should be multiplied in parallel with the corresponding rows of matrix B , where the row index of matrix B is determined by the column index of the non-zero value in matrix A . The product of the multiplication is produced row-by-row, as follows:

$$C[i, :] = \sum_k A[i, k]B[k, :] \quad (1)$$

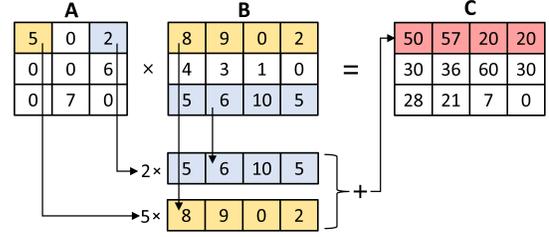


Fig. 2. Row-wise matrix multiplication to compute output row $C[0, :]$.

Fig. 2 illustrates a simple example of how row 0 of the result matrix C is produced. Row-wise matrix multiplication can be easily vectorized, since each element of A is multiplied with *all* the elements of a row of matrix B . This yields a vector of partial results for a row in the result matrix C .

Algorithm 1 depicts the vectorized form of the row-wise matrix-multiplication algorithm, *ignoring for clarity* any sparseness in matrix A and assuming an arbitrarily large vector length. In line 3, all the elements of row i of matrix A are loaded as a vector into the vector register file. Similarly, all the elements of the corresponding row of matrix B are also loaded as a vector in line 5. The vector multiplication between the first element of the loaded row of A and the entire row of matrix B is performed in line 7, which also accumulates the partial results. The row of A is shifted to the right by one element in line 8 to enable the repetition of the above process for the remaining elements of the same row. The final values of all the elements of the corresponding row of the result matrix C are stored back in memory in line 10. The process is repeated until all the rows of matrix C are produced.

Algorithm 1 Row-wise vector matrix multiplication

- 1: set vector length
 - 2: **for** $i=0$ **until** $i=\text{num_of_rows_of_A}-1$ **do**
 - 3: vload $A[i, :]$ ▷ load the i th row of A
 - 4: **for** $j=0$ **until** $j=\text{num_of_columns_of_A}-1$ **do**
 - 5: vload $B[j, :]$ ▷ load the corresponding row of B
 - 6: $s0 = A[i, 0]$ ▷ move $A[i, 0]$ to a scalar reg
 - 7: $C[i, :] += s0 * B[j, :]$ ▷ scalar-vector mul-acc
 - 8: $A[i, :] \gg 1$ ▷ vector slide to the right
 - 9: **end for**
 - 10: vstore $C[i, :]$ ▷ store to mem the i th row of C
 - 11: **end for**
-

The baseline vectorized row-based matrix multiplication algorithm can be reformulated to support a *structured-sparse* matrix A . This is shown in Algorithm 2, where all mandatory changes/additions are highlighted with color. The first essential difference pertains to the loading of matrix A . Instead of whole rows of A , the `values` and `col_idx` (index) vectors that refer *only* to the *non-zero* elements of A are loaded (lines 3 and 4). The second fundamental difference concerns the selection of the corresponding rows of B . In this case, only the rows that refer to the column indexes of the non-zero elements of A are selected (lines 7 and 8) to participate in the multiply-and-accumulate operation (line 10). To move to the next non-zero element, the `values` and `col_idx` vectors are slid to the right (lines 11 and 12).

