Are SNNs Truly Energy-efficient? — A Hardware Perspective

Abhiroop Bhattacharjee*, Ruokai Yin*, Abhishek Moitra*, and Priyadarshini Panda Department of Electrical Engineering, Yale University, USA Email: {abhiroop.bhattacharjee, ruokai.yin, abhishek.moitra, priya.panda}@yale.edu

Abstract-Spiking Neural Networks (SNNs) have gained attention for their energy-efficient machine learning capabilities, utilizing bio-inspired activation functions and sparse binary spike-data representations. While recent SNN algorithmic advances achieve high accuracy on large-scale computer vision tasks, their energyefficiency claims rely on certain impractical estimation metrics. This work studies two hardware benchmarking platforms for large-scale SNN inference, namely SATA and SpikeSim. SATA is a sparsity-aware systolic-array accelerator, while SpikeSim evaluates SNNs implemented on In-Memory Computing (IMC) based analog crossbars. Using these tools, we find that the actual energyefficiency improvements of recent SNN algorithmic works differ significantly from their estimated values due to various hardware bottlenecks. We identify and addresses key roadblocks to efficient SNN deployment on hardware, including repeated computations & data movements over timesteps, neuronal module overhead and vulnerability of SNNs towards crossbar non-idealities.

Index Terms—Spiking Neural Networks, Systolic-arrays, Inmemory Computing, Crossbars, Energy-efficiency

I. INTRODUCTION

Spiking Neural Networks (SNNs) have garnered significant attention as a power-efficient solution for machine learning [1], [2]. SNNs process data over multiple time steps using biologically inspired non-linear activation functions, such as Leaky-integrate-and-Fire (LIF) neurons. During each time step, input data is represented as either a spike (binary 1) or no-spike (binary 0), creating a sparsely encoded temporal spike-data representation. This representation potentially offers several hardware advantages: (1) Multiplier-less Computation: SNNs use computation units that rely solely on accumulators for dot-products, eliminating the need for multipliers used in Artificial Neural Networks (ANNs) for Multiply-and-Accumulate (MAC) operations [3], [4]. (2) Reduced On-chip **Memory:** The binary nature of SNNs significantly reduces the on-chip memory required to store intermediate layer activations during SNN processing. These features add to the energy-efficiency of SNN algorithms.

Fortunately, over the last few years, there have been huge advances in the SNN training algorithms [5]–[14] leading to state-of-the-art classification accuracy at low timesteps on large-scale image datasets such as CIFAR10, CIFAR100, Tiny-ImageNet and ImageNet. However, being alogrithm-focused, the energy-efficiency claimed by these works are based on primitive metrics such as FLOPs, timesteps and spike-data

TABLE I: Summary of Different Works								
Work	Training (T) or Inference (I)	Platform	Hardware Benchmarking					
Small-scale Optimization Tasks								
BrainScale [24]	T & I	Analog	X					
Loihi [18]	I	Digital	X					
TrueNorth [25]	I	Digital	X					
Large-scale Computer Vision Tasks								
SpinalFlow [19], PTB [20]	I	Digital	x					
RESPARC [21]	I	Analog	X					
H2Learn [22]	Т	Digital	X					
SATA [26]	T & I	Digital	✓					
SpikeSim [27]	Ι	Analog	1					

sparsity. Such energy evaluation is impractical as metrics such as FLOPs do not account for hardware overheads like memory access and data communication. Additionally, real systolic-array [15] and In-memory Computing (IMC) [16] accelerators are ineffective in handling the spike-data sparsity, particularly during the memory fetches. Therefore, there is a need for realistic SNN hardware benchmarking platforms. As shown in Table I, there are several SNN-specific hardware accelerators. Works such as BrainScale [17], Loihi [18] and TrueNorth [3] are geared towards small-scale optimization tasks for SNNs. In more recent years, there have been several hardware co-design works [19]–[22] that cater to large-scale SNN implementations. However, they lack several practical considerations such as the data communication overhead, LIF activation storage and hardware non-idealities [23].

To this end, we study two hardware accelerators SATA and SpikeSim. Unlike prior SNN hardware platforms [19]–[22], both SATA and SpikeSim support end-to-end hardwarerealistic benchmarking of large-scale SNNs, during inference. SATA [26] is a sparsity-aware systolic-array based training and inference accelerator for SNNs. While SATA evaluates SNN workloads on a fully-digital CMOS platform, SpikeSim [27] performs hardware-realistic accuracy, energy, latency and area evaluation of SNN workloads on IMC analog crossbars based on Resistive Random-access Memories (RRAMs) [28].

