

# Design-Technology Co-Optimization for NVM-based Neuromorphic Processing Elements

SHIHAO SONG, ADARSHA BALAJI, ANUP DAS, and NAGARAJAN KANDASAMY, Drexel University, USA

An emerging use-case of machine learning (ML) is to train a model on a high-performance system and deploy the trained model on energy-constrained embedded systems. Neuromorphic hardware platforms, which operate on principles of the biological brain, can significantly lower the energy overhead of a machine learning inference task, making these platforms an attractive solution for embedded ML systems. We present a design-technology tradeoff analysis to implement such inference tasks on the processing elements (PEs) of a Non-Volatile Memory (NVM)-based neuromorphic hardware. Through detailed circuit-level simulations at scaled process technology nodes, we show the negative impact of technology scaling on the information-processing latency, which impacts the quality-of-service (QoS) of an embedded ML system. At a finer granularity, the latency inside a PE depends on 1) the delay introduced by parasitic components on its current paths, and 2) the varying delay to sense different resistance states of its NVM cells. Based on these two observations, we make the following three contributions. First, on the technology front, we propose an optimization scheme where the NVM resistance state that takes the longest time to sense is set on current paths having the least delay, and vice versa, reducing the average PE latency, which improves the QoS. Second, on the architecture front, we introduce isolation transistors within each PE to partition it into regions that can be individually power-gated, reducing both latency and energy. Finally, on the system-software front, we propose a mechanism to leverage the proposed technological and architectural enhancements when implementing a machine-learning inference task on neuromorphic PEs of the hardware. Evaluations with a recent neuromorphic hardware architecture show that our proposed design-technology co-optimization approach improves both performance and energy efficiency of machine-learning inference tasks without incurring high cost-per-bit.

CCS Concepts: • **Hardware** → **Neural systems; Emerging languages and compilers; Emerging tools and methodologies**; • **Computer systems organization** → **Data flow architectures; Neural networks**.

Additional Key Words and Phrases: neuromorphic computing, design-technology co-optimization (dtco), non-volatile memory (NVM), oxide-based resistive random access memory (OxRRAM)

## ACM Reference Format:

Shihao Song, Adarsha Balaji, Anup Das, and Nagarajan Kandasamy. 2021. Design-Technology Co-Optimization for NVM-based Neuromorphic Processing Elements. *ACM Trans. Embedd. Comput. Syst.* 37, 4, Article 111 (August 2021), 28 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

Neuromorphic computing systems are integrated circuits that implement the architecture of central nervous system in primates [14, 22, 65]. These systems facilitate energy-efficient computations using Spiking Neural Networks (SNN) [63] for power-constrained embedded devices. To this end, the design workflow is to train a machine learning model (commonly on a backend server)

---

Authors' address: Shihao Song; Adarsha Balaji; Anup Das; Nagarajan Kandasamy, [anup.das@drexel.edu](mailto:anup.das@drexel.edu), Drexel University, 3141 Chestnut Street, Philadelphia, PA, USA, 19104.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2021 Association for Computing Machinery.

1539-9087/2021/8-ART111 \$15.00

<https://doi.org/10.1145/1122445.1122456>

and subsequently, convert the trained model to spike-based computations and deploy it on the neuromorphic hardware of an embedded system. The quality of inference (e.g., accuracy) is assessed in terms of the inter-spike interval (ISI) (see Section B). Therefore, any deviation from its expected value will lead to a degradation of the inference quality.

A typical neuromorphic system such as Loihi [28], DYNAPs [66] and  $\mu$ Brain [93] consists of processing elements (PEs) that communicate spikes using a shared interconnect. Each PE implements neuron and synapse circuitries. A common technique to implement a neuromorphic PE is using an analog crossbar where bitlines and wordlines are organized in a grid with memory cells connected at their crosspoints to store synaptic weights [2, 32, 40, 45, 46, 51, 61, 69, 100]. Neuron circuitries are implemented along bitlines and wordlines. Figure 1 (left) shows the architecture of an  $N \times N$  analog crossbar with  $N$  bitlines and  $N$  wordlines.

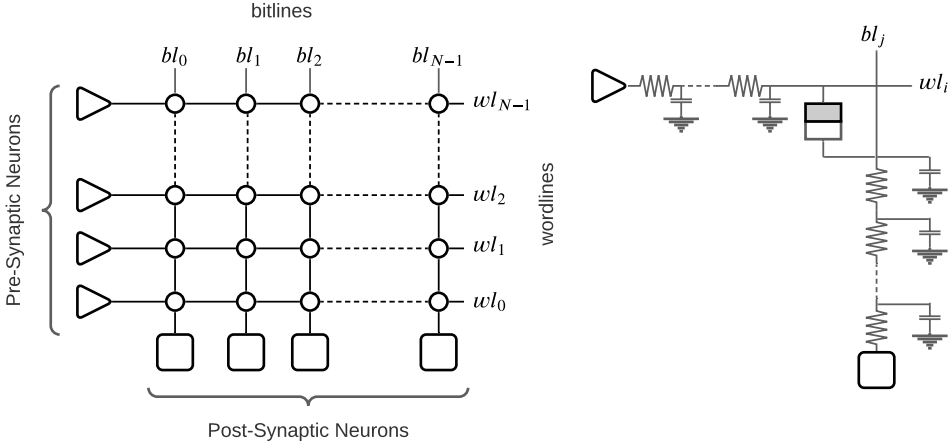


Fig. 1. An  $N \times N$  crossbar showing the parasitic components within.