Algorithm 2 Row-wise vector sparse-dense matrix multiplication

```

1: set vector length
2: for i=0 until i=num_of_rows_of_A-1 do
3:   vload values[i, :] ▷ load ith row of values of A
4:   vload col_idx[i, :] ▷ load ith row of col. indexes of A
5:   col_idx[i, :] += B_address ▷ adjust load addresses
6:   for j=0 until j=nonzero_elems_per_block do
7:     row=col_idx[i, 0]
8:     vload B[row, :] ▷ load the selected row of B
9:     s0 = values[i, 0] ▷ move values[i, 0] to scalar reg
10:    C[i, :] += s0 * B[row, :] ▷ scalar-vector mul-acc
11:    values[i, :] >> 1 ▷ vector slide to the right
12:    col_idx[i, :] >> 1 ▷ vector slide to the right
13:   end for
14:   vstore C[i, :] ▷ store to mem the ith row of C
15: end for

```

III. OPTIMIZING SPARSE-DENSE MATRIX MULTIPLICATION

The implementation of sparse-dense multiplication eliminates unnecessary multiplications, due to the structured-sparse format of A . However, one crucial bottleneck of the computation is the abundance of *vector loads* from memory for the rows of matrix B , as shown in lines 8 and 9 of Algorithm 2. To tackle this issue, we leverage the structured sparsity of matrix A to reduce memory traffic and allow the computations to use *local* data that already reside in the vector register file.

The proposed optimization combines: (a) the pre-loading of tiles of matrix B in the vector register file, where they are kept stationary for as long as they are needed by the computations; and (b) a custom index-multiply-accumulate instruction that replaces the vector loads of matrix B with lower-cost indirect reads of the vector register file.

The key attribute that enables the pre-loading of tiles of matrix B in the register file is the well-defined, regular structure in the format of matrix A . Since the sparsity of A is – by construction – *structured*, the blocks within said matrix have a constant and known size. In turn, this implies that the column indexes of the non-zero element values in A are ‘bounded’ by the block size M , i.e., all `col_idx` values reside within the range $[0, M - 1]$. Recall that the block size M is the number of consecutive elements within a row of A that can contain up to a specific number (N) of non-zero elements. Exploiting this trait of structured sparsity, we may pre-load as many rows of matrix B as our vector register file can accommodate (with some restrictions, as will be explained shortly) and be sure that the column indexes in matrix A will only point to those local rows. On the contrary, with unstructured sparsity, the column indexes are, essentially, ‘unbounded’ and could point to any row of B (thereby rendering the pre-loading of specific rows of B in the register file futile).

In almost all practical cases, the sizes of matrices A and B are larger than the hardware-supported vector length. Therefore, matrix multiplication is executed in *tiles*. The tile size of matrix B that is pre-loaded in the vector register file is $L \times Vector_Length$; i.e., L rows, with each row having $Vector_Length$ columns. The number of rows L must be a multiple of M . Since a vector register can hold at

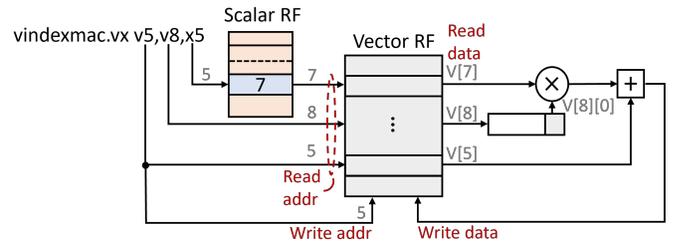


Fig. 3. The operation of the proposed `vindexmac` instruction. The contents of the scalar register are used to address a specific vector register. The vector read is multiplied with the least significant element of another vector register that is read in parallel. The result of the multiplication is accumulated with the previous contents of the vector destination register.

most $Vector_Length$ elements of a row of A , this means that, for $N:M$ structured sparsity, these elements refer to $Vector_Length/N$ blocks of this row, where each block includes M columns. Thus, the total number of columns of a row of a A and, effectively, the total number of rows of B that can be addressed is *at most* equal to $M \times Vector_Length/N$. This sets the upper limit on how many rows L of B makes sense to pre-load into the vector register file. If more rows of B are pre-loaded, then they would simply not be accessed. Of course, pre-loading fewer rows is possible, as long as their number is a multiple of M .