Table II shows the estimated and hardware-realistic (SATA and SpikeSim implemented) energy-efficiency improvements of state-of-the-art SNN algorithms during inference. The estimated energy is proportional to the product of FLOPs, timesteps and the sparsity (as shown in footnote of Table II). Evidently, there is a significant difference between the estimated and hardware-realistic energy-efficiency improvements. To this end, in this work, we perform realistic SNN benchmarking on the SATA and SpikeSim platforms. With this, we bring forth the key bottlenecks that SNNs exhibit on hardware and propose effective mitigation strategies. Essentially, we address the following key bottlenecks: (1) repeated

^{*}Equal contribution.

This work was supported in part by CoCoSys, a JUMP2.0 center sponsored by DARPA and SRC, the NSF CAREER Award, TII (Abu Dhabi), and the DoE MMICC center SEA-CROGS (Award #DE-SC0023198).

TABLE II: Table showing energy comparisons for recent SNN algorithm works using the CIFAR10 dataset. Est. Energy denotes the qualitative estimated energy (nJ) of the SNNs calculated using the equation specified in the footnote. $(N \times)$ denotes the estimated or actual improvements in energy as compared to the corresponding baselines in each work. E_{AC} denotes the energy expended by a single INT8 Accumulation (AC) operation. All energy values are reported in 28 nm CMOS technology node.

Work	Accuracy	Sparsity	Т	Est. Energy ² (nJ)	Actual Er SATA ³	nergy (nJ) SpikeSim ⁴
S-BP [5]	89.3%	$90.0\%^{1}$	50	$11.7e^4(10\times)$	$1.1e^{7}(4.5\times)$	$2.3e^{5}(5.2\times)$
BNTT [7]	90.3%	90.5%	20	$4.2e^{4}(20\times)$	$0.6e^7(2.8\times)$	$0.9e^{5}(4\times)$
Direct [8]	90.5%	$90.0\%^{1}$	10	$11.7e^4(20\times)$	$3.8e^7(2.3\times)$	$2.6e^{5}(4\times)$
TSSL [9]	91.4%	90.1%	5	$5.8e^4(80 \times)$	$3.1e^{7}(4.9\times)$	$1.3e^{5}(16\times)$
LTH [11]	93.2%	$97.5\%^{1}$	5	$2.7e^4(15 \times)$	$5.7e^{7}(1.3\times)$	$2.6e^{5}(2\times)$
TDBN [10]	92.9%	85.0%	4	$13.3e^{4}(83 \times)$	$5.4e^{7}(6.8\times)$	$2.1e^5(25\times)$

¹ We use 90% sparsity for [5], [8]. For [11], we use the weight sparsity.

 $^{2}E_{est} = FLOPs \times Timesteps \times (1 - Sparsity) \times E_{AC}$ ³ Codes are available at: https://github.com/Intelligent-Computing-Lab-Yale/SATA

⁴ Codes are available at: https://github.com/Intelligent-Computing-Lab-Yale/SpikeSim

memory accesses and computations over multiple timesteps, (2) overhead of the LIF neuronal module, and (3) vulnerability of IMC-implemented SNNs towards analog crossbar nonidealities. This work encapsulates the above key hardware challenges overlooked by the SNN research community, and motivates future works aimed towards efficient hardwareaware SNN algorithm design.

II. BACKGROUND

Spiking Neural Networks: The distinguishing feature of SNNs lies in their utilization of a different neuronal activation function (most commonly, LIF) for temporal signal processing, as opposed to the ReLU activation commonly used in ANNs. The LIF neuron *i* associated with a membrane potential u_i^t , that accumulates a train of spike inputs as follows:

$$u_i^t = \lambda u_i^{t-1} + \sum_j w_{ij} o_j^t.$$

$$\tag{1}$$

Here, t stands for the timestep, w_{ij} for weight connections between neuron i and neuron j and λ denotes the leak factor. The LIF neuron *i* generates an output spike o_i^t at the end of timestep t if the membrane potential exceeds a threshold θ :

$$o_i^t = \begin{cases} 1, & \text{if } u_i^t > \theta, \\ 0 & \text{otherwise.} \end{cases}$$
(2)

Upon firing, the membrane potential is reset to zero. The integrate-and-fire behavior exhibited by an LIF neuron results in a non-differentiable function, making it challenging to employ standard backpropagation for training SNNs. To this end, Surrogate gradient learning or Backpropagation Through Time (BPTT) addresses the non-differentiability problem of a LIF neuron by approximating the backward gradient function [29] and offers a means to directly learn from spikes using fewer timesteps. Further, BPTT can be implemented using popular machine learning frameworks like PyTorch [30].