We investigate the internal architecture of a crossbar and find that parasitic components introduce delay in propagating current from a pre-synaptic neuron to a post-synaptic neuron as illustrated in Figure 1 (right). This delay depends on the specific current path activated in a crossbar. Higher the number of parasitic components on a current path, larger is its propagation delay [70, 86, 88, 89, 91]. Parasitic components on bitlines and wordlines are a major source of latency at scaled process technology nodes and they create significant **latency variation** in a crossbar [33, 34, 49, 53, 56, 67, 79, 83, 87]. Such variation can introduce ISI distortion (Section B), which may impact the quality of an inference task [9, 27].

To increase the energy efficiency of a neuromorphic system, Non-Volatile Memory (NVM) such as oxide-based random access memory (OxRRAM), phase-change memory (PCM), ferroelectric RAM (FeRAM), and spin-based magnetic RAM (MRAM) is used to implement the memory cells in a crossbar [15, 64, 94, 95, 98]. An NVM cell can be programmed to a high-resistance state (HRS) or one of many low-resistance states (LRS), implementing multi-bit synaptic weights [59, 64, 97]. To implement a synaptic weight on a memory cell of a crossbar, the synaptic weight is programmed as the conductance of the cell.

A crossbar can accommodate only a fixed number of pre-synaptic connections per post-synaptic neuron. To give an example, the crossbar in Figure 1 (left) has  $N$  pre-synaptic neurons,  $N$  post-synaptic neurons, and  $N^2$  memory cells. Each post-synaptic neuron can have a maximum of  $N$  pre-synaptic connections. To mitigate the negative impact of technology scaling, e.g., increase in the value of parasitic components on current paths and increase in the power density,  $N$  is constrained to a lower value, typically between 128 and 256 (see our tradeoff analysis in Section 2). To map a

large SNN model on a multi-PE hardware, a system software framework such as NEUTRAMS [50], NeuroXplorer [12], SentryOS [93], LAVA [60] and DFSynthesizer [82] is commonly used. These frameworks first partition a model into clusters, where a cluster is a subset of neurons and synapses of the model that can be implemented on the architecture of a crossbar. Subsequently, the partitioned clusters are implemented on different crossbars of a neuromorphic hardware.

We make the following two key **observations** related to a neuromorphic PE.

**Observation 1:** *The latency within a crossbar is a function of the length (i.e., the number of parasitic components) of current paths and the delay to sense the NVM cell activated on a current path.*

**Observation 2:** *Due to how memory cells are organized in a crossbar, a significant fraction of these memory cells remains unutilized when implementing machine learning inference tasks.*

Based on these two observations (which we elaborate in Sections 2-4), we present a design-technology tradeoff analysis to implement machine learning inference tasks on different PEs of an NVM-based neuromorphic system. We make the following four key **contributions**.

- Through detailed circuit-level simulations at scaled process technology nodes, we show that bitline and wordline parasitics are the primary sources of long latency in a crossbar and they create asymmetry in inference latency. With technology scaling, the absolute latency increases and the latency asymmetry becomes increasingly more significant. In addition, different resistance states of a multi-level NVM cell take varying latencies to sense during an inference operation (see Section 2).
- We propose to optimize the implementation of synaptic weights on NVM cells such that the resistance state that takes the longest time to sense is programmed on the NVM cell that has the least parasitic delay in a crossbar. This lowers the latency of a crossbar. (see Section 3)
- We propose an architectural change of introducing isolation transistors in a crossbar to partition it into regions that can be individually power-gated based on their utilization. In this way, we improve energy efficiency. In addition, by isolating the unutilized region of a crossbar from the active region, parasitics of only the active region contribute to latency, rather than both as in a baseline non-partitioned crossbar architecture. This reduces the latency of a crossbar (see Section 4).
- We show that our technological and architectural optimizations can only deliver on its latency and energy improvement promises if they are exploited efficiently by the system software. Therefore, we propose a mechanism to expose our proposed design changes to the system-level, allowing the system software to improve both latency and energy when implementing machine-learning inference tasks on hardware (see Section 5).

We evaluate our design-technology co-optimization approach for a recent neuromorphic hardware using 10 machine learning inference tasks. Results show 12% reduction in average PE latency and 22% lower application energy compared to current state-of-the-art.

To the best of our knowledge, this is the first work that demonstrates the energy and latency improvement of power gating crossbar-based neuromorphic hardware designs.

## 2 DESIGN-TECHNOLOGY TRADEOFF ANALYSIS

Without loss of generality, we demonstrate the design-technology tradeoff for an OxRRAM-based neuromorphic PE, where each NVM cell can be programmed to the following four resistance levels (i.e., 2-bit per synapse) – 1.5 K $\Omega$ , 5.78 K $\Omega$ , 13.6 K $\Omega$ , and 73 K $\Omega$  [19, 20, 31, 64, 74]. Furthermore, we show our analysis for four process technology nodes – 16nm, 22nm, 32nm, and 45nm, which are obtained from our technology provider. The analysis can be easily extended to other NVM types and also to other process technology nodes.