A. The proposed index-multiply-accumulate instruction

The purpose of the vector index-multiply-accumulate (`vindexmac`) instruction is to transform the vector load and multiply-accumulate operations found in lines 8–11 of Algorithm 2 into a new combined operation that would, instead, read the corresponding pre-loaded rows of matrix B directly from the vector register file and operate on them. The `vindexmac` instruction is defined as follows:

```

vindexmac.vx vd, vs2, rs
vd[i] += vs2[0] * vrf[rs|4:0][i]

```

where `vrf` refers to the vector register file. The instruction has three operands: one vector destination register (`vd`) and two source registers. One source register is scalar (`rs`) and the other source register is a vector (`vs2`). To execute `vindexmac`, the scalar register `rs` is read and its contents (only the 5 least significant bits are actually needed) are used as an address to the vector register file. The contents of the vector register that are read via the address contained in `rs` are multiplied with the least significant element of vector register `vs2` and accumulated with the contents of `vd`. The operation of `vindexmac` is visualized in Fig. 3.

The usage of the new `vindexmac` instruction to significantly speed up the sparse-dense matrix multiplication of Algorithm 2 is depicted in Algorithm 3, for one tile of matrix B . To complete the matrix multiplication, Algorithm 3 should be repeated for all tiles. As shown in lines 2–4, a tile of matrix B is pre-loaded into the vector register file. In lines 6 and 7, the non-zero values and column indexes of a row of matrix A are loaded also into the vector register file, similar to Algorithm 2. The essential differences between Algorithm 3 and the previous algorithms are observed in lines 10 and 11. In line 10, the value that is transferred to the scalar register is the column index needed by `vindexmac` instruction, which is executed in line 11 and it replaces the standard scalar-vector multiply-accumulate operation of the previous

algorithms. Since the tile of matrix B remains stationary in the vector register file, part of output matrix C is reloaded, updated for each non-zero element of matrix A , and stored back to memory in lines 8, 11, and 15, respectively.

Algorithm 3 Vector sparse-dense matrix multiplication with the use of the new `vindexmac` instruction for a tile of B

```

1: set vector length
2: for k=0 until k=L-1 do
3:   vload B[k,:] ▷ preload L rows of B to vector regs
4: end for
5: for i=0 until i=num_of_rows_of_A-1 do
6:   vload values[i,:] ▷ load i-th row of values of A
7:   vload col_idx[i,:] ▷ load i-th row of col. indexes of A
8:   vload C[i,:] ▷ Load the i-th row of C from mem
9:   for j=0 until j=nonzero_elems_per_block do
10:    s0 = col_idx[i, 0] ▷ move index to scalar reg
11:    vindexmac.vx C[i,:], values[i,:], s0
12:    values[i,:] >> 1 ▷ vector slide to the right
13:    col_idx[i,:] >> 1 ▷ vector slide to the right
14:   end for
15:   vstore C[i,:] ▷ Store to mem the i-th row of C
16: end for

```

Overall, the `vindexmac` instruction replaces the three instructions (one of which is a vector load from memory) of lines 8–10 of Algorithm 2 with the two instructions of lines 10 and 11 in Algorithm 3. The vector load from memory has now been eliminated. The `vindexmac` instruction itself calculates a vector of partial results that correspond to a row of result matrix C .

B. Hardware support for `vindexmac` execution

To execute the `vindexmac` instruction, we need to access both the scalar and the vector register files. This is an inherent attribute of all RISC-V scalar-vector instructions that have the letter ‘x’ in the suffix, e.g., `.vx`, `.vxm`, `.wx`, etc.

In this work, we target vector processors that follow a so called *decoupled* architecture, which consists of a scalar core that is responsible for instruction fetching and orchestration of execution, and a vector engine that executes the vector operations received from the scalar core [22]–[27]. Since, at any given time, there may be many vector instructions that require the values of scalar registers, the scalar core is responsible to transfer these values to the vector processor, together with the vector instructions themselves.

