# Full-Stack Optimization for CAM-Only DNN Inference

João Paulo C. de Lima[*†‡], Asif Ali Khan[*], Luigi Carro[‡], Jeronimo Castrillon[*†]

[*]Chair for Compiler Construction, Technische Universität Dresden, Dresden, Germany
[†]Center for Scalable Data Analytics and Artificial Intelligence (ScaDS.AI), Dresden, Germany
[‡]Informatics Institute, Federal University of Rio Grande do Sul, Porto Alegre, Brazil
Email: {joao.lima, asif_ali.khan, jeronimo.castrillon}@tu-dresden.de, carro@inf.ufrgs.br

*Abstract*—The accuracy of neural networks has greatly improved across various domains over the past years. Their ever-increasing complexity, however, leads to prohibitively high energy demands and latency in von-Neumann systems. Several *computing-in-memory* (CIM) systems have recently been proposed to overcome this, but trade-offs involving accuracy, hardware reliability, and scalability for large models remain a challenge. Additionally, for some CIM designs, the activation movement still requires considerable time and energy. This paper explores the combination of algorithmic optimizations for ternary weight neural networks and *associative processors* (APs) implemented using *racetrack memory* (RTM). We propose a novel compilation flow to optimize convolutions on APs by reducing their arithmetic intensity. By leveraging the benefits of RTM-based APs, this approach substantially reduces data transfers within the memory while addressing accuracy, energy efficiency, and reliability concerns. Concretely, our solution improves the energy efficiency of ResNet-18 inference on ImageNet by 7.5× compared to crossbar in-memory accelerators while retaining software accuracy.

*Index Terms*—Associative memory; racetrack memory; neural network, compiler optimizations;

## I. Introduction

Neural network (NN) models have experienced growth, leading to enhanced accuracy and wider applicability across diverse domains. However, this progress has also brought forth computational challenges due to the substantial data movement between memory and processing units. To mitigate this, algorithmic approaches for efficient NN acceleration, such as pruning, extreme quantization, and optimized convolution algorithms like Winograd, can reduce memory requirements and computational costs, while, in the architectural front, the *computing-in-memory* (CIM) paradigm exploits the physical attributes of the memory cells to compute in place [1].

Matrix multiplication is a crucial operation in NNs, and CIM-based accelerators often use resistive crossbar arrays for analog/mixed-signal matrix-vector multiplication (MVM) in constant time, achieving significant performance gains over digital CMOS methods. However, their widespread usage is hindered by reliability issues and energy-hungry digital-to-analog (DAC) and analog-to-digital (ADC) conversions. Despite numerous advances in ADCs/DACs [1], crossbars still exhibit high inaccuracy, are limited to small-scale use cases, and face reliability issues. As alternative to crossbars, promising CIM paradigms like *content-addressable memories* (CAMs) have also been explored [2, 3, 4]. Despite their excellent

performance and no need for costly peripherals, current CAM designs struggle to maintain accuracy for larger networks as they rely on methods like binarized networks (BNNs) [2], approximate MACs [5], or dot-product approximations [3, 4].

Recent research has shown that ternary weights networks (TWNs), i.e., $w_i \in \{-1, 0, 1\}$, combined with reduced-precision activations (often integers with 8 to 4 bits) can retain accuracy while drastically reducing computational complexity on specialized hardware [6, 7]. In bulk-bitwise CIM designs, if executed intelligently, such ternary models would rely solely on $O(m)$ addition/subtraction, eliminating the need for $O(m^2)$ multiplication, where $m$ is the bitwidth of operands. However, existing CIM designs prioritize efficient handling of weights, disregarding the significant cost of moving activations. The cost of transferring activations can become significant, depending on the activation precision, possibly outweighing the cost of weights. In modern NNs, this movement can consume 30% to 70% of the total CIM system energy [8].

In this paper, we present a full-stack solution, comprising a compilation flow and CAM-based accelerator, that considers accuracy as a first-class optimization metric, along with performance and energy consumption. By consolidating previously scattered methods, we investigate their joint impact and propose a sequence of steps to co-optimize bulk-bitwise convolutions for performance and to reduce activation transfer. Specifically, we employ *racetrack memory* (RTM)-based CAMs to implement *associative processors* (APs) for DNN inference. Compared to other NVMs, RTMs [9] provide an order of magnitude higher density, comparable performance, and improved energy efficiency and endurance. Concretely, we make the following novel contributions:

- An RTM-based AP for efficient bulk-bitwise processing. By sequentially storing feature maps in RTM cells, we effectively leverage the sequential accesses of RTMs for bit-serial word-parallel convolution in the AP model.
- A compilation flow for the proposed AP with optimizations to reduce the arithmetic intensity of convolution and transform weights into AP instructions. Our transformation focuses on eliminating multiplications, removing redundant additions/subtractions and optimizing the bitwidth of partial sums.
- We evaluate our full-stack solution on various NN architectures - VGG9, VGG11, and ResNet18, using CIFAR10

and ImageNet datasets. Ensuring software accuracy, our RTM-APs accelerate inference for the largest model (ResNet-18/ImageNet) by $3\times$ with $2.5\times$ lower energy consumption compared to crossbar-based accelerators, resulting in $7.5\times$ energy efficiency improvement.

## II. BACKGROUND

This section provides a concise overview of the RTM technology, CAMs and APs, as well as neural networks.

### A. DNNs: Quantization and sparsity

Convolutional layers, depicted in Fig. 1, are the most computationally intensive in modern DNNs. They take multiple $C_{in}$ input feature maps (IFMs) and convolve them with $C_{out}$ filters to generate $C_{out}$ output feature maps (OFMs).
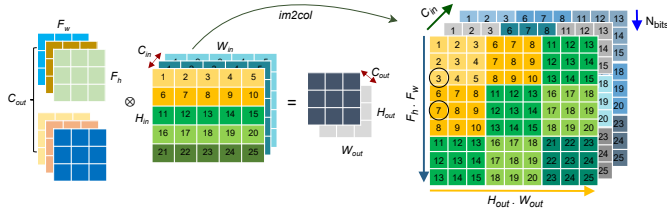


Fig. 1: Direct convolution and im2col transformation

Most DNNs are overparameterized and contain several sub-networks with comparable accuracy and substantially less memory and computational intensity. This concept backs popular methods like pruning and quantization. The recently proposed multi-prize lottery ticket hypothesis [7] also supports the existence of an optimal binary sub-network inside a super-network. This finding establishes a necessary condition for accurately implementing NN hardware using low-cost bit-wise operations, like the XNOR-popcount implementation of BNNs [2]. However, the activation requirements tell a different story: binary activations strongly limit the representational capacity of the network, even when taking batch normalization into account [7]. As a consequence, this results in an unavoidable loss in accuracy. Recent research shows that a moderate activation quantization, such as 4-bit, is sufficient to maintain the same accuracy as its full-precision counterpart, even in complex classification problems [6].

### B. Content addressable memories and associative processing

CAMs represent a class of CIM that enables fast and energy-efficient memory search operations. CAMs take an input pattern (search key), compare it with all the contents of the CAM array, and return locations where the contents match. For the search operation, all match lines are precharged to a high voltage and the key is applied to the CAM input lines. If a single cell mismatch occurs, the corresponding match line discharges. Otherwise, the match line will stay high when all bits are matched. The match lines are connected to a sensing circuitry to decode a corresponding address. This way, parallel searches are performed on all CAM content in $O(1)$.
**Associative processors (APs):** In various studies, CAMs have been used for exact and approximate search [3, 4], as well as

TABLE I: LUT tables for both in-place and out-of-place 1-bit addition (left) and subtraction (right).

| 1b Adder | | | In-place | | | Out-of-place | | | 1b Sub | | | In-place | | | Out-of-place | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Cr | B | A | Cr | B | Comment | Cr | R | Comment | Br | B | A | Br | B | Comment | Br | R | Comment |
| 0 | 0 | 0 | 0 | 0 | NC | 0 | 0 | NC | 0 | 0 | 0 | 0 | 0 | NC | 0 | 0 | NC |
| 0 | 0 | 1 | 0 | 1 | 2nd | 0 | 1 | 1st | 0 | 0 | 1 | 1 | 1 | 1st | 1 | 1 | 1st |
| 0 | 1 | 0 | 0 | 1 | NC | 0 | 1 | 2nd | 0 | 1 | 0 | 0 | 1 | NC | 0 | 1 | 2nd |
| 0 | 1 | 1 | 1 | 0 | 1st | 1 | 0 | NC | 0 | 1 | 1 | 0 | 0 | 2nd | 0 | 0 | NC |
| 1 | 0 | 0 | 0 | 1 | 3rd | 0 | 1 | 3rd | 1 | 0 | 0 | 1 | 1 | 4th | 1 | 1 | 3rd |
| 1 | 0 | 1 | 1 | 0 | NC | 1 | 0 | NC | 1 | 0 | 1 | 1 | 0 | NC | 1 | 0 | NC |
| 1 | 1 | 0 | 1 | 0 | 4th | 1 | 0 | 4th | 1 | 1 | 0 | 0 | 0 | 3rd | 0 | 0 | 4th |
| 1 | 1 | 1 | 1 | 1 | NC | 1 | 1 | 5th | 1 | 1 | 1 | 1 | 1 | NC | 1 | 1 | 5th |