## 2.1 Cost-per-Bit Analysis for a Neuromorphic PE

The computer memory industry has thus far been primarily driven by the **cost-per-bit metric**, which provides the maximum capacity for a given manufacturing cost. As shown in recent works [56, 67, 68, 77, 79, 81, 83], manufacturing cost can be estimated from the area overhead. To estimate the cost-per-bit of a neuromorphic PE, we investigate the internal architecture of a crossbar and find that a neuron circuit can be designed using 20 transistors and a capacitor [48], while an NVM cell is a 1T-1R arrangement with a transistor used as an access device for the cell. Within an  $N \times N$  crossbar, there are  $N$  pre-synaptic neurons,  $N$  post-synaptic neurons, and  $N^2$  synaptic cells. The total area of all the neurons and synapses of a crossbar is

$$\begin{aligned} \text{neuron area} &= 2N(20T + 1C) \\ \text{synapse area} &= N^2(1T + 1R) \end{aligned} \quad (1)$$

where  $T$  stands for transistor,  $C$  for capacitor, and  $R$  for NVM cell. The total synaptic cell capacity is  $N^2$ , with each NVM cell implementing 2-bit per synapse. The total number of bits (i.e., synaptic capacity) in the crossbar is

$$\text{total bits} = 2N^2 \quad (2)$$

Therefore, the cost-per-bit of an  $N \times N$  crossbar is

$$\text{cost-per-bit} = \frac{2N(20T + 1C) + N^2(1T + 1R)}{2N^2} \approx \frac{F^2(27 + 2N)}{N}, \quad (3)$$

where the cost-per-bit is represented in terms of the crossbar dimension  $N$  and the feature size  $F$ . Equation 3 provides a **back-of-the-envelope** calculation of cost-per-bit. Figure 2 plots the normalized cost-per-bit for four different process technology nodes, with the crossbar dimension ranging from 16 to 256. We make the following two observations.

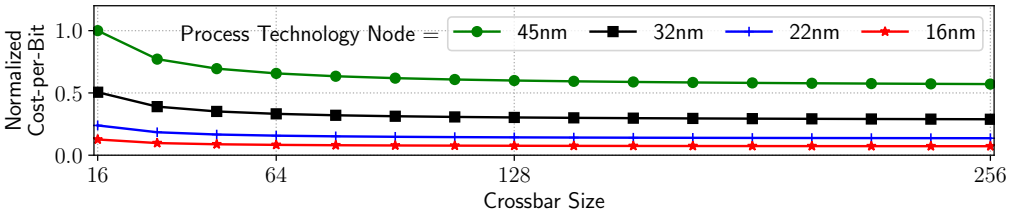


Fig. 2. Cost-per-bit analysis of a crossbar.

First, the cost-per-bit reduces with increase in the dimension of a crossbar, i.e., larger-sized crossbars can accommodate more bits for a given cost. However, both the absolute latency and latency variation increases significantly for larger-sized crossbars, which increases inference latency and reduces the quality of machine learning inference due to an increase in the ISI distortion (see our analysis in Section 2.2). Second, the cost-per-bit reduces considerably with technology scaling. This is due to higher integration density at smaller process technology nodes.

The formulation for the cost-per-bit (Equation 3) depends on the specific neuron architecture of [48] and the one transistor (1T)-based OxRRAM design of [64]. This formulation can be easily extended to other neuron and synapse designs. Furthermore, system designers can use our formulation to configure their neuromorphic hardware, without having to access and plug-in technology-related data.

## 2.2 Latency Variation in a Neuromorphic PE

Figure 3 shows the difference between the best-case and worst-case latency in a crossbar (expressed as a fraction of  $1\mu\text{s}$  spike duration) for five different crossbar configurations at 45nm, 32nm, 22nm, and 16nm process technology nodes. See our experimental setup using NeuroXplorer [12] in Section 6.1, which incorporates software, architecture, circuit, and technology. All NVM cells are programmed to the HRS state, i.e.,  $73\text{ K}\Omega$  (see Section 2.3 for the dependency on resistance states).

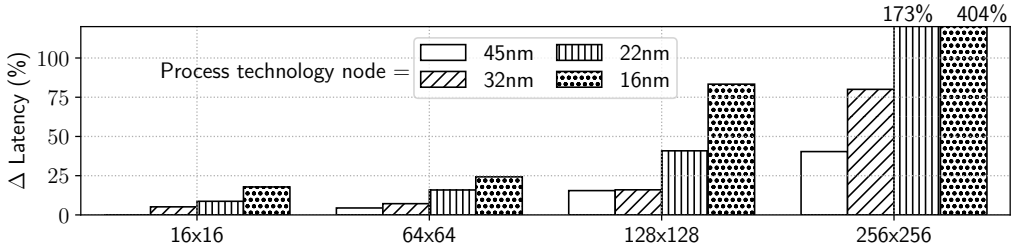


Fig. 3. Variance in latency within a crossbar, expressed as a fraction of a single spike duration.

We make two key observations. First, the latency difference increases with crossbar size due to an increase in the number of parasitic components on current paths. The average latency difference for  $256 \times 256$  crossbar is higher than  $16 \times 16$ ,  $64 \times 64$ , and  $128 \times 128$  crossbar by 16.5x, 13.4x, and 4.5x, respectively. This average is computed across the four process technology nodes. Therefore, smaller-sized crossbars lead to a smaller variation in latency, which is good for performance. However, smaller-sized crossbars also lead to higher cost-per-bit, which we have analyzed in Section 2.1. For most neuromorphic PE designs,  $128 \times 128$  crossbars achieve the best tradeoff in terms of latency variation and cost-per-bit [42, 57, 64, 66, 70, 102].

Second, the latency difference increases significantly for scaled process technology nodes due to an increase in the value of the parasitic component. The average latency difference for 32nm, 22nm, and 16nm process technology nodes is higher than 45nm by 1.3x, 3x, and 6.6x, respectively. The unit wordline (bitline) parasitic resistance ranges from approximately  $2.5\Omega$  ( $1\Omega$ ) at 45nm node to  $10\Omega$  ( $3.8\Omega$ ) at 16nm node. The value of these unit parasitic resistances is expected to scale further reaching  $\approx 25\Omega$  at 5nm node [33, 35, 36, 73, 92]. The unit wordline and bitline capacitance values also scale proportionately with technology. Latency variation increases ISI distortion, which degrades the quality of machine learning inference.