Therefore, the value of the scalar register `rs` required by `vindexmac` to address the vector register file is already provided by the scalar core. The given address drives one of the read ports of the vector register file, which outputs the requested vector operand. Another read port is used to read the elements of `vs2` and the third on reads `vd` as required by all multiply-and-accumulate scalar-vector instructions already present in RISC-V. Therefore, the only hardware requirement of the new `vindexmac` instruction is the addition of a multiplexer in front of the address bus of one of the read ports of the vector register file, which selects between `vs1` of normal vector arithmetic operations, or the 5 least significant bits of `rs` (as required by `vindexmac`). In other words, the new instruction by re-using the hardware infrastructure

TABLE I
SIMULATED PROCESSOR CONFIGURATION

Scalar core	<ul style="list-style-type: none"> RISC-V ISA (RV64GC), 8-way-issue out-of-order, 16-entry LSQ, 90 physical integer and 90 physical floating-point registers, 60-entry ROB L1I cache: 1-cycle hit latency, 4-way, 64KB L1D cache: 2-cycle hit latency, 4-way, 64KB
Vector engine	<ul style="list-style-type: none"> 512-bit vector engine with 16-lane configuration (32-bit elements \times 16 execution lanes) The vector engine is connected directly to the L2 cache through 16 store queues and 16 load queues
L2 cache	<ul style="list-style-type: none"> 8-way, 8-bank 8-cycle hit latency, 512KB shared by both the big core and the vector engine
Main Memory	DDR4-2400

of scalar-vector multiply-add instructions does *not* require an additional read port in the vector register file, but only a 5-bit 2-to-1 multiplexer in front of an existing read port.

Similarly, in fully integrated scalar-vector architectures [28], [29], the `vindexmac` instruction would be implemented in exactly the same manner as any of the other `.vx` instructions in a fully integrated scalar-vector setup. The only difference is that the scalar value provided would be used to drive one of the read ports of the vector register file rather than participating directly in computation.

Being a `.vx` instruction, `vindexmac` follows the standard encoding dictated by the RISC-V ISA for scalar-vector instructions [30]. However, in terms of the instruction functionality, there is a slight deviation from typical scalar-vector operations. The vector register identifier `vs2` is used only for its least significant element `vs2[0]`, i.e., it plays the role of the scalar value. The actual vector is fetched via an indirect read using the 5 least significant bits of `rs` as an address into the vector register file.

IV. EXPERIMENTAL EVALUATION

The experimental evaluation presented in this section aims to demonstrate the effectiveness of the new `vindexmac` instruction when executing state-of-the-art CNNs such as Resnet50 [31], DenseNet121 [32] and InceptionV3 [33]. The three CNNs were pruned to structured block sparsities of 1:4 and 2:4. This pruning and the appropriate fine-tuning (i.e., re-training) were performed using the TensorFlow library and the ImageNet dataset [34]. Our goal is to quantify the impact of transforming vector loads from memory into indirect reads from the vector register file – through the use of the new instruction – in terms of achieved speedups and reduction in the total number of memory accesses.

A. Experimental setup

For all experiments, we utilize a fully implemented decoupled vector unit connected to an out-of-order superscalar processor, i.e., model 1bDV in [24]. This design was modeled in the Gem5 simulator [18], [19] and the salient design parameters of the simulated processor setup are summarized in Table I. The new `vindexmac` instruction was incorporated in the RISC-V GNU toolchain and its operation was implemented in the decoupled vector engine in Gem5.

The convolutions of each layer of the examined CNNs are mapped to sparse-dense matrix multiplications $A \times B$ [5].