for bulk-bitwise arithmetic and logic operations, also known as APs [10]. APs perform in-place computation on stored inputs using CAM search/write, specifically masked search and parallel column write operations. In this execution model, the truth table of a Boolean expression is broken down into simpler search and write primitives, generating a lookup table (LUT) of masked keys to implement the required function. Fig. 2c illustrates an AP consisting of a CAM array, which allows searching a key in the stored content and writing a data pattern (write key) in all tagged rows (see Sec. III). Depending on the operation (search/write), the AP controller sets the mask and key values iteratively by referencing the corresponding LUT (Table I). In the *search phase*, the search key and mask fields are set according to the left side of the LUT and compared with the CAM content in the specified/masked columns and all rows in parallel. The matched lines are stored in the *tag register*, which drives updating the data pattern in all tagged rows in the next phase. In the *write phase*, the mask and key values are set by observing the LUT's right side, and the CAM content in the tagged rows and selected columns is written accordingly. Thanks to the parallel search across all rows, any function performed on a sequential processor can be implemented on the AP in a SIMD fashion.

### C. Racetrack memory

RTM is a magnetic NVM technology where a single cell comprises a magnetic nanowire, also referred to as a track, with the capacity to store up to 100 data bits (domains) (see Fig. 2e) [11]. Each track has one or more access ports that enable the read/write operations. The domain walls must be first shifted and aligned with the access port to access a specific domain within a track. Typically, multiple tracks are grouped into domain-wall block clusters (DBCs) and can be accessed in parallel. While RTMs offer higher storage density and lower power operations than other memory technologies [9], their read/write speed can vary with different access patterns due to the movement of magnetic domain walls in nanowires. For this reason, data in RTM is typically stored in a bit-interleaved fashion to reduce the number of shifts and sequential accesses.

## III. RTM-AP: ACCELERATOR ARCHITECTURE

As illustrated in Fig. 2a-c, the accelerator is organized into a three-level hierarchy comprising banks, tiles, and APs. Each AP is an independent processing unit operated and accessed in parallel. The hierarchical organization enables fine-grained control over the compute resources by allocating banks, tiles, and APs based on the computational needs of each layer.
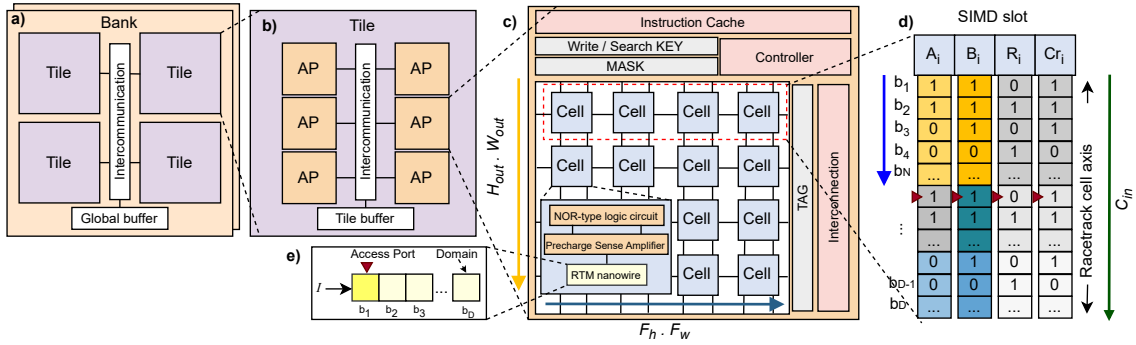
Fig. 2: RTM-AP. a) Hierarchical accelerator architecture consisting of banks, tiles, APs, buffers and interconnection network, b) tile showing array of APs, c) AP organization consisting of CAM array, registers, instruction cache and control unit, d) SIMD slot for two-operand addition and mapping of inputs to racetrack domains, e) an RTM nanowire.

The bitwise implementation of addition or subtraction can be in-place or out-of-place, meaning that the result is written into one of the inputs or a new location. Table I illustrates the LUTs for both versions. The run order indicated in the comment column in Table I must be followed to ensure the correctness of bitwise computations. Note that keys marked as NC (no change) do not alter any content in the CAM, hence no search/write is needed. In all tables, $R$ and $C_r/B_r$ represent the result and carry/borrow, respectively, and A and B are the inputs. At a coarser level, the AP instruction cache drives the execution of application kernels by providing the sequence of instructions, operand locations, and their bitwidth.