## 2.3 Varying Latency to Sense NVM Resistance States

The latency (i.e., the delay) on a current path from a pre-synaptic neuron to a post-synaptic neuron within a crossbar is proportion to  $R_{eff} \cdot C_{eff}$ , where  $R_{eff}$  ( $C_{eff}$ ) is the effective resistance (capacitance) on the path. This delay increases the time it takes for the membrane potential of a post-synaptic neuron to rise above the threshold voltage causing a delay in spike generation.

The effective resistance on a current path depends on the value of parasitic resistances and the resistance of the NVM cell. We analyze the latency impact due to different resistance states. Figure 4 plots the increase in latency (expressed as a fraction of  $1\mu\text{s}$  spike duration) to sense three NVM resistance states (LRS2, LRS3, and HRS) with respect to LRS1 at 45nm, 32nm, 22nm, and 16nm process technology nodes. These results are for a neuromorphic PE with a  $128 \times 128$  crossbar.

We observe that the latency to sense the HRS state is considerably higher than all three LRS states at all process technology nodes (consistent with [17, 64, 99]). The latency difference increases

with technology scaling due to an increase in the size of parasitic components on bitlines and wordlines of a crossbar, which we have analyzed in Section 2.2.

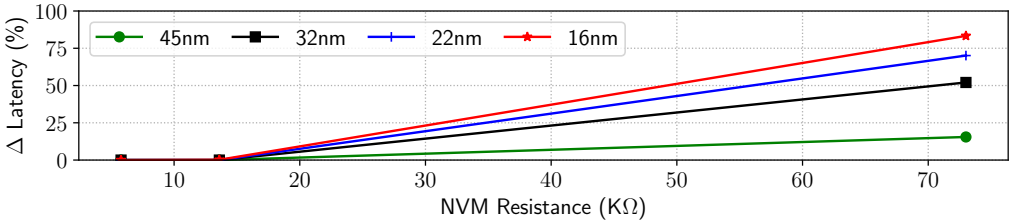


Fig. 4. Latency to sense various NVM resistance states, expressed as a fraction of a single spike duration.

### 3 PROPOSED TECHNOLOGICAL IMPROVEMENTS

Based on the design-technology tradeoff analysis of Section 2, we now present our technology-related optimization. Without loss of generality, we present our optimization for a  $128 \times 128$  crossbar-based neuromorphic hardware designed at 16nm node. We exploit the following two observations from Section 2: 1) HRS resistance state in an NVM cell takes higher latency to sense than LRS states, and 2) spike propagation latency in a crossbar depends on the number of parasitic components on its current path. The left sub-figure of Figure 5 shows the proposed technological changes. A crossbar is partitioned into three regions. The number of parasitic components on current paths in region A is considerably lower than in the rest of the crossbar. Therefore, all NVM cells in this region (4 in this example) implement only the HRS resistance state, which takes the longest time to sense. Conversely, NVM cells in region B have longer propagation delay due to higher number of parasitic components. Therefore, all NVM cells in this region (9 in this example) implement only the LRS resistance state, which takes the shortest time to sense. Finally, all other NVM cells (i.e., those in region C) are programmable, i.e., these cells can implement all four resistance states. The overall objective is to balance the latency on different current paths within a crossbar. This minimizes the latency variation in a crossbar, which reduces ISI distortion and improves the quality of machine learning inference tasks.

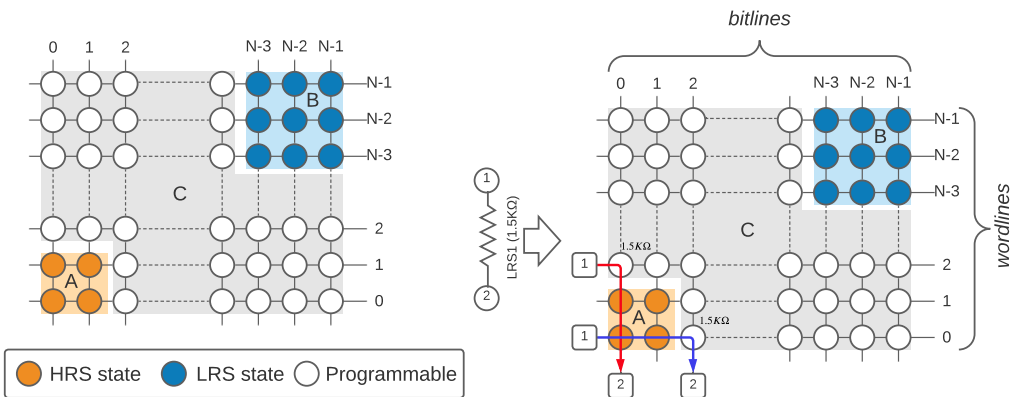


Fig. 5. Our proposed technological change.

The right sub-figure shows a pre- and a post-synaptic neuron connected via a synapse that is programmed to the LRS state. The synaptic connection can be implemented on NVM cells in

region B (with only LRS1 state) and region C (with programmable states). The figure illustrates two alternative implementations of these neurons. If the pre-synaptic neuron is implemented on wordline 0, then the post-synaptic neuron *cannot* be implemented on bitlines 0 and 1. This is because NVM cells in region A are all in HRS. In this example, we show the implementation on bitline 2 (see the blue implementation). Conversely, if the post-synaptic neuron is implemented on bitline 0, then the pre-synaptic neuron *cannot* be implemented on wordline 0 and 1 (to avoid using region A). We show the implementation on wordline 2 (see the red implementation).