Further, following previous work [31], we use the direct encoding method to encode the input tensor into spike trains with total timesteps T. To get the final prediction, we repeat the inference process T times (t = 1, 2, ..., T) and average the output from the SNN output classifier.

Systolic-array Accelerators: Systolic-array architecture is popular among the digital von-Neumann accelerator designs for SNNs [19], [20], [26]. With a regular and dataflow centric design, systolic-arrays can efficiently process the matrixmatrix multiplications in parallel with high spatio-temporal locality and compute density [32] (see Fig. 2). In this work,

we will evaluate the SATA design in two dataflow modes: 1) output-stationary (OS) mode, where the partial sums will remain inside each Processing Engine (PE) of the array during the dot-product operations; 2) weight-stationary (WS) mode, where, the weights are pre-stored into the PEs of the array before the dot-product operation starts.

Analog Crossbars: Analog crossbars comprise of a 2D array of IMC devices, interfaced with Digital-to-Analog Converters (DACs), Analog-to-Digital Converters (ADCs), and write circuits dedicated towards programming the IMC devices [16], [33], [34]. The SNN's spike inputs are encoded as analog voltages V_i to each row of the crossbar by the DACs, while weights are programmed as synaptic device conductances (G_{ij}) at the cross-points, as shown in Fig. 1.



Fig. 1: An analog crossbar array with input voltages V_i , synaptic devices bearing conductances G_{ij} and output currents I_j .

For emulating dot-product operations in case of an ideal N×M crossbar during inference, the voltages interact with the device conductances, resulting in a current governed by Ohm's Law. Finally, adhering to Kirchoff's current law, the net output current sensed at each column j by the ADCs represents the sum of currents flowing through each device, expressed as $I_{j(ideal)}$

 $\sum_{i=1}^{N} G_{ij} * V_i$. In practical scenarios, the analog nature of computation introduces various non-idealities, including interconnect parasitic resistances and synaptic device variations [35]–[38]. Thus in a non-ideal scenario, the net output current sensed at each column j deviates from the ideal value $I_{j(ideal)}$. These deviations manifest as significant accuracy losses for SNNs on crossbars [23].

III. TOOLS FOR BENCHMARKING SNNS ON HARDWARE

A. SATA: A Systolic-array Benchmarking Accelerator

SATA [26] is a sparsity-aware training accelerator designed for benchmarking the state-of-the-art BPTT-based SNN training on a fully digital von-Neumann architecture. Unlike the prior SNN training accelerators, which have numerous complex engines to boost performance, SATA adopts a simple and reconfigurable systolic-array design with a three-level memory hierarchy [39]. This makes it straightforward for the SNN algorithm designers to deploy their workloads on SATA and get an estimation of the hardware energy cost. Though designed as a training accelerator, SATA can be used to benchmark the inference performance of pre-trained SNNs by detaching the training-related components. This work will focus on using SATA to evaluate BPTT-trained SNNs during inference. The architecture design of SATA for inference is shown in Fig. 2. The analyses performed on SATA help identify some major bottlenecks, such as repetitive data movements across timesteps, that hinder SNNs from being energy-efficient on hardware. SATA shows that spike-data sparsity in SNNs can only be leveraged inside the PE computation unit (that performs weighted-accumulation and LIF operations). Outside the computation unit, even sparse input and weight data incur



Fig. 2: Architecture design of the inference version of the SATA tool. Spad refers to scratchpad memory (registers) and PE stands for Processing Engine that performs the weighted-accumulation and LIF operations. The red wire is used to indicate skipping of operations in case of zero operands as inputs.



Fig. 3: SpikeFlow Architecture. GB (GA), TB (TA) and PB (PA) denote global, tile and PE buffers (Accumulators). PO and LIF/IF denote the pooling and LIF/IF neuronal module, respectively.

memory fetches from the on-chip buffers, adding to significant energy overhead.