Fig. 2c shows the CAM cell incorporating a NOR-type logic circuit, a precharge sense amplifier, and a memory device (in this case, an RTM nanowire). APs support only bitwise operations and require storing the input so that each row stores bits from different operands at the same bit position. In prior research [10], each row also stores all operands' bits across multiple columns due to limitations in their CAM design. In the case of NVM technologies like RRAM, where the resistance window is divided into ranges to represent different bits, their multi-level cell operation is incompatible with bitwise operations. Consequently, they must operate in a single bit per cell mode, significantly losing capacity.

However, RTMs can effectively provide the mechanism needed for AP implementation. Earlier research has successfully implemented ultra-dense, low-power CAMs using RTM and skyrmion devices [11, 12, 13] based on manipulation and movement of domain walls. Since operands in APs are processed in a bit-serial fashion, in this work, we store the operands sequentially in the nanowires. In the example in Fig. 2d, each column $(A, B, R, C_r)$ represents an operand, and the row value "1101" is the readout bits given by the state of the current domains aligned to access ports. Unlike prior research that employed interleaving to work around RTMs' sequential nature [9], our novel execution model harnesses this characteristic for true multi-bit storage and bitwise processing.

## IV. COMPILATION FRAMEWORK FOR RTM-APs

In bulk bitwise CIM architectures, we need to prioritize simpler and fewer operations, while other aspects important for von-Neumann systems, such as code size and temporal locality, may have less significance. We propose an automatic compilation flow that optimizes convolutions at the application and arithmetic levels. This is achieved by automatically reducing the arithmetic intensity of operations and by efficiently mapping them to reduce data movement. Fig. 3a depicts a high-level overview of our framework that transforms trained TWNs into an optimized sequence of AP instructions. In contrast to conventional convolution, where weights are fetched from memory and multiplied with inputs, our approach involves statically compiling weights into AP instructions. Additionally, our implementation stands apart from other CIM approaches that pass feature maps through multiple memory arrays, incurring high costs for peripheral circuits. Instead, we adopt a data-centric approach, where feature maps are primarily computed in place, significantly reducing data movement.

### A. Data-flow graph generation

Consider the standard loop nest of convolutional layers shown in Fig. 3b. In inference tasks, $W$ can usually be known and fixed at compile time. Assuming trained TWNs, we can statically replace multiplications with expressions involving only additions and subtractions. Depending on the model's sparsity, this can also eliminate many expressions involving multiplication with zero. To enable this optimization, we unroll loops $kh$ and $kw$ and apply constant weight folding, where weights limited to $\{-1, 0, 1\}$ replace multiplications. This optimization effectively reduces computational complexity by eliminating multiplication from the convolution kernel. However, it comes at the cost of fully unrolling all loops, leading to a substantial increase in code size overhead.

To further reduce arithmetic complexity and code size, we must expose the loop body with the highest redundancy in additions before unrolling it. This region of high redundancy refers to the weight slices convolved on the same input patch. This is achieved by applying loop interchange on the naïve loop (Fig. 3b) to make $ofm$ the third innermost loop, followed by unrolling the three innermost loops—namely, the variables $ofm$, $kh$, and $kw$. After applying constant weight folding, the resultant loop has a larger body but a higher number
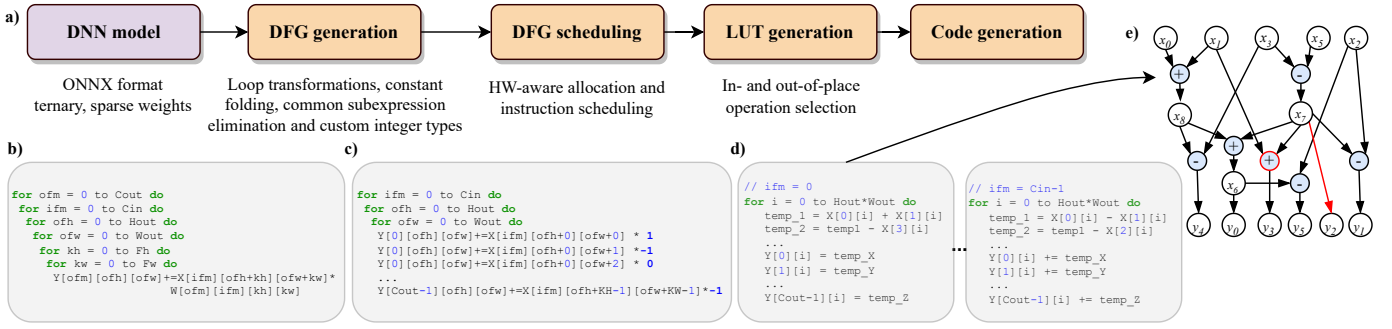
Fig. 3: a) Compilation flow and optimization techniques used in each step, b) naïve loop in convolutional layers, c) loop after applying loop interchange, unrolling and constant folding of ternary weights, d) loop after loop fission and common subexpression elimination, e) optimized data-flow graph (DFG) for Equation 1.