Formally, the proposed neuromorphic PE is represented by a tuple  $\langle N, N_h, N_l \rangle$ , where  $N$  is the dimension of its crossbar. All NVM cells at crosspoints of wordlines  $0, 1, \dots, N_h - 1$  and bitlines  $0, 1, \dots, N_l - 1$  (i.e., region A) can implement only HRS. All NVM cells at crosspoints of wordlines  $N - N_l, N - N_l + 1, \dots, N - 1$  and bitlines  $N - N_l, N - N_l + 1, \dots, N - 1$  (i.e., region B) can implement only LRS. All other NVM cells in the PE's crossbar can implement all four resistance states.

Figure 6 plots the variation of latency in the proposed  $128 \times 128$  crossbar, normalized to a baseline architecture [9], where any NVM cell can be programmed to any of the four resistance states. See Section 6.4 for a description of this baseline architecture and Section 6.1 for the simulation setup. The variation in latency is measured as ratio of the best-case and worst-case latency in the crossbar. The figure reports latency variation for  $N_h$  ranging from 2 to 64 with  $N_l$  set to 16, 32, and 64.

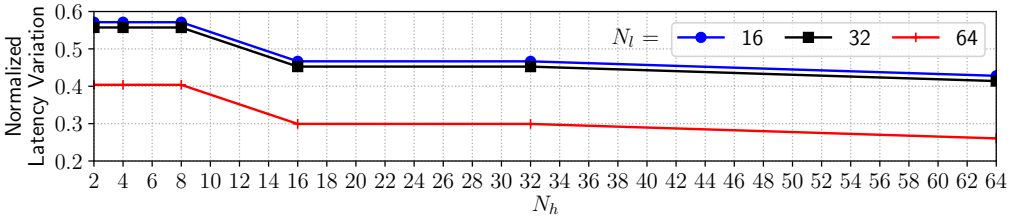


Fig. 6. Latency variation in the proposed crossbar architecture for different settings of  $N_h$  and  $N_l$ .

We observe that latency variation decreases with an increase in  $N_h$ . This is due to an increase in the size of region A, which increases the (worst-case) latency due to an increase in the number of parasitic components on current paths via the HRS state. However, the (best-case) latency of current paths via the LRS state remains the same. Therefore, the latency variation reduces which improves inference quality by lowering the ISI distortion. To illustrate this concept, Figure 7 provides an example where two synapses are mapped to a  $4 \times 4$  crossbar. In Figure 7a, the red synapse (in HRS state) is mapped to the bottom left corner of the crossbar, while the blue synapse (in LRS state) to the top right corner. The figure shows the timing of two spikes. The input spike on the red and blue synapses are  $t_1$  and  $t_2$ , respectively. Without loss of generality let  $t_2 > t_1$ . The ISI of these two spikes is  $t_2 - t_1$ . Due to the delay in current propagation through bitlines and wordlines, these two spikes arrive at the output terminal at different times – red synapse with a delay of  $x$  and blue synapse with a delay of  $y$ . Here,  $y > x$ . Therefore, ISI of the output spikes is  $(t_2 + y) - (t_1 + x)$ . The ISI distortion (difference of ISI between input and output) is

$$\text{ISI distortion} = \left( (t_2 + y) - (t_1 + x) \right) - (t_2 - t_1) = y - x \quad (4)$$

Figure 7b illustrates a scenario where region A is increased to include more cells that are programmed to the HRS state. The mapping process will map the red synapse using the farthest cell of region A. The delay on this synapse is  $x + \Delta$ , where  $\Delta$  is the additional delay due to routing

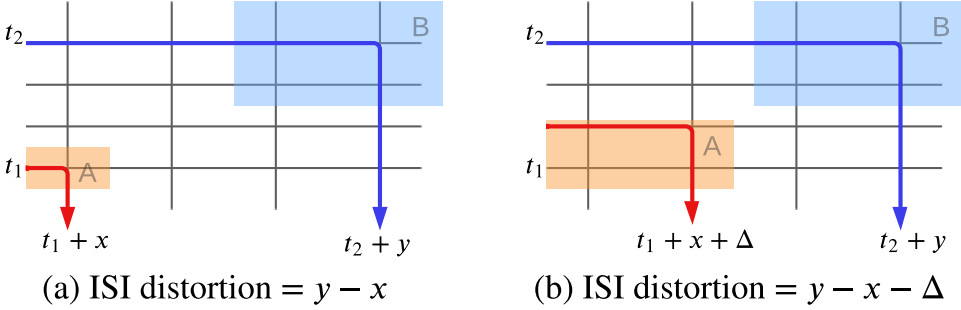


Fig. 7. ISI improvement due to increase in the size of region A.

spikes on the red synapse via a longer route compared to that in Figure 7a. Therefore, ISI of the output spikes is  $(t_2 + y) - (t_1 + x + \Delta)$ . The ISI distortion is

$$\text{ISI distortion} = \left( (t_2 + y) - (t_1 + x + \Delta) \right) - (t_2 - t_1) = y - x - \Delta \quad (5)$$

Comparing Equations 4 & 5, we observe that the ISI distortion reduces due to an increase in the size of region A. ISI distortion also reduces with an increase in  $N_I$  due to a reduction in the worst-case latency. We also note that, large  $N_h$  may lead to higher average crossbar latency, which impacts real-time performance. Finally, we see that going from  $N_I = 16$  to 32, there is no significant reduction in the latency variation. Although the size of region B increases with an increase in  $N_I$ , we observe only marginal reduction of the best-case latency. Overall, with  $N_h = N_I = 64$ , the latency variation is 74% lower than baseline. This is chosen based on the tradeoff between latency variation and average latency for  $128 \times 128$  crossbar at 16nm. The tradeoff point can change for other technology nodes and for other crossbar configurations.