Challenges: Based on our benchmarking using SATA, we identify two key challenges for deploying SNNs on a systolicarray architecture. Firstly, increased timesteps bring in extra computation and data movement costs. As shown in Fig. 4(b), the energy costs for both PE computation and data movement from the on-chip buffers scales with increasing timesteps. These time-repetitive costs (specifically the data movement cost) reduce the energy-efficiency of the deployed SNNs. A prior work [19] has shown that the inference energy gap between an SNN with 16 timesteps and its ANN counterpart can be as large as $16 \times$ across various workloads on a systolicarray architecture, due to repetitive evaluations across multiple time-steps. The other challenge is attributed to the hardware cost of the LIF units which are used to generate the output spikes and store the membrane potential across timesteps.

B. SpikeSim: An IMC-based Benchamarking Accelerator

SpikeSim [27] is a IMC crossbar-based hardware evaluation tool for benchmarking SNN inference. SpikeSim maps BPTTtrained SNN workloads on a monolithic IMC-based weightstationary tiled architecture, called SpikeFlow (see Fig. 3), and performs hardware-realistic accuracy (incorporating crossbar non-idealities), energy, latency and area evaluations. Spike-Flow incorporates a digital Leaky-Integrate-Fire/Integrate-Fire (LIF/IF) neuronal activation unit to store the intermediate membrane potentials (u^t) and generate spike outputs during SNN inference. Furthermore, the analog crossbars in the SpikeFlow architecture are based on RRAM devices [40]. For crossbar-realistic dot-product operations, SpikeSim emulates the impact of RRAM device read noise and IR-drop nonidealities due to crossbar interconnect parasitics. Another unique characteristic of SpikeSim is a fully digital DIFF module in the SpikeFlow architecture that eliminates the conventional double-crossbar approach for performing signed dot-products, thereby bringing in energy and area savings. SpikeSim uses H-trees to communicate partial sums emerging from the crossbars to the digital peripherals inside a tile, and an inter-tile Network-on-Chip (NoC) architecture to communicate spikes to and from the neuronal unit.

Challenges: Based on our benchmarking using SpikeSim, we bring forth three key challenges to SNN inference on IMC crossbar-based hardware. First, SNNs are highly vulnerable on analog crossbars owing to the impact of the non-idealities which lead to accumulation of errors in the dot-product operations over multiple timesteps (see Fig. 6(b)). Second, unlike ANNs that have a simple ReLU activation unit, SNNs entail a high LIF/IF neuronal area overhead on SpikeSim. This is due to large u^t SRAM cache to store intermediate membrane potentials of different layers of an SNN model during inference. Third, unlike ANNs, the number of timesteps in SNNs plays a crucial role in the hardware performance and inference accuracy. As shown in Fig. 4(a), both the inference energy and latency on SpikeSim scale linearly with timesteps, similar to the trend obtained using SATA in Fig. 4(b).

IV. MITIGATION STRATEGIES

Experimental Setup: For all our experiments, we use BPTTtrained SNNs—VGG9, VGG16 and ResNet18 models, on the CIFAR10 and Tiny ImageNet datasets. We obtain SNN models using code provided in [7]. Unless otherwise mentioned, all the pre-trained SNNs have 8-bit weight-precision. SATA [26] is calibrated in 28 nm CMOS technology node, while SpikeSim [27] is in 65 nm CMOS technology node with RRAM crossbars of size 64×64 in the PEs.

Dynamic Timestep Reduction: To improve the energyefficiency of SNNs by reducing timesteps, we analyse an input-aware Dynamic Timestep SNN (DT-SNN) methodology [41]. DT-SNN dynamically determines the least number of timesteps required for a confident prediction in an inputdependent basis during inference of a pre-trained SNN. This is done by simply appending a digital entropy-computation module with our SNN-based hardware accelerators (SATA or SpikeSim). For every input, the calculated value of entropy of the SNN's predicted output at the end of every timestep is compared against a predefined threshold. An early temporal exit (termination of inference) or a prediction is carried out if the entropy is lower than the set threshold at any given timestep. Note, the entropy-computation module incurs negligible energy overhead to the overall SNN inference energy on SATA or SpikeSim. On SpikeSim, Fig. 4(a) shows a $10.4 \times$ higher Energy-Delay-Product (EDP) on increasing the number of timesteps from 1 to 4 for the inference of a standard VGG16 SNN on the Tiny ImageNet dataset. We find that DT-SNN can reduce the overall EDP by $2.54\times$, while maintaining similar inference accuracy, compared to a standard SNN inference with 4 timesteps across all inputs. It turns out a large fraction of the test images in the Tiny ImageNet dataset are classified



Fig. 4: (a) Impact of DT-SNN on SpikeSim using VGG16 SNN on the TinyImageNet dataset. The EDPs are normalized with respect to the value at timestep = 1. (b) Computation and memory movement costs on SATA across multiple timesteps for VGG9 SNN on the CIFAR10 dataset. Impact of DT-SNN on SATA energy costs.