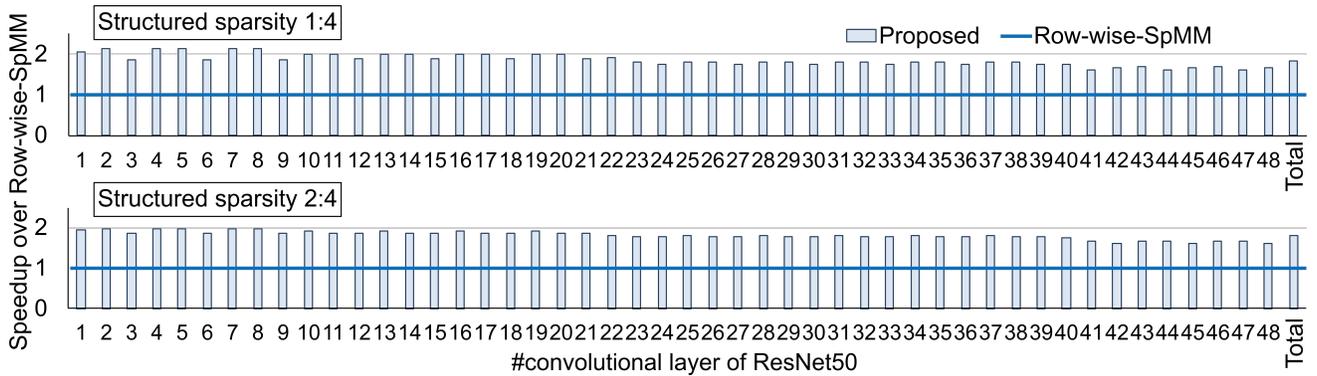


Fig. 4. The speedup achieved by the proposed approach in the *per-layer* execution times of ResNet50 [31], for (a) 1:4 and (b) 2:4 structured sparsity. The speedup is normalized to the execution time of ‘Row-Wise-SpMM’ in each respective layer.

Matrix A includes the structured-sparse weights and matrix B the input features of the corresponding CNN layers. Input features are considered dense, since there is no clear statistical attribute that can be exploited for them. Even if part of the input features contain zero values generated by the corresponding ReLU activation functions in each layer, their number is highly sensitive to the actual input values and filter weights.

The two designs under comparison are (a) ‘**Row-wise-SpMM**’: the simulated processor setup executing the row-wise sparse-dense matrix multiplications *without* the new instruction, i.e., executing Algorithm 2; (b) ‘**Proposed**’: the simulated processor setup executing the row-wise sparse-dense matrix multiplications using Algorithm 3, i.e., by *employing the new vindexmac instruction*.

To increase the performance of both algorithms, we applied loop unrolling, as proposed in [17], to produce four output results with the corresponding multiply-accumulate instructions in the same loop iteration. Both approaches benefit equally from loop unrolling.

Furthermore, as previously explained, the proposed approach follows a B -stationary dataflow to leverage the structured sparsity of A and the new `vindexmac` instruction. In all cases, we assume that the tile of B that is pre-loaded in the vector register file consists of $L=16$ rows. To ensure fairness in the comparisons, we tested *all three* dataflow types for ‘Row-Wise-SpMM,’ i.e., A -, B -, and C -stationary [17]. The experimental results show that the B -stationary dataflow (used by ‘Proposed’) also yields the best total execution times for ‘Row-Wise-SpMM.’ Therefore, all the experimental results assume a B -stationary dataflow for both approaches under comparison.

B. Evaluation results

The first set of results, shown in Fig. 4, refer to the execution latencies of each of the CNN layers of ResNet50 [31], for the two examined structured sparsity scenarios (1:4 and 2:4). The obtained execution times are normalized to the performance of ‘Row-wise-SpMM’. Under 1:4 sparsity, the proposed approach using the new `vindexmac` instruction achieves speedups in the range of $1.60\times$ – $2.15\times$ across all executed CNN layers. Similarly, under 2:4 sparsity, the obtained speedup is $1.63\times$ – $1.99\times$ across all layers.