### 3.1 Reduction in Latency Variation

To understand the reduction of latency variation within a crossbar as a result of our technological changes, we provide a simple example. Consider there are only two current paths in a crossbar. The parasitic delay on the shortest and longest current paths are  $D$  and  $(D + \Delta)$ , respectively. The time to sense LRS and HRS NVM states are  $S$  and  $(S + \delta)$ , respectively. Without any optimization, the worst-case condition is triggered when the HRS state is programmed on the longest path and the LRS state on the shortest path. The minimum and maximum latencies are  $(D + S)$  and  $(D + S + \Delta + \delta)$ , respectively. The latency variation is  $(\Delta + \delta)$ . Using our technology optimization, HRS state is programmed on the shortest path and LRS state on the longest path. The two latencies are  $(D + S + \Delta)$  and  $(D + S + \delta)$ . The latency variation reduces to  $(|\Delta - \delta|)$ .

Within a crossbar, there are many current paths ( $N^2$  current paths in a  $N \times N$  crossbar). The precise reduction in latency variation depends on the specific current paths activated for a synaptic connection, which is controlled during the mapping of a machine learning application to the crossbars of the hardware. In Figure 6, we show a 74% reduction comparing only the shortest and the longest paths in a  $128 \times 128$  crossbar. In Section 7.2, we evaluate the general case considering the mapping process. We report an average 22% reduction of latency variation.

Reducing the latency variation helps reduce the ISI distortion, which improves the inference quality. In Section 7.4, we report an average 4% increase of inference quality.



### 3.2 Impact on Latency

While latency variation impacts inference quality, the average crossbar latency impacts the real-time performance. To understand the impact of our technological optimization on the average crossbar latency, we consider the same example of two current paths. Consider there are  $m$  synapses with LRS states and  $n$  synapses with HRS state. The average latency in the worst-case condition is  $\frac{m \cdot (D+S) + n \cdot (D+S+\Delta+\delta)}{m+n}$ . Using the technological improvement, the average latency is  $\frac{m \cdot (D+S+\Delta) + n \cdot (D+S+\delta)}{m+n}$ . Therefore, the change in latency is  $(\frac{n-m}{n+m}) \Delta$ . This change in latency depends on 1) current paths activated in a crossbar and 2) the value of  $n$  and  $m$ , i.e., the number of synaptic connections with HRS and LRS states, respectively. In Section 7.3, we show an average 3% reduction of the average crossbar latency for all the evaluated applications.

## 4 ARCHITECTURAL ENHANCEMENTS TO NEUROMORPHIC PE

To understand the motivation of the proposed architectural changes, Figure 8 reports the average synapse utilization of  $128 \times 128$  crossbars in neuromorphic PEs for 10 machine learning models implemented using the spatial decomposition technique of [11], which is a **best-effort approach** to improve the utilization of crossbars in a neuromorphic hardware.

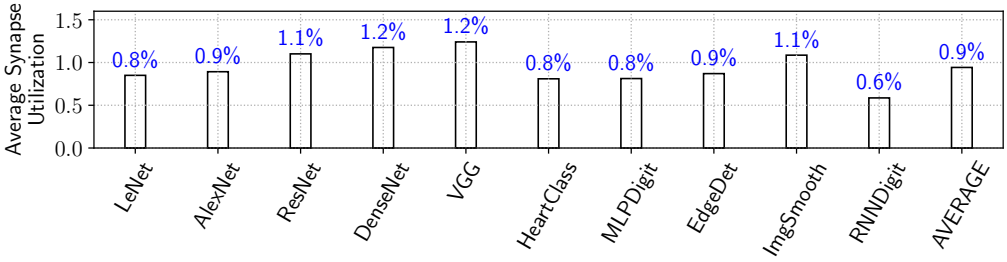


Fig. 8. Average synapse utilization of neuromorphic PEs.

We observe that the average synapse utilization is only 0.9%. This is because a crossbar can accommodate only a limited number of pre-synaptic connections per post-synaptic neuron. To illustrate this, Figure 9 shows three examples of implementing neurons on a  $4 \times 4$  crossbar. The synapse utilization of the three example scenarios are (a) 25% (4 out of 16), (b) 18.75% (3 out of 16), and (c) 25% (4 out of 16). As the crossbar dimension increases, the utilization drops significantly. For instance, if a  $128 \times 128$  crossbar is used to implement a single 128-input neuron, i.e., generalization of Fig. 9a, the utilization is only 0.78% (128 utilized synapses out of a total of  $128^2 = 16,384$  synapses). Lower synapse utilization leads to lower energy efficiency.

To improve energy efficiency, we propose to partition a neuromorphic PE into regions that can be dynamically power-gated based on its utilization for a given machine-learning inference task. Figure 10 shows the use of isolation transistors in a neuromorphic PE to partition a  $4 \times 4$  crossbar into active and unutilized regions. Figure 10a illustrates the implementation of only a single neuron function  $y_1$  in the crossbar. To improve energy efficiency, isolation transistors are needed on every bitline (between wordlines 3 and 4) and on every wordline (between bitlines 1 and 2). Figure 10b illustrates the implementation of two neuron functions  $y_1$  and  $y_2$  in the crossbar. In this scenario, isolation transistors are only needed on every wordline (between bitlines 2 and 3). To implement inference on a neuromorphic system, each crossbar may have different utilization of its memory cells. Therefore, to improve energy efficiency in every crossbar, isolation transistors are needed on every bitline (and between every pair of wordlines) and on every wordline (and between every