Fig. 5: (a) Data movement energy cost difference between a standard OS dataflow and the SNN-tailored dataflow for SATA using VGG9 SNN on Tiny ImageNet dataset. (b) The power cost breakdown of computation units on SATA with 128 PEs. The baseline has 128 LIF units. To deploy C#n EfficientLIF-Net, SATA requires 1/n of the baseline LIF units.

with 1 timestep, and overall DT-SNN requires 2.14 timesteps on average across all test images resulting in lower EDP. On SATA, a similar trend of inference energy increase is observed for both computation and memory access energy when the timestep increases from 1 to 4. With DT-SNN, the total energy cost is reduced by 46.5% compared to the standard 4-timestep VGG9 SNN inference on CIFAR10 dataset.

Data-movement Cost Mitigation: As we discussed in Section III-A, one of the major challenges for SNN deployment on digital hardware platforms like SATA is the repetitive data movement costs. We design an SNN-tailored dataflow for SATA that can significantly reduce the repetitive data movement cost for SNNs. In SATA's dataflow design, we adopt the tick-batch method [19] and maximally reuse the weights at the PE level by having scratch-pad memories inside every PE to hold the weights stationary [39] throughout the PE computation process. By utilizing such a dataflow, SATA will only read out the data once from the higher memory hierarchies (DRAM and SRAM) to the PE array for each layer across all timesteps. In Fig. 5 (a), we compare SATA's SNNtailored dataflow with the standard output-stationary dataflow on a VGG9 SNN on the Tiny ImageNet dataset with different timesteps. Utilizing the SNN-tailored dataflow can save 62.5%memory movement energy with a timestep of 4. The benefits will increase when a larger timestep is used. Besides redesigning the dataflow for the hardware, model compression techniques like quantization [42]–[44] and pruning [11], [44], [45] will also help in reducing the data movement costs.

Mitigating LIF Overhead: LIF units are energy-hungry components on the hardware, which can take up to 61.6% total power of the computation units [46]. This translates to approximately $2\times$ higher energy cost for LIF operations compared to other operations. To mitigate the overhead of LIF units, recent work EfficientLIF-Net [46] shares the LIF



Fig. 6: (a) Impact of LIF-sharing on neuronal area using SpikeSim for ResNet18 SNN on Tiny ImageNet dataset. (b) Non-ideal accuracy improvement on SpikeSim for a pre-trained VGG16 SNN (4-bit weights) on Tiny ImageNet dataset with non-ideality-aware weight-encoding & BN adaptation.

neurons across layers and channels. We use the notation of C#n to represent the EfficientLIF-Nets that share 1 LIF neuron for n post-synaptic neurons on the output channel dimension. On SATA, as shown in Fig. 5 (b), we can reduce 75.1% of the power cost of LIF units by having a C#4 LIF-sharing. Quantization of membrane potential [42] can also help to reduce the LIF unit cost by having smaller registers for the membrane potentials. The LIF-sharing and membrane potential quantization methods are orthogonal techniques for mitigating the LIF-units cost. On SpikeSim, we find that LIF sharing with C#2 and C#4 results in $1.38 \times$ and $2.41 \times$ reductions in LIF area, respectively, for a ResNet18 SNN on the TinyImagenet dataset (see Fig. 6(a)).

Crossbar Non-ideality Mitigation: To address the accuracy degradation of SNNs on non-ideal crossbars, two training-less approaches are studied in Fig. 6(b): (1) SpikeSim supports a non-ideality aware weight-encoding scheme on the RRAM crossbars to increase the proportion of high resistance synapses during SNN inference. Prior works have shown that the impact of crossbar non-idealities decreases upon increasing the proportion of high resistance synapses in the crossbars [36], [47], [48]. Thus, the non-ideal SNN accuracy is improved by 40.13%. (2) Non-ideality aware adaptation [23], [49] of the SNN's batchnorm (BN) parameters prior to inference can mitigate the impact of crossbar non-idealities, specifically the interconnect parasitics. During BN adaptation, we forward a number of training image samples through the SNN deployed on crossbars, adapting the moving average & variance of the batchnorm layers with respect to noisy activations (while keeping the learnable parameters or weights frozen). Consequently, the measured accuracy loss on SpikeSim due to non-idealities is reduced to 1.22% compared to the software baseline.

