

# An Event-Based Digital Compute-In-Memory Accelerator with Flexible Operand Resolution and Layer-Wise Weight/Output Stationarity

Nicolas Chauvaux<sup>1</sup>, Adrian Kneip<sup>1,2</sup>, Christoph Posch<sup>3</sup>, Kofi Makinwa<sup>1</sup>, and Charlotte Frenkel<sup>1</sup>

<sup>1</sup> Department of Microelectronics, Delft University of Technology, The Netherlands

<sup>2</sup> ESAT Department, KU Leuven, Belgium <sup>3</sup> Prophesee, 75012 Paris, France

**Abstract**—Compute-in-memory (CIM) accelerators for spiking neural networks (SNNs) are promising solutions to enable  $\mu\text{s}$ -level inference latency and ultra-low energy in edge vision applications. Yet, their current lack of flexibility at both the circuit and system levels prevents their deployment in a wide range of real-life scenarios. In this work, we propose a novel digital CIM macro that supports arbitrary operand resolution and shape within a unified CIM storage for weights and membrane potentials. These circuit-level techniques enable a hybrid weight- and output-stationary dataflow at the system level to maximize operand reuse, thereby minimizing costly on- and off-chip data movements during the SNN execution. Measurement results of a fabricated FlexSpIM prototype in 40-nm CMOS demonstrate a  $2\times$  increase in 1-bit-normalized energy efficiency compared to prior fixed-precision digital CIM-SNNs, while providing resolution reconfiguration with bitwise granularity. Our approach can save up to 90% energy in large-scale systems, while reaching a state-of-the-art classification accuracy of 95.8% on the IBM DVS gesture dataset.

**Index Terms**—Digital compute-in-memory, spiking neural networks, flexible operand resolution, hybrid-stationary dataflow.

## I. INTRODUCTION

Deploying convolutional neural networks (CNNs) to enable vision at the edge has unlocked applications ranging from user recognition to decentralized object detection and classification. Recently, event-based vision has emerged as a promising opportunity for reduced system-level energy and latency [1]. Indeed, pixels of event-based cameras fire events independently, generating sparse event streams at a  $\mu\text{s}$ -level temporal resolution (Fig. 1(a)), which calls for dedicated algorithms [2]. Among them, spiking neural networks (SNNs) exploit bio-inspired neuron models (e.g., integrate-and-fire (IF) in Fig. 1(b)) to sparsely process information using binary spikes: they rely on an internal state called membrane potential to retain the evolution of information across time. While SNNs are well suited to map per-timestep processing scenarios for low-latency decisions at the edge (Fig. 1(c)), large-scale models that cannot fit entirely in the on-chip memory suffer from significant data movement overheads, which prevents their efficient deployment on conventional hardware architectures.

To that end, compute-in-memory (CIM) hardware for SNNs has recently been proposed [3]–[12], but their limited reconfigurability makes them fundamentally unsuited to the large diversity of CNN layer specifications (Fig. 1(a)). This either implies sub-optimal workload mapping to CIM hardware, leading to latency and energy penalties, or increased time-to-

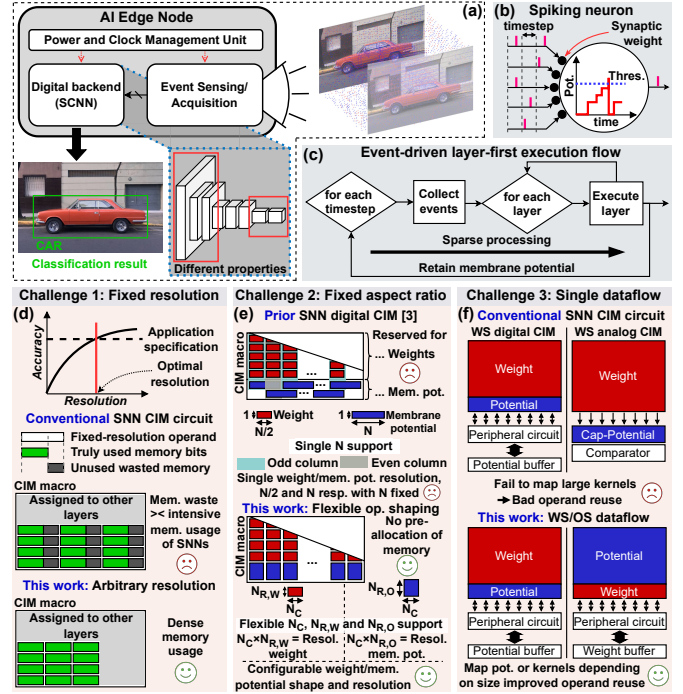


Fig. 1. (a) Event-based edge vision system and workload example. (b) Integrate-and-fire spiking neuron model. (c) Adopted execution flow targeting low latency execution. (d-f) Three core challenges in the state of the art and corresponding innovations in the proposed design.