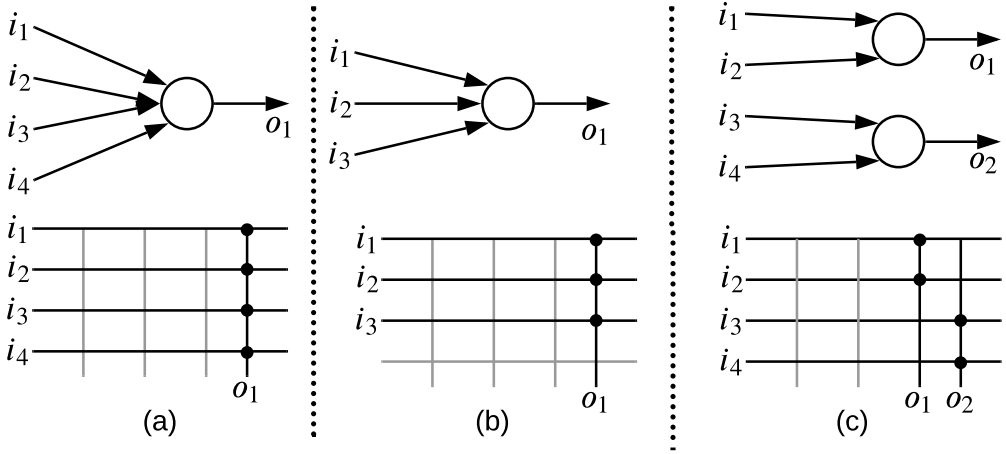


Fig. 9. Implementation of a) one 4-input, b) one 3-input, and c) two 2-input neurons to a  $4 \times 4$  crossbar.

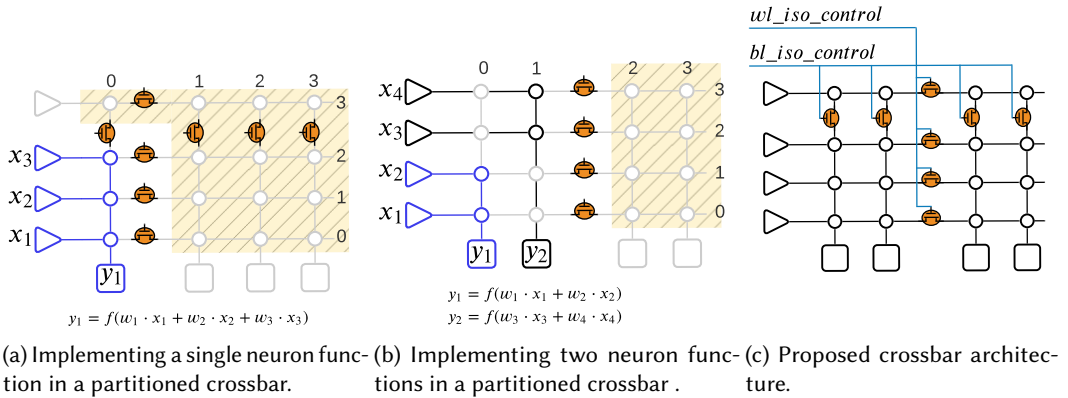


Fig. 10. Proposed neuromorphic PE architecture partitioned using isolation transistors.

pair of bitlines) – a total of 24 isolation transistors for this example  $4 \times 4$  crossbar (in general,  $2N(N-1)$  for an  $N \times N$  crossbar). This fine-grained partitioned PE architecture offers flexibility in energy management incorporating crossbar utilization, but leads to a significant increase in the area, latency, and system overhead to control isolation transistors.

To overcome these limitations while improving energy efficiency, we enable a coarse-grained partitioning in a crossbar as illustrated in Figure 10c. In this example, isolation transistors are inserted selectively on every bitline (between wordlines 3 and 4) and on every wordline (between bitlines 2 and 3). This coarse-grained partitioned PE architecture requires a total of 8 isolation transistors (in general,  $2N$  for an  $N \times N$  crossbar). To reduce the control overhead, isolation transistors on wordlines of a crossbar are controlled using a single control signal `wl_iso_ctrl` and those on bitlines using the signal `bl_iso_ctrl`. Through these two control signals, we enable four distinct configurations of the crossbar, which are summarized in Table 1.

Table 1. Different PE configurations enabled using the two new crossbar control signals.

Crossbar Control		Key Parameters		
wl_iso_ctrl	bl_iso_ctrl	Dimension	Energy	Latency
<b>Baseline PE Architecture</b>				
-	-	$4 \times 4$	$\propto 4^*4$	Best-case: $t_{1,1}$ Worst-case: $t_{4,4}$
<b>Proposed Partitioned PE Architecture</b>				
0	0	$3 \times 2$	$\propto 3^*2$	Best-case: $t_{1,1}$ Worst-case: $t_{3,2} - \Delta$
0	1	$4 \times 2$	$\propto 4^*2$	Best-case: $t_{1,1}$ Worst-case: $t_{4,2} - \Delta + t_{ON}$
1	0	$3 \times 4$	$\propto 3^*4$	Best-case: $t_{1,1}$ Worst-case: $t_{3,4} - \Delta + t_{ON}$
1	1	$4 \times 4$	$\propto 4^*4$	Best-case: $t_{1,1}$ Worst-case: $t_{4,4} + 2 \cdot t_{ON}$

In a baseline PE architecture, a crossbar dimension is fixed to  $4 \times 4$ . Its static energy is proportional to the number of memory cells, which is  $4^*4 = 16$  in this example. Latency in the crossbar varies from  $t_{1,1}$  (nearest cell or best-case) to  $t_{4,4}$  (farthest cell or worst-case).

In the proposed partitioned PE architecture, there are four configurations.

In configuration '00', the crossbar is configured as a  $3 \times 2$  array with its static energy proportional to  $3 \times 2 = 6$  memory cells. This is when the unutilized region is power-gated. The best-case latency is  $t_{1,1}$  and the worst-case latency is  $t_{3,2} - \Delta$ , where  $\Delta$  is the reduction in parasitic delay due to shorter bitlines and wordlines.