Under both sparsities, the speedup tends to slightly decrease as the execution reaches the latter stages of the CNN. As the

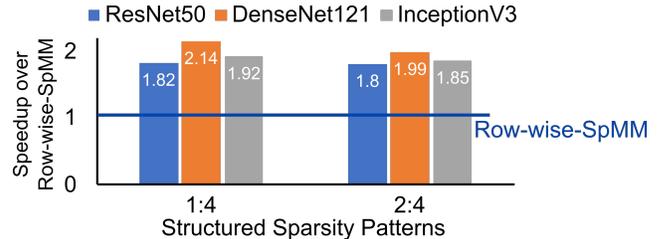


Fig. 5. The speedup achieved by the proposed approach in the *total* execution times of the three examined CNNs, for (a) 1:4 and (b) 2:4 structured sparsity. The results are normalized to ‘Row-Wise-SpMM’ of the respective sparsity.

stages progress, the matrix containing the input features (i.e., matrix B) becomes smaller, thus the benefit of pre-loading its rows in the vector register file becomes less pronounced. An interesting observation is that the speedup achieved by the proposed approach is slightly lower (across all layers) in the 2:4 sparsity scenario, as compared to 1:2. This is due to the fact that all the operations for matrix A are now twice as many as before, whereas the operations for matrix B remain the same. Thus, the operations for A now contribute substantially more to the total execution time and decrease the percentage contribution of the operations on B , which are the optimization target of `vindexmac`.

Similar behavior is observed in the *per-layer* execution times of the other two examined CNNs, DenseNet121 [32] and InceptionV3 [33]. Thus, those results are omitted for brevity.

The speedups achieved in each of the CNN layers translate to significant reductions in the *total* execution times of all three CNNs. Fig. 5 shows the achieved speedup of the proposed approach over ‘Row-Wise-SpMM’ in the three evaluated CNNs, for 1:4 and 2:4 sparsities. In both sets of bars, the results are normalized to ‘Row-Wise-SpMM’ of the respective sparsity. Clearly, the performance of ‘Proposed’ is substantially better in all cases. Across all three CNNs, the average speedup for 1:4 sparsity is $1.95\times$, while the average speedup for 2:4 sparsity is $1.88\times$.

The above-mentioned improvements in performance when using the ‘Proposed’ approach are reaped from the reduction of vector operations per iteration as well as the elimination of unnecessary vector loads from memory and their transformation into indirect reads from the vector register file. By pre-loading tiles of matrix B into the vector register file, the proposed approach exploits data locality very effectively,

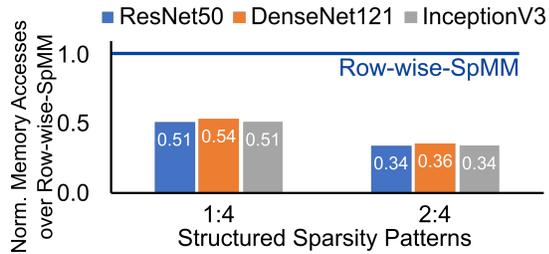


Fig. 6. The normalized number of *total* memory accesses observed when using the proposed approach in the three examined CNNs, for (a) 1:4 and (b) 2:4 structured sparsity. The normalization is with respect to ‘Row-Wise-SpMM’ of the respective sparsity.

thereby lowering the memory traffic. The results presented in Fig. 6 quantify the reduction in total memory accesses when using the proposed `vindexmac` instruction. The presented results are normalized to the number of memory accesses observed with ‘Row-Wise-SpMM’ of the respective sparsity. As can be seen, the total memory accesses decrease markedly. For sparsity 1:4, the memory accesses are reduced by 48%, on average, while the average reduction for 2:4 sparsity is 65%. Of course, if ‘Row-Wise-SpMM’ were to employ a *C*-stationary (instead of a *B*-stationary) dataflow, its total number of memory *stores* would decrease significantly. However, as explained at the end of Section IV-A, this reduction in *store* instructions does not improve the total execution time.