of subexpressions depending only on an input patch of size $F_h \times F_w$, as illustrated in Fig. 3c. Loop interchange is an enabler for a later optimization step, as using it alone is not beneficial due to poor locality. Initially, all three steps might appear counterintuitive, especially considering that real-world networks often have $C_{out}$ ranging from 64 to 512 and their overhead in code size and poor locality in CPUs. However, these transformations target in-memory architectures, which differ significantly from von-Neumann ones, and existing compilers would not naturally explore these optimization steps.

The outermost loop from Fig. 3c is also fully unrolled, resulting in $C_{in}$ loop bodies depicted in Fig. 3d. Each of these bodies handles a single IFM and iterates over $ofh$ and $ofw$ variables, which are then simplified in a single loop over $H_{out} * W_{out}$. They can be further optimized by Common Subexpression Elimination (CSE), reducing the number of additions. CSE aims to identify and eliminate common expressions within a program, reducing the total number of computations required and improving the program's overall performance. In this approach, CSEs are found within the weight slice ($C_{out} \times 1 \times F_h \times F_w$), which is the slice that allows the greatest potential for reuse of a single input-channel across all output channels. Equation 1 illustrates a simple implementation of MVM with ternary weights.

$$y = \begin{pmatrix} 1 & -1 & 0 & 1 & 0 & -1 \\ 0 & 0 & -1 & 1 & 0 & -1 \\ 0 & 0 & 0 & -1 & 0 & 1 \\ 0 & -1 & 0 & -1 & 0 & 1 \\ 1 & -1 & 0 & -1 & 0 & 0 \\ 1 & -1 & -1 & 1 & 0 & -1 \end{pmatrix} * \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} x_6 \\ x_7 - x_2 \\ -x_7 \\ -x_1 - x_7 \\ x_8 - x_3 \\ x_6 - x_2 \end{pmatrix} \quad (1)$$

with $x_6 = x_7 + x_8$, $x_7 = x_3 - x_5$, $x_8 = x_0 + x_1$. The MVM operation in Eq. 1 originally involves 19 operations and can be reduced to 7 when removing redundant expressions. The Data-flow graph (DFG) for optimized computation is shown in Fig. 3e. The red operator and arrow are similar operations to Fig. I but combined with negative output. After applying loop fission, explicit array access, and CSE, the final unrolled code consists of IFM loops, as illustrated in Fig. 3d.

To enable vectorization, a sliding window must be stored as a column, achieved by expanding the input channel into a column matrix using the im2col operation (Fig. 2b). Thus, each iteration of loops in Fig. 3d generates $C_{out}$ partial dot-products and can be vectorized based on the CAM size.

Finally, the flow annotates the minimum required bitwidth for each level of the DFG. As AP supports arbitrary integer types, we leverage narrow-precision data types for faster addition and reduced energy consumption without affecting correctness.

### B. Input mapping and DFG scheduling

Using arrows, Fig. 2c shows the dimensions of a typical input mapped to a CAM and its multi-bit RTM cells shown in Fig. 2d. $F_h * F_w$ is distributed along the CAM columns, while $H_{out} * W_{out}$ is distributed along the CAM rows. Multiple APs can be used to meet the requirements of each layer so that all points in the $H_{out} * W_{out}$ dimension are computed in parallel. Fig. 2d shows how $N$ bits of the input are stored in a single RTM cell and $C_{in}$ channels can be stored contiguously, sharing the same nanowire and favoring density. Given the limited number of bits per cell, multiple CAMs are needed to hold a realistic number of channels, thus adding parallelism.

Finally, instructions are scheduled in two phases: the *channel-wise* DFG, where $C_{in}$ partial OFMs are generated, and *accumulation phase*, where $C_{in}$ partial results are accumulated to form the output channel. In the first step, the operands in AP instructions (i.e., CAM columns) are treated as registers, similar to the register allocation problem, and solved through graph coloring. In the second step, we first accumulate partial OFMs locally in the same APs, and we then move partial results from one AP to another in a regular adder tree fashion to favor the parallelism of multiple CAMs. The last step also fuses the activation function and stores the OFMs to facilitate the computation of the next layer.