V. CONCLUSION

Our study unveils critical challenges in efficient deployment of large-scale SNNs on hardware, highlighting discrepancies between estimated and hardware-realistic energy-efficiency improvements. We base our study using two state-of-the-art hardware benchmarking tools for SNNs (SATA and SpikeSim) which help identify and address key hardware bottlenecks such as- repeated computations & data movements over timesteps, LIF neuronal module overhead and SNN's vulnerability to crossbar non-idealities. We identify some key mitigation strategies that help address the hardware overheads. The findings from the benchmarking tools underscore the importance realistic hardware-aware SNN algorithm-design in the future for driving low-power neuromorphic applications at the edge.

REFERENCES

- [1] K. Roy et al., "Towards spike-based machine intelligence with neuromorphic computing," *Nature*, 2019. K. Yamazaki *et al.*, "Spiking neural networks and their applications: A
- [2] review," Brain Sciences, 2022.
- [3] F. Akopyan et al., "Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip," IEEE transactions on computer-aided design of integrated circuits and systems, 2015.
- [4] C. Tang and J. Han, "Hardware efficient weight-binarized spiking neural networks," in 2023 Design, Automation & Test in Europe Conference & Exhibition (DATE). IEEE, 2023.
- [5] C. Lee et al., "Enabling spike-based backpropagation for training deep neural network architectures," Frontiers in neuroscience, p. 119, 2020.
- [6] S. S. Chowdhury et al., "Spatio-temporal pruning and quantization for low-latency spiking neural networks," in 2021 International Joint Conference on Neural Networks (IJCNN). IEEE, 2021.
- [7] Y. Kim and P. Panda, "Revisiting batch normalization for training low-latency deep spiking neural networks from scratch," Frontiers in neuroscience, 2021.
- [8] Y. Wu et al., "Direct training for spiking neural networks: Faster, larger, better," in Proceedings of the AAAI conference on artificial intelligence, 2019
- [9] W. Zhang and P. Li, "Temporal spike sequence learning via backpropagation for deep spiking neural networks," Advances in Neural Information Processing Systems, vol. 33, pp. 12022-12033, 2020.
- [10] H. Zheng et al., "Going deeper with directly-trained larger spiking neural networks," in Proceedings of the AAAI conference on artificial intelligence, 2021.
- [11] Y. Kim et al., "Exploring lottery ticket hypothesis in spiking neural networks," in European Conference on Computer Vision. Springer, 2022.
- [12] Y. Venkatesha et al., "Federated learning with spiking neural networks," IEEE Transactions on Signal Processing, 2021.
- [13] Y. Li et al., "Seenn: Towards temporal spiking early-exit neural networks," arXiv:2304.01230, 2023.
- [14] Li et al., "Efficient human activity recognition with spatio-temporal spiking neural networks," Frontiers in Neuroscience, vol. 17, p. 1233037.
- R. Xu et al., "A survey of design and optimization for systolic array [15] based dnn accelerators," ACM Computing Surveys, 2023.
- [16] N. Verma et al., "In-memory computing: Advances and prospects," IEEE Solid-State Circuits Magazine, 2019.
- [17] C. Pehle et al., "The brainscales-2 accelerated neuromorphic system with hybrid plasticity," Frontiers in Neuroscience, 2022.
- [18] M. Davies et al., "Loihi: A neuromorphic manycore processor with onchip learning," Ieee Micro, 2018.
- [19] S. Narayanan et al., "Spinalflow: An architecture and dataflow tailored for spiking neural networks," in International Symposium on Computer Architecture (ISCA). IEEE, 2020.
- [20] J.-J. Lee et al., "Parallel time batching: Systolic-array acceleration of sparse spiking neural computation," in International Symposium on High-Performance Computer Architecture (HPCA). IEEE, 2022.
- [21] A. Ankit et al., "Resparc: A reconfigurable and energy-efficient architecture with memristive crossbars for deep spiking neural networks," in Design Automation Conference, 2017.
- L. Liang et al., "H2learn: High-efficiency learning accelerator for highaccuracy spiking neural networks," IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 2021.
- [23] A. Bhattacharjee et al., "Examining the robustness of spiking neural networks on non-ideal memristive crossbars," International Symposium on Low Power Electronics and Design (ISLPED), 2022.
- [24] I. L. van Soelen et al., "Brain scale: brain structure and cognition: an adolescent longitudinal twin study into the genetic etiology of individual differences," Twin Research and Human Genetics, 2012.
- [25] P. A. Merolla et al., "A million spiking-neuron integrated circuit with a scalable communication network and interface," Science, 2014.
- [26] R. Yin et al., "Sata: Sparsity-aware training accelerator for spiking neural networks," IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 2022.
- [27] A. Moitra et al., "Spikesim: An end-to-end compute-in-memory hardware evaluation tool for benchmarking spiking neural networks," IEEE TCAD, 2023.
- [28] I. Chakraborty et al., "Pathways to efficient neuromorphic computing with non-volatile memory technologies," Applied Physics Reviews, 2020.
- [29] Y. Wu et al., "Spatio-temporal backpropagation for training highperformance spiking neural networks," Frontiers in neuroscience, 2018.
- [30] A. Paszke et al., "Automatic differentiation in pytorch," 2017.