market due to the need for application-specific CIM hardware. In this work, we propose FlexSpIM, a digital CIM-based accelerator for SNN inference with high flexibility at the circuit and system levels. It solves three core challenges compared to prior CIM-based SNNs:

- 1) while previous works only support a fixed resolution or a few pre-defined options, FlexSpIM supports a fully reconfigurable resolution for weights and membrane potentials, thereby expanding the exploration of the trade-off landscape between accuracy, energy efficiency, and memory footprint for SNN workloads (Fig. 1(d));
- 2) while operand mapping is usually restricted to either fully bit-serial row-wise or fully bit-parallel column-wise, FlexSpIM allows for reconfigurable operand shapes to support different, non-proportional resolution values for weights and membrane potentials, which were otherwise constrained to fixed ratios [3] (Fig. 1(e));
- 3) while maximizing operand stationarity is key to reduce

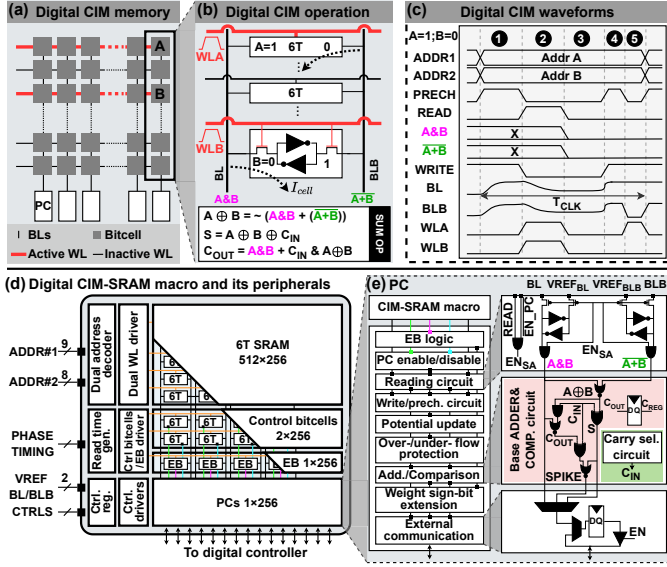


Fig. 2. (a) Digital CIM memory overview: array of bitcells and peripheral circuit (PC) attached to each bitline (BL) for operation handling. Following digital CIM operations, two wordlines (WLs) are simultaneously activated. (b) Example of a digital CIM operation between two bitcells storing  $A=1$  and  $B=0$ . Two boolean operations are obtained and can be used to obtain a 1-bit full adder by following the equations provided. (c) Example of waveforms and phases to perform the digital CIM operation illustrated in (b). (d) Architecture of the proposed FlexSpIM digital CIM-SRAM macro. (e) Decomposition of a PC into modules with their detailed schematics.

external memory accesses, prior CIM-based SNNs only support weight stationarity, and thus are ill-suited for layers that are bottlenecked by membrane-potential data movement (e.g., first layers of ResNet [13]). FlexSpIM introduces a unified weight/membrane potential memory that allows for a hybrid dataflow, making the best out of both weight stationarity (WS) and output (i.e., membrane potential) stationarity (OS) on a per-layer basis. This minimizes operand replacement in the CIM macro for the selected workload, thereby directly alleviating the data movement efficiency bottleneck of previous work (Fig. 1(f)).