### C. Lookup table generation

We have two options for adder/subtractor implementation: a) in-place, requiring eight cycles, and b) out-of-place, requiring ten cycles. Given that the number of cycles/passes needed for in-place operations is smaller than for out-of-place, we maximize the use of in-place operations. This is done by exploiting the liveliness of a variable in a DFG or by making copies of the result so that in-place operations can be used in the future. Fig. 3e shows an example where the operands $x_6, x_7$ and $x_8$ are used more than once. Thus, their defining operations must store their results in as many columns as they

are used by writing in the selected columns in parallel, allowing the subsequent operations to be in place. Additionally, LUTs for generating negative output are provided at the same cost as the standard ones.

## V. EXPERIMENTAL SETUP AND EVALUATION RESULTS

As baseline RTM TCAM, we use the $45\,\text{nm}$ $256{\times}256$ design in [12]. It features a search delay under $200\,\text{ps}$ and per-bit search energy of around $3\,\text{fJ}$. We assume 64 domains per nanowire based on [9]. We built a functional simulator that models the architecture presented in Fig. 2 and estimates performance and energy consumption based on the figure of merits of this CAM design. The functional simulation generates outputs that are analyzed to assess accuracy. We consider a conservative $1\,\text{pJ/bit}$ energy consumption for internal data movement at the tile, bank, and global level [14].

**Benchmarks:** We evaluate the ImageNet and CIFAR10 datasets on ResNet18 and VGG-9/VGG-11 models, respectively, which were trained with BIPROP [7] for 100 epochs to identify TWNs that achieve accuracy comparable to a dense full-precision network. We use learned step quantization (LSQ) [6] for quantizing activations to 8 and 4 bits.

**Evaluated configurations:** As a comparison baseline, we use an RRAM-based CIM design simulated on DNN+NeuroSim [14], comprising 8-bit weights, $256{\times}256$ arrays and 5-bit ADCs. To quantify the individual effect of the proposed transformations on RTM-AP, we evaluate two configurations: *unroll*, implementing loop unrolling, constant weight folding, and custom integer types optimizations, and *unroll+CSE*, including all optimizations from Fig. 3a.

### A. Results summary and comparison to state-of-the-art

Table II summarizes the impact of the sparsity and activation precision on the accuracy, energy consumption, and latency per inference of our proposed AP compared to DNN+NeuroSim [14] (crossbar-based) and DeepCAM [4] (CAM-based). Notably, the 4-bit activation maintains software accuracy while achieving optimal inference performance and energy efficiency. Compared to [14], our *unroll+CSE* reduces the per inference latency and energy consumption by up to ($\sim3\times$, $2.5\times$) and ($1.4\times$, $1.8\times$) for ResNet18 and VGG-9, respectively. The improvements in performance and energy are partially due to the reduced search/write latency and energy consumption of the RTM-CAMs but are largely attributed to our compiler transformations. For instance, the CSE optimization alone reduces the number of additions by an average of 31%. Despite outperforming RTM-AP, DeepCAM faces scalability issues, rendering this comparison unfair. DeepCAM depends on large arrays, up to 512 rows and 1024 columns, and the energy efficiency of deeper networks like ResNet18 does not scale as effectively as with LeNet and small VGGs. Also, the accuracy of complex tasks, like ImageNet, is more sensitive to approximation. A thorough comparison requires further investigation into the necessary peripheral circuits required by DeepCAM to convert the timing of the match line discharge into digital Hamming distance values, which is beyond the scope of this work.

### B. Impact on performance and energy consumption

This section provides a detailed layer-by-layer comparison of the *unroll* and *unroll+CSE* configurations to the baseline DNN+NeuroSim for a more fine-grained analysis. We only focus on ResNet18 (the largest model) for space reasons. Fig. 4 presents the energy consumption and latency of the three configurations, considering the contributions of the individual system components. Our solution (*unroll+CSE*) effectively reduces the number of additions (from 1,499K to 931K), thus significantly improving the total energy consumption and latency per image. The reductions are more significant in the first layer because larger kernels (i.e., $7{\times}7$) allow for more subexpression elimination. Also, one can notice that deeper into the network, latency tends to increase due to reduced CAM row utilization as $H_{out} * W_{out}$ decreases. Therefore, the layers 16-20 are slower on RTM-AP compared to [14], but still more energy-efficient. This could be alleviated, however, by processing multiple images per layer.