- [31] Y. Kim et al., "Rate coding or direct coding: Which one is better for accurate, robust, and energy-efficient spiking neural networks?" in ICASSP. IEEE, 2022.
- [32] H.-T. Kung, "Why systolic architectures?" Computer, 1982.
- [33] M. Hu et al., "Dot-product engine for neuromorphic computing: Programming 1t1m crossbar to accelerate matrix-vector multiplication," in ACM/EDAC/IEEE DAC, 2016.
- [34] Q. Xia and J. J. Yang, "Memristive crossbar arrays for brain-inspired computing," Nature materials, 2019.
- Sun et al., "Impact of non-ideal characteristics of resistive synaptic [35] devices on implementing convolutional neural networks," IEEE JETCAS, 2019.
- [36] A. Bhattacharjee et al., "Neat: Non-linearity aware training for accurate, energy-efficient and robust implementation of neural networks on 1t-1r crossbars," IEEE TCAD, 2021.
- [37] S. Jain et al., "Rxnn: A framework for evaluating deep neural networks on resistive crossbars," IEEE TCAD, 2020.
- [38] I. Chakraborty et al., "Geniex: A generalized approach to emulating non-ideality in memristive xbars using neural networks," in ACM/IEEE DAC, 2020.
- [39] Y.-H. Chen et al., "Everiss: A spatial architecture for energy-efficient dataflow for convolutional neural networks," ACM SIGARCH Computer Architecture News, 2016.
- [40] B. Hajri et al., "Rram device models: A comparative analysis with experimental validation," IEEE Access, 2019.
- [41] Y. Li et al., "Input-aware dynamic timestep spiking neural networks for efficient in-memory computing," Design and Automation Conference, 2023.
- [42] R. Yin et al., "Mint: Multiplier-less integer quantization for spiking neural networks," arXiv:2305.09850, 2023.
- [43] P.-Y. Tan and C.-W. Wu, "A low-bitwidth integer-stbp algorithm for efficient training and inference of spiking neural networks," in Proceedings of the 28th Asia and South Pacific Design Automation Conference, 2023, pp. 651-656.
- [44] L. Deng et al., "Comprehensive snn compression using admm optimization and activity regularization," IEEE transactions on neural networks and learning systems, 2021.
- [45] R. Yin et al., "Workload-balanced pruning for sparse spiking neural networks," arXiv:2302.06746, 2023.
- [46] Y. Kim et al., "Sharing leaky-integrate-and-fire neurons for memoryefficient spiking neural networks," arXiv:2305.18360, 2023.
- [47] A. Bhattacharjee and P. Panda, "Switchx: Gmin-gmax switching for energy-efficient and robust implementation of binarized neural networks on reram xbars," ACM TODAES, 2021.
- [48] Bhattacharjee et al., "Examining and mitigating the impact of crossbar non-idealities for accurate implementation of sparse deep neural networks," in Design, Automation & Test in Europe Conference & Exhibition (DATE). IEEE, 2022.
- A. Bhattacharjee et al., "Examining the role and limits of batchnorm op-[49] timization to mitigate diverse hardware-noise in in-memory computing,' GLSVLSI, 2023.