## II. PROPOSED RECONFIGURABLE DIGITAL CIM MACRO

SRAM-based CIM is a computing paradigm where operations are carried out directly inside the SRAM, harnessing a low-cost reuse of the data stored in the array. In the case of boolean digital CIM-SRAMs [14], multiple bitwise operations can be performed in parallel on the different vertical bitlines (BLs) by enabling two shared horizontal wordlines (WLs) at a time (Figs. 2(a) and (b)). Applied to SNNs, these architectures can execute the XNOR-and-accumulate operation of IF neurons (Fig. 1(b)) by sequentially accessing corresponding bits of a weight  $A$  and a membrane potential  $B$ , and updating the potential accordingly [3]. In this work, we split this operation in five phases (Fig. 2(c)): 1) precharge of BL/complementary BL (BLB) to VDD, 2) AND/NOR operation between the stored  $A$  and  $B$  operand values on BL/BLB by activating the corresponding WLs, 3) 1-bit sum and carry-out generation in the peripheral circuit (PC) yielding the updated membrane potential, 4) half-select-prevention BL/BLB precharge, and

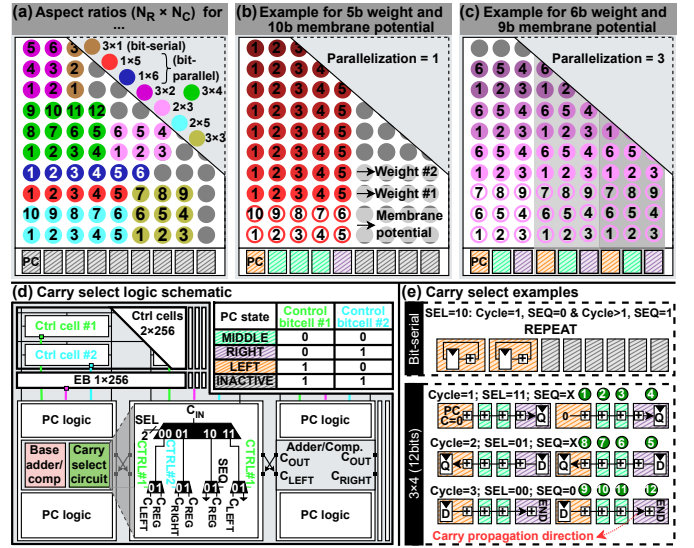


Fig. 3. (a) Arbitrary resolution and operand shaping principle. (b) Example of operand shaping for 5-bit weight and 10-bit membrane potential with a selected parallelization of one neuron. (c) Example of operand shaping for 6-bit weight and 9-bit membrane potential with a selected parallelization of three neurons. (d) Carry-selection logic and PC state configuration modes. (e) PC configurations for bit-serial and  $4 \times 3$  operand shaping.

5) write-back of the new membrane potential bit into  $A$ . These steps are repeated until all bits of the membrane potential and weight have been processed. Then, a comparison between the resulting membrane potential and a threshold conditionally generates an output spike, transferred to the next layer.

Reconfigurability in the FlexSpIM macro (Fig. 2(d)) is achieved by combining a dense 6T SRAM array storing both the weights and membrane potentials with a modular PC per column (Fig. 2(e)). Two additional control bitcells define the per-PC state, while emulation bits (EBs) allow for sign-bit extension and write-free CIM operation during data broadcasting in the macro. Each pitch-matched PC consists of a dual sense amplifier (SA) for the individual readout of BL/BLB, a 1-bit full-adder adapted from [14] with a comparison and a carry-in selection circuits, and logic for I/O communication.

### A. Arbitrary Operand Resolution and Shape

FlexSpIM leverages arbitrary operand resolutions (i.e., 1-to-512 $\times$ 256-bit with bitwise granularity), which can be selected on a per-layer basis for both weights and membrane potentials (Fig. 3(a)). This degree of flexibility overcomes the limited set of resolutions supported in previous works [3]–[12], thereby preventing any waste of storage space. To support weight and membrane potential operands that may take different and non-proportional resolutions, the carry-select circuit (Fig. 2(e)) allows FlexSpIM to map operand bits using any shape in the unified 6T SRAM array (Figs. 3(b) and (c)). The number of columns occupied by each multi-bit operand is defined using the 2-bit control bitcells, which allow chaining multiple 1-bit adders of neighboring PCs for multi-bit computation by changing the carry-in origin (Fig. 3(d)). The multi-bit CIM operation then occurs in parallel over the columns, and sequentially from the LSB row to the MSB row (Fig. 3(e)). For