Compared to the crossbar-based baseline, *unroll+CSE* can significantly improve end-to-end energy consumption and latency while retaining accuracy. It is worth mentioning that the DNN+NeuroSim model of ResNet-18 does not count with a precise model for weights, activations and ADC quantization, while our modeling requires only activation quantization. In DNN+NeuroSim[14], *peripherals* account for buffers, digital logic modules (e.g., decoder, switch matrix, mux), and interconnect, while *accumulation* accounts for shift and adders and accumulation units at subarray, tile and global levels. Since our design relies on additional accumulation units, the *accumulation* label denotes the AP energy for computing the *accumulation phase*, while the energy spent in the channel-wise DFG phase is indicated as *DFG*.

### C. Impact on data movement and write endurance

Reducing data movement of partial sums and activations is challenging in current CIM accelerators. We solve this problem by storing activations in CAMs and computing them exclusively with AP operations. Data movement of partial results is not eliminated but it accounts for only 3% of the total energy consumption. Most of the computation is done locally in each AP, and communication is required only during the *accumulation phase*. In contrast, in crossbar-based design, communication accounts for 41% of the total energy [14].

As for write endurance, RTM [9] offers $10^{16}$ write cycles (best compared to other NVMs). In the worst-case scenario, at most two columns are written only once for each in-place or out-of-place operation (Fig. I), which takes 0.8 or 1 ns, respectively. Given that the execution flow is distributed across 256 columns, it is reasonable to rewrite the same column only after approximately 128 operations. This implies writing to the same location roughly every 100 ns, on average, hence resulting in an estimated lifespan of $\sim31$ years.

TABLE II: Accuracy, energy consumption, latency, required memory and number of operations comparison.

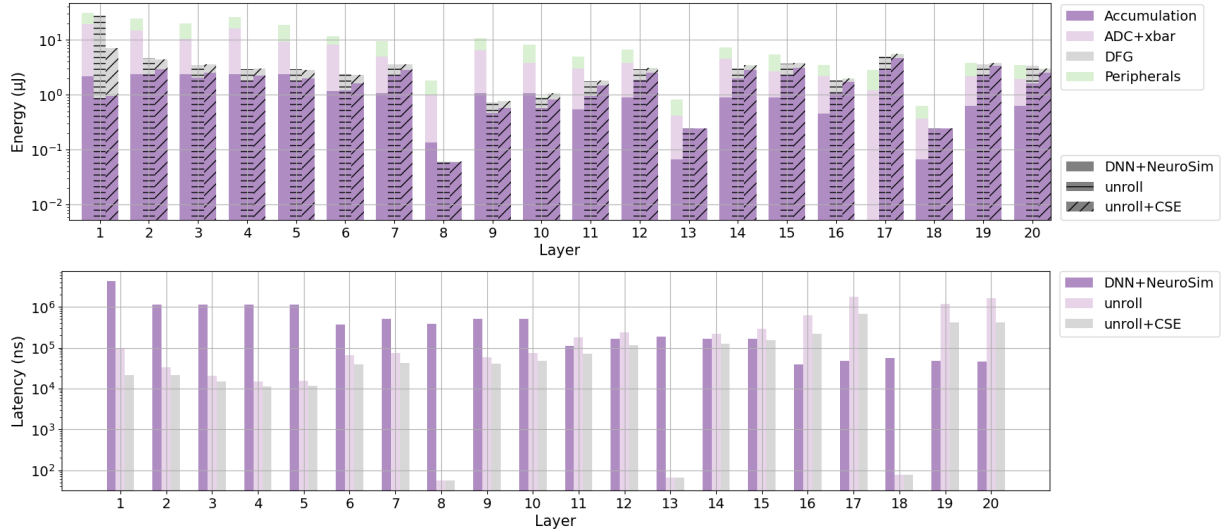| Network / Dataset | Sparsity Actv. prec. | Top-1 Accuracy (%) | | | Energy / Inference ($\mu$J) | | Latency (ms) | | # Arrays $256 \times 256$ | #Adds/Subs ($1e^3$) | |
| | | FP | 4-bit | 8-bit | 4-bit | 8-bit | 4-bit | 8-bit | | unroll | unroll+CSE |
| ResNet18/ImageNet | 0.8 | 70.5 | 70.6 | 70.6 | 55.04 | 78.56 | 2.46 | 4.10 | 49 | 1499 | 931 |
| ResNet18/ImageNet [14] | - | 70.5 | 70.0 | 70.0 | 104.92 | 199.90 | 9.56 | 12.2 | 41 | n/a | n/a |
| VGG-9/CIFAR10 | 0.85 | 93.2 | 93.5 | 93.5 | 22.80 | 30.34 | 1.24 | 2.14 | 4 | 696 | 542 |
| | 0.9 | | 93.0 | 93.0 | 16.13 | 22.11 | 0.71 | 1.25 | 4 | 470 | 402 |
| VGG-9/CIFAR10 [14] | - | 93.2 | 90.2 | 89.7 | 19.55 | 41.37 | 1.06 | 1.18 | 17 | n/a | n/a |
| VGG-11/CIFAR10 | 0.85 | 93.6 | 93.6 | 93.6 | 24.83 | 36.62 | 2.47 | 4.24 | 4 | 1390 | 1069 |
| | 0.9 | | 93.5 | 93.5 | 18.35 | 26.86 | 1.41 | 1.94 | 4 | 929 | 797 |
| VGG-11/CIFAR10 [4] | n/a | 93.6 | 90.0 | 90.0 | 0.49 | | 0.87 | | 24 | n/a | n/a |