V. CONCLUSIONS

Vector processors can efficiently handle the abundant data-level parallelism available in modern ML applications. The scalability of ML models calls for appropriate model pruning that reduces their memory footprint and makes them amenable to applications at the edge. In this context, we aimed to seamlessly integrate the simplicity of structured sparsity with mainstream vector architectures enhanced with a proposed new instruction. The `vindexmac` instruction operates on local data that is pre-loaded into the vector register file, thereby reducing the number of instructions executed and eliminating unnecessary vector loads. Most importantly, the new instruction can be implemented with negligible hardware cost. The evaluation results demonstrate substantial speedups in the execution latency of representative layers of state-of-the-art CNN applications.

REFERENCES

- [1] D. Dabbelt *et al.*, “Vector processors for energy-efficient embedded systems,” in *ACM Intern. Workshop on Many-core Embedded Systems*, 2016, pp. 10–16.
- [2] G. Jeong *et al.*, “Vegeta: Vertically-integrated extensions for sparse/dense gemm tile acceleration on cpus,” in *IEEE Int. Symp. on High-Performance Comp. Arch. (HPCA)*, Feb. 2023, pp. 259–272.
- [3] T. Hoefler, D. Alistarh, T. Ben-Nun, N. Dryden, and A. Peste, “Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks,” *The Journal of Machine Learning Research*, vol. 22, no. 1, pp. 10 882–11 005, 2021.
- [4] U. Evci, T. Gale, J. Menick, P. S. Castro, and E. Elsen, “Rigging the lottery: Making all tickets winners,” in *Inter. Conf. on Machine Learning*, Jul. 2020, pp. 2943–2952.
- [5] A. Mishra *et al.*, “Accelerating sparse deep neural networks,” *arXiv preprint arXiv:2104.08378*, 2021.
- [6] A. Zhou, Y. Ma, J. Zhu, J. Liu, Z. Zhang, K. Yuan, W. Sun, and H. Li, “Learning N:M fine-grained structured sparse neural networks from scratch,” in *Inter. Conf. on Learning Representations*, May 2021.