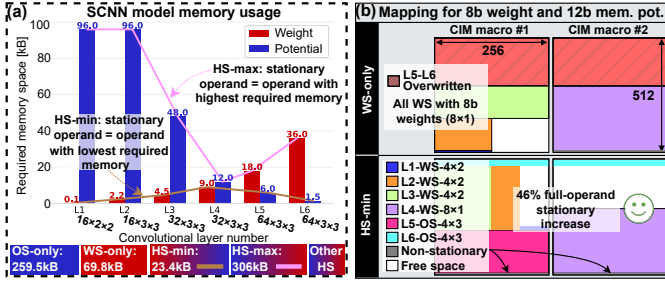


Fig. 4. (a) Layer-level memory requirements for weights and membrane potentials of a spiking CNN composed of six convolutional layers (i.e., L1 to L6) and three FC layers (not shown), with two different HS dataflow situations highlighted by the brown and pink lines. (b) Mapping of the model on two CIM macros for WS-only and HS-min dataflows.

multi-cycle operations, a ping-pong left/right sum direction is used between subsequent cycles to keep inter-PC data movement bounded to their direct neighbors, ensuring design scalability to any macro dimensions.

### B. Hybrid-Stationary Dataflow

At the system level, the unified memory of FlexSpIM allows to support hybrid stationarity (HS) by choosing between WS and OS on a per-layer basis, contrary to previous WS-only CIM designs for SNNs [3]–[6], [9]–[12]. To illustrate this point, we consider a typical spiking CNN workload (SCNN), made of six convolution layers (defined in Fig. 4(a)) followed by three fully-connected (FC) ones (not shown). Exploiting the known memory requirements of both weight and membrane potential operands in each of the SCNN’s layers, the HS flow selects each layer’s dataflow type in order to maximize the overall utilization of the CIM storage space (Figs. 4(a) and (b)), thereby increasing the overall operand stationarity across the multi-timestep execution of a model. This principle can be extended to multiple macros, where a full HS scenario requires at least two macros to ensure the full stationarity of at least one of the operands of every layer when targeting our SCNN workload. Fig. 4(a) depicts two HS dataflows in which the stationary operand is either the one requiring the least or the most memory, respectively referred to as HS-min and HS-max. Compared to the conventional WS-only dataflow, HS-min increases the amount of stationary operands by 46% with an optimal layer mapping across both macros (Fig. 4(b)). Further efficiency gains can be unlocked by scaling up the number of macros, where additional CIM storage avoids frequent external memory accesses by ensuring the stationarity of the operands with the largest memory footprint (Fig. 4(a)).

### III. IMPLEMENTATION AND MEASUREMENT RESULTS

The proposed CIM macro is integrated within the complete accelerator shown in Fig. 5(a). Beyond the proposed 16kB CIM-SRAM macro, it is composed of a 4.25kB memory for per-timestep input spike buffering, and  $4 \times 4$  banks of 2kB SRAMs buffering the SNN weights (resp. membrane potentials) in OS (resp. WS) mode. A 32-to-256-bit bandwidth-adaptive merge-and-shift unit ensures correctly aligned data transfers for arbitrary CIM configurations. The microphotograph of the chip in bulk 40-nm CMOS is shown in Fig. 5(b).

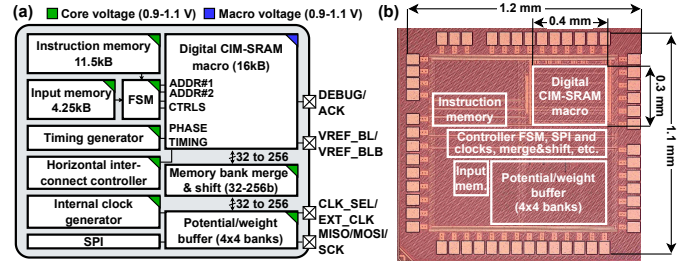
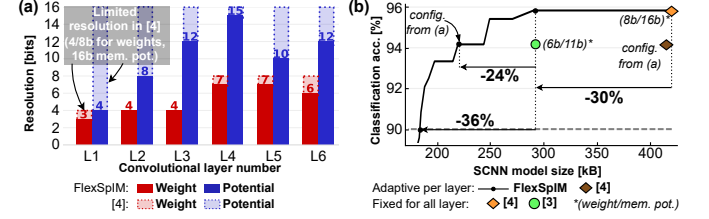


Fig. 5. (a) Overall FlexSpIM system architecture and (b) chip microphotograph in bulk 40-nm CMOS.