Fig. 4: Layer-by-layer comparison with DNN+NeuroSim [14]

## VI. CONCLUSION

We introduced a full-stack solution for efficient NN processing in CAMs. On the software side, our approach leverages compiler transformations applicable to TWNs, ensuring lower arithmetic intensity without compromising accuracy in various large networks. On the hardware front, we enhance existing AP models by integrating RTM cells, harnessing multi-level cell capability and sequential access for bit-serial word-parallel convolution with reduced data transfers within the accelerator. This combined approach achieves a remarkable 7.5× energy efficiency improvement over crossbar-based accelerators.

## ACKNOWLEDGMENTS

## REFERENCES

[1] K. Smagulova et al., "Resistive neural hardware accelerators," *JPROC*, vol. 111, no. 5, pp. 500–527, 2023.

[2] W. Choi et al., "Content addressable memory based binarized neural network accelerator using time-domain signal processing," in *DAC*, 2018, pp. 1–6.

[3] J. Ran et al., "Pecan: A product-quantized content addressable memory network," in *DATE*. IEEE, 2023, pp. 1–6.

[4] D.-T. Nguyen et al., "Deepcam: A fully cam-based inference accelerator with variable hash lengths for energy-efficient deep neural networks," in *DATE*. IEEE, 2023, pp. 1–6.

[5] M. Imani, D. Peroni, Y. Kim, A. Rahimi, and T. Rosing, "Efficient neural network acceleration on gpgpu using content addressable memory," in *DATE*. IEEE, 2017, pp. 1026–1031.

[6] S. K. Esser, J. L. McKinstry, D. Bablani, R. Appuswamy, and D. S. Modha, "Learned step size quantization," *arXiv preprint arXiv:1902.08153*, 2019.

[7] J. Diffenderfer and B. Kailkhura, "Multi-prize lottery ticket hypothesis: Finding accurate binary neural networks by pruning a randomly weighted network," *arXiv preprint arXiv:2103.09377*, 2021.

[8] M. Scherer, G. Rutishauser, L. Cavigelli, and L. Benini, "Cutie: Beyond petaop/s/w ternary dnn inference acceleration with better-than-binary energy efficiency," *TCAD*, vol. 41, no. 4, pp. 1020–1033, 2021.

[9] R. Bläsing et al., "Magnetic racetrack memory: From physics to the cusp of applications within a decade," *JPROC*, vol. 108, no. 8, pp. 1303–1321, 2020.

[10] Y. Zha and J. Li, "Hyper-ap: Enhancing associative processing through a full-stack optimization," in *ISCA*. IEEE, 2020, pp. 846–859.

[11] P. Junsangsri, J. Han, and F. Lombardi, "A non-volatile low-power tcam design using racetrack memories," in *IEEE-NANO*. IEEE, 2016, pp. 525–528.

[12] K. P. Gnawali et al., "Low power spintronic ternary content addressable memory," *TNANO*, vol. 17, no. 6, pp. 1206–1216, 2018.

[13] R. Zhang, C. Tang, X. Sun, M. Li, W. Jin, P. Li, X. Cheng, and X. S. Hu, "Sky-tcam: Low-power skyrmion-based ternary content addressable memory," *TED*, 2023.

[14] X. Peng, S. Huang, Y. Luo, X. Sun, and S. Yu, "Dnn+ neurosim: An end-to-end benchmarking framework for compute-in-memory accelerators with versatile device technologies," in *IEDM*. IEEE, 2019, pp. 32–5.