- [7] Z.-G. Liu, P. N. Whatmough, Y. Zhu, and M. Mattina, “S2TA: Exploiting structured sparsity for energy-efficient mobile CNN acceleration,” in *IEEE Inter. Symp. on High-Performance Comp. Arch. (HPCA)*, Apr. 2022, pp. 573–586.
- [8] M. Zhu, T. Zhang, Z. Gu, and Y. Xie, “Sparse tensor core: Algorithm and hardware co-design for vector-wise sparse neural networks on modern GPUs,” in *IEEE Int. Symp. on Microarchitecture*, 2019, pp. 359–371.
- [9] J. R. Gilbert, C. Moler, and R. Schreiber, “Sparse matrices in MATLAB: Design and implementation,” *SIAM journal on matrix analysis and applications*, vol. 13, no. 1, pp. 333–356, 1992.
- [10] N. Bell, S. Dalton, and L. N. Olson, “Exposing fine-grained parallelism in algebraic multigrid methods,” *SIAM Journal on Scientific Computing*, vol. 34, no. 4, pp. C123–C152, 2012.
- [11] J. Li, F. Wang, T. Araki, and J. Qiu, “Generalized sparse matrix-matrix multiplication for vector engines and graph applications,” in *IEEE Workshop on Memory Centric High Perf. Comp. (MCHPC)*, 2019, pp. 33–42.
- [12] V. Le Fèvre and M. Casas, “Efficient execution of SpGEMM on long vector architectures,” in *Intern. Symp. on High-Performance Parallel and Distributed Computing (HPDC)*, 2023, p. 101–113.
- [13] J. Pavon *et al.*, “VIA: A smart scratchpad for vector units with application to sparse matrix computations,” in *IEEE Intern. Symp. on High-Performance Computer Archit. (HPCA)*, 2021, pp. 921–934.
- [14] N. Srivastava, H. Jin, J. Liu, D. Albonese, and Z. Zhang, “Matraptor: A sparse-sparse matrix multiplication accelerator based on row-wise product,” in *IEEE Int. Symp. on Microarchitecture*, 2020, pp. 766–780.
- [15] Z. Liu, M. Mattina, P. Whatmough, and J. G. BEU, “Processor for sparse matrix computation,” US Patent 11,392,376, 2022.
- [16] D. A. Iliescu and F. Petrogalli, “ARM scalable vector extension and application to machine learning,” ARM whitepaper, Tech. Rep., 2017.
- [17] G. Alaejos *et al.*, “Micro-kernels for portable and efficient matrix multiplication in deep learning,” *The Journal of Supercomputing*, vol. 79, no. 7, pp. 8124–8147, 2023.
- [18] N. Binkert *et al.*, “The gem5 simulator,” *SIGARCH Comput. Archit. News*, vol. 39, no. 2, Aug. 2011.
- [19] J. Lowe-Power *et al.*, “The gem5 simulator: Version 20.0+,” arXiv cs.AR 2007.03152, 2020.
- [20] A. Morad, L. Yavits, and R. Ginosar, “Efficient dense and sparse matrix multiplication on gp-simd,” in *Intern. Workshop on Power and Timing Modeling, Optimization and Simulation (PATMOS)*, 2014.
- [21] F. G. Gustavson, “Two fast algorithms for sparse matrices: Multiplication and permuted transposition,” *ACM Trans. on Mathematical Software (TOMS)*, vol. 4, no. 3, pp. 250–269, 1978.
- [22] R. Espasa *et al.*, “Tarantula: a vector extension to the alpha architecture,” in *Intern. Sympos. on Comp. Arch. (ISCA)*, May 2002, pp. 281–292.
- [23] C. Chen *et al.*, “Xuantie-910: A commercial multi-core 12-stage pipeline out-of-order 64-bit high performance RISC-V processor with vector extension : Industrial product,” in *ACM/IEEE Intern. Symp. on Comp. Arch. (ISCA)*, 2020, pp. 52–64.
- [24] T. Ta *et al.*, “big.VLITTLE: On-demand data-parallel acceleration for mobile systems on chip,” in *IEEE Intern. Symp. on Microarch. (MICRO)*, 2022, pp. 181–198.
- [25] 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 Trans. on VLSI Systems*, vol. 28, no. 2, pp. 530–543, feb 2020.
- [26] K. Patsidis, C. Nicopoulos, G. C. Sirakoulis, and G. Dimitrakopoulos, “RISC-V²: A Scalable RISC-V Vector Processor,” in *IEEE Intern. Symp. on Circuits and Systems (ISCAS)*, 2020.
- [27] F. Minervini *et al.*, “Vitruvius+: an area-efficient RISC-V decoupled vector coprocessor for high performance computing applications,” *ACM Trans. on Arch. and Code Optim. (TACO)*, vol. 20, no. 2, 2023.
- [28] J. Rupley, B. Burgess, B. Grayson, and G. D. Zuraski, “Samsung m3 processor,” *IEEE Micro*, vol. 39, no. 2, pp. 37–44, 2019.
- [29] N. Stephens *et al.*, “The ARM scalable vector extension,” *IEEE Micro*, vol. 37, no. 2, pp. 26–39, 2017.
- [30] “RISC-V vector extension (Version 1.0),” <https://github.com/riscv/riscv-v-spec>, accessed: 09-14-2023.
- [31] K. He *et al.*, “Deep residual learning for image recognition,” in *IEEE Conf. on Comp. Vision and Pattern Recogn. (CVPR)*, 2016, pp. 770–778.
- [32] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *IEEE Conf. on Comp. Vision and Pattern Recogn. (CVPR)*, 2017, pp. 4700–4708.
- [33] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *IEEE Conf. on Comp. Vision and Pattern Recogn. (CVPR)*, 2016, pp. 2818–2826.
- [34] O. Russakovsky *et al.*, “ImageNet Large Scale Visual Recognition Challenge,” *Intern. Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.