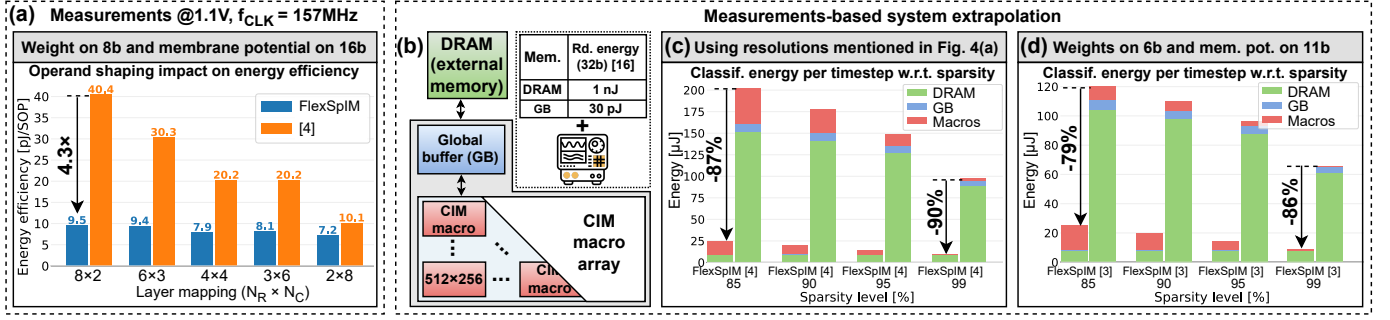


Fig. 7. (a) Measurement results of shape-dependent energy efficiency. (b) System architecture used for the many-macro extrapolation. (c-d) System-level many-macro extrapolation comparing FlexSpIM to [4] and [3], respectively.

TABLE I  
COMPARISON TO THE STATE OF THE ART OF SNN ACCELERATORS.

|  | This work                | SSC-L'21 [3]             | ISSCC'24 [4]              | JSSC'23 [5]         | A-SSCC'22 [6]               | ISSCC'22 [15]              |
|--|--------------------------|--------------------------|---------------------------|---------------------|-----------------------------|----------------------------|
| Technology                               | 40nm                     | 65nm                     | 22nm                      | 28nm                | 65nm                        | 28nm                       |
| Implementation                           | Digital (CIM)            | Digital (CIM)            | Analog CIM                | Analog CIM          | Analog CIM                  | Digital                    |
| Core area (mm <sup>2</sup> )             | 1.37                     | 0.089 <sup>a</sup>       | 2.28                      | 2.9                 | 0.25 <sup>a</sup>           | 0.45                       |
| Macro memory capacity (kB)               | 16                       | 1.37                     | 4                         | 20                  | 4                           | N/A                        |
| Bitcell type                             | 6T                       | 10T                      | 6T                        | 8T                  | 2×6T+6T                     | N/A                        |
| Spiking network type                     | CNN                      | Modified LeNet5          | Residual CNN              | ResNet-12           | CNN                         | RNN                        |
| Representative dataset                   | DVS gesture <sup>b</sup> | MNIST/IMDB               | DVS gesture <sup>b</sup>  | CIFAR-10            | MNIST/CIFAR10               | DVS gesture <sup>b</sup>   |
| Accuracy on DVS gesture                  | 95.8%                    | N/A                      | 94%                       | N/A                 | N/A                         | 87.3%                      |
| Multi-aspect ratio support               | ✓                        | ×                        | ×                         | ×                   | ×                           | ×                          |
| HS support                               | ✓                        | ×                        | ×                         | ×                   | ×                           | ×                          |
| Mem. pot. resolution                     | Any                      | 11b                      | 16b                       | 8b                  | Analog                      | 16b                        |
| Weight resolution                        | Any                      | 6b                       | 4/8b                      | 1/4/8b              | 1.5b                        | 8b                         |
| Weight/Mem. pot. location                | Not fixed                | Fixed                    | Fixed                     | Fixed               | Fixed                       | Fixed                      |
| Supply range (V)                         | 0.9 – 1.1                | 0.7 – 1.2                | 0.55 – 0.9                | 1.1                 | N/A                         | 0.5 – 0.8                  |
| Frequency (MHz)                          | 75.5 – 157               | 66.7 – 500               | 51 – 280                  | 200                 | N/A                         | 13 – 115                   |
| Peak throughput (GSOPS)*                 | 1.2 – 2.5 <sup>c</sup>   | 0.07 – 0.5 <sup>d</sup>  | N/A                       | N/A                 | 163.8 <sup>f</sup>          | 0.013 – 0.115 <sup>c</sup> |
| 1b-norm. throughput (GSOPS) <sup>‡</sup> | 154 – 320                | 4.62 – 33                | N/A                       | N/A                 | N/A                         | 1.67 – 14.7                |
| Power (mW)                               | 6.8 – 17.9 <sup>a</sup>  | 0.1 – 0.9 <sup>a</sup>   | 0.524 – 6.4               | 15.84 <sup>a</sup>  | 0.56 <sup>a</sup>           | 0.077 – N/A                |
| Efficiency (pJ/SOP)*                     | 5.7 – 7.2 <sup>c</sup>   | 1.09 – 1.74 <sup>d</sup> | 3.78 – 10.01 <sup>c</sup> | 0.0016 <sup>e</sup> | 3.45 × 10 <sup>-3</sup> (f) | 5.3 – 12.8 <sup>c</sup>    |
| 1b-norm. eff. (fJ/SOP) <sup>†</sup>      | 44.5 – 56.3              | 16.5 – 26.4              | 29.5 – 78.2               | 0.025               | N/A                         | 41.4 – 100                 |

<sup>a</sup>CIM macro only <sup>b</sup> 10 classes

<sup>c</sup> 8-bit weight and 16-bit mem. pot. <sup>d</sup> 6-bit weight and 11-bit mem. pot.

<sup>f</sup> 1.5-bit weight and N/A for mem. pot. <sup>‡</sup>GSOPS × weight-bit × pot.-bit

<sup>†</sup>fJ/SOP/(weight-bit × pot.-bit) \*1 SOP = 1 addition + mem. pot. update

mainly originates from the PC standby mode that decreases the energy of inactive columns by 87%. A detailed comparison between FlexSpIM and state-of-the-art SNN accelerators is provided in Table I, highlighting 2× better 1-bit-normalized energy efficiency compared to past digital CIM while supporting bitwise granularity on resolution reconfiguration.

### B. System Level

At the system level, accurate evaluation can be obtained by considering a many-macro CIM array architecture with a global on-chip buffer and an external DRAM (Fig. 7(b)). By extracting the energy efficiency of the complete system, accounting for macro-level measurements, the impact of FlexSpIM's system-level flexibility is assessed on the six-layer SCNN workload. First, considering the optimum-resolution mappings of FlexSpIM and [4] (Fig. 3(a)), a FlexSpIM-based system composed of 16 CIM macros with HS maximizing the amount of stationary operands achieves an 87 – 90% energy efficiency gain in the 85 – 99% input sparsity range. Second, considering the fixed 6-bit weight and 11-bit membrane potential resolutions of IMPULSE [3], a FlexSpIM-based system

with 18 CIM macros achieves a 79 – 86% energy efficiency gain in the 85 – 99% sparsity range.

### IV. CONCLUSION

With FlexSpIM, we introduced arbitrary resolution and operand shaping, together with a unified weight/membrane potential memory, to enable CIM-based SNN hardware with the highest resolution and dataflow configurability. With a prototype fabricated in 40-nm CMOS, we demonstrated based on silicon measurements that this flexibility, enabling a 79 – 90% reduction of energy per operation at the system level, is achieved while maintaining a competitive macro-level trade-off between peak throughput and energy per operation, especially compared to past digital CIM solutions. This work thus underlines the importance of macro-level flexibility to enable significant system-level gains in scalable architectures for real-world edge-vision tasks.

### ACKNOWLEDGMENT

This project was co-funded by Prophesee and by the Dutch government as an HTSM-TKI project. The authors would like to thank Douwe den Blanken and Martin Lefebvre, from tapeout support to fruitful discussions.

## REFERENCES

- [1] A. Amir *et al.*, "A Low Power, Fully Event-Based Gesture Recognition System", in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, 2017, pp. 7388-7397.
- [2] G. Chen, H. Cao, J. Conradt, H. Tang, F. Rohrbein, and A. Knoll, "Event-Based Neuromorphic Vision for Autonomous Driving: A Paradigm Shift for Bio-Inspired Visual Sensing and Perception", *IEEE Signal Processing Magazine*, vol. 37, no. 4, pp. 34-49, July 2020.
- [3] A. Agrawal, M. Ali, M. Koo, N. Rath, A. Jaiswal, and K. Roy, "IMPULSE: A 65-nm Digital Compute-in-Memory Macro With Fused Weights and Membrane Potential for Spike-Based Sequential Learning Tasks", *IEEE Solid-State Circuits Letters*, vol. 4, pp. 137-140, June 2021.
- [4] Y. Liu *et al.*, "30.2 A 22nm 0.26nW/Synapse Spike-Driven Spiking Neural Network Processing Unit Using Time-Step-First Dataflow and Sparsity-Adaptive In-Memory Computing", in *2024 IEEE International Solid-State Circuits Conference (ISSCC)*, San Francisco, CA, USA, 2024, pp. 484-486.
- [5] S. Kim, S. Kim, S. Um, S. Kim, K. Kim, and H. -J. Yoo, "Neuro-CIM: ADC-Less Neuromorphic Computing-in-Memory Processor With Operation Gating/Stopping and Digital-Analog Networks", *IEEE Journal of Solid-State Circuits*, vol. 58, no. 10, pp. 2931-2945, Oct. 2023.
- [6] J. Song *et al.*, "Spike-CIM: A 290TOPS/W Spike-Encoding Sparsity-Adaptive Computing-in-Memory Macro with Differential Charge-Domain Integrate-and-Fire", in *2022 IEEE Asian Solid-State Circuits Conference (A-SSCC)*, Taipei, Taiwan, 2022, pp. 1-3.
- [7] S. Kim *et al.*, "A Reconfigurable 1T1C eDRAM-based Spiking Neural Network Computing-In-Memory Processor for High System-Level Efficiency", in *2023 IEEE International Symposium on Circuits and Systems (ISCAS)*, Monterey, CA, USA, 2023, pp. 1-5.
- [8] Y. Liu *et al.*, "An 82-nW 0.53-pJ/SOP Clock-Free Spiking Neural Network With 40- $\mu$ s Latency for AIoT Wake-Up Functions Using a Multilevel-Event-Driven Bionic Architecture and Computing-in-Memory Technique", *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 70, no. 8, pp. 3075-3088, Aug. 2023.
- [9] H. Fu *et al.*, "DS-CIM: A 40nm Asynchronous Dual-Spike Driven, MRAM Compute-In-Memory Macro for Spiking Neural Network", *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 71, no. 4, pp. 1638-1650, Apr. 2024.
- [10] F. Moro *et al.*, "Hardware Calibrated Learning to Compensate Heterogeneity in Analog RRAM-based Spiking Neural Networks", in *2022 IEEE International Symposium on Circuits and Systems (ISCAS)*, Austin, TX, USA, 2022, pp. 380-383.
- [11] B. Wang *et al.*, "SNNIM: A 10T-SRAM based Spiking-Neural-Network-In-Memory Architecture with Capacitance Computation", in *2022 IEEE International Symposium on Circuits and Systems (ISCAS)*, Austin, TX, USA, 2022, pp. 3383-3387.
- [12] A. Singh, M. A. Lebdeh, A. Gebregiorgis, R. Bishnoi, R. V. Joshi, and S. Hamdioui, "SRIF: Scalable and Reliable Integrate and Fire Circuit ADC for Memristor-Based CIM Architectures", *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 68, no. 5, pp. 1917-1930, May 2021.
- [13] K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016, pp. 770-778.
- [14] C. Eckert *et al.*, "Neural Cache: Bit-Serial In-Cache Acceleration of Deep Neural Networks", in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*, Los Angeles, CA, USA, 2018, pp. 383-396.
- [15] C. Frenkel and G. Indiveri, "ReckOn: A 28nm Sub-mm<sup>2</sup> Task-Agnostic Spiking Recurrent Neural Network Processor Enabling On-Chip Learning over Second-Long Timescales", in *2022 IEEE International Solid-State Circuits Conference (ISSCC)*, San Francisco, CA, USA, 2022, pp. 1-3.
- [16] M. Horowitz, "1.1 Computing's Energy Problem (and What We Can Do About It)", in *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, San Francisco, CA, USA, 2014, pp. 10-14.