In configuration '01', the crossbar is configured as a  $4 \times 2$  array with its static energy proportional to  $4 \times 2 = 8$  memory cells. The best-case latency is  $t_{1,1}$  and worst-case latency is  $t_{4,2} - \Delta + t_{ON}$ , where  $t_{ON}$  is the delay of the isolation transistor on current paths.

In configuration '10', the crossbar is configured as a  $3 \times 4$  array with its static energy proportional to  $3 \times 4 = 12$  memory cells. The best-case latency is  $t_{1,1}$  and the worst-case latency is  $t_{3,4} - \Delta + t_{ON}$ .

In configuration '11', the crossbar is configured as the baseline  $4 \times 4$  array with its static energy proportional to  $4 \times 4 = 16$  memory cells. The best-case latency is  $t_{1,1}$  while the worst-case latency is  $t_{4,4} + 2 \cdot t_{ON}$ . Observe that on the longest current path there are now two isolation transistors, resulting in higher worst-case latency than in the baseline design.

Our proposed system software (which we discuss in Section 5) minimizes the use of configuration '11', improving both performance and energy efficiency.

**Single Control:** The proposed partitioned PE architecture also supports using a single control signal for all the isolation transistors in a crossbar. When using a single control, only the configurations '00' and '11' are used, implementing a  $3 \times 2$  and a  $4 \times 4$  array, respectively.

To generalize the discussion for an  $N \times N$  crossbar, assume that isolation transistors are inserted on every bitline (between wordlines  $P$  and  $P + 1$ ) and on every wordline (between bitlines  $Q$  and  $Q + 1$ ). Then, the four configurations are: '00': a  $P \times Q$  array; '01': a  $N \times Q$  array; '10': a  $P \times N$  array; and '11': a  $N \times N$  array. Formally,  $\langle N, N_h, N_l, P, Q \rangle$  represents the proposed partitioned PE architecture. Equation 6 summarizes the notations.

$$\begin{aligned}
\langle N \rangle &= \text{a baseline } N \times N \text{ crossbar} \\
\langle N, N_h, N_l \rangle &= N \times N \text{ crossbar with tech. enhancement (Sec. 3)} \\
\langle N, N_h, N_l, P, Q \rangle &= N \times N \text{ crossbar with tech. \& arch. enhancements} \\
&\quad \text{(see Sec. 3 \& 4)}
\end{aligned} \tag{6}$$

We introduce the following four terminologies: 1) **expanded mode**: in this mode, a crossbar is operated in configuration ‘11’, 2) **collapsed mode**: in this mode, a crossbar is operated in configurations ‘00’, ‘01’, and ‘10’, 3) **collapsed region**, this is the reduced dimension of the crossbar when operating in configurations ‘00’, ‘01’, and ‘10’, and 4) **far region**, this is the region of the crossbar excluding the collapsed region.

In our design methodology, the far region of a crossbar is power-gated using the two control signals at design-time considering the crossbar’s utilization. This is achieved during mapping of neurons and synapses to the hardware. Since neuron and synapse mapping does not change during inference, there is no dynamic power management needed. Consequently, there is also no latency and energy overhead involved in switching the far region on/off at run-time.

#### 4.1 Placing Isolation Transistors in a Crossbar

To illustrate the design space exploration involved in placing isolation transistors in a crossbar, Figure 11(a) illustrates a baseline crossbar with four current path that are activated during mapping of neurons and synapses. Figures 11(b)-(d) show three alternative placements of isolation transistors in the crossbar. In Figure 11(b), P and Q values are kept small. The size of the far region is large. In this figure, only two of the current paths (1 & 2) stays within the collapsed region of the crossbar, while the other two current paths (3 & 4) traverse via the far region. This means that the latency of paths 3 & 4 increases due to the delay of the isolation transistors on current paths. Additionally, the far region cannot be power gated, so there is a limited scope for energy reduction using power gating. Increasing P and Q values further (Figure 11(c)), the far region reduces in size as illustrated in the figure. Although three of the four current paths stay in the collapsed region, the far region still cannot be power-gated due to the presence of path 4 in this region. Finally, Figure 11(d) illustrates a possibility where all current paths stay in the collapsed region. The far region can therefore be power-gated. However, because of the small size of the far region, the energy benefits may not be significant. We explore this latency and energy tradeoffs.

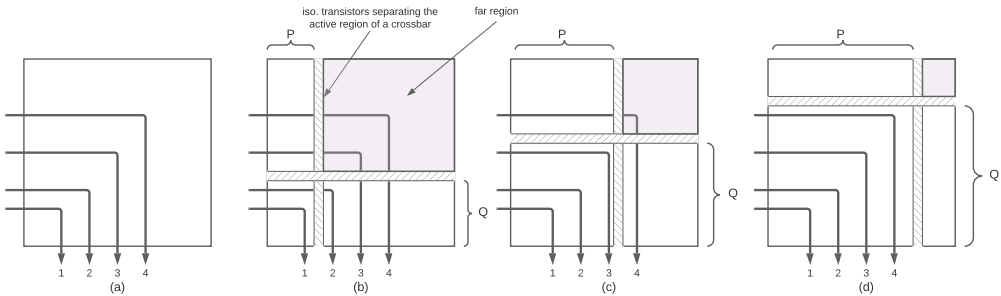


Fig. 11. Placing isolation transistors in a crossbar.

Figure 12 shows the latency and energy tradeoffs in selecting the values of  $P$  and  $Q$  for the ResNet inference workload implemented on  $128 \times 128$  crossbars in a neuromorphic hardware. Latency and energy numbers are normalized to baseline. We make the following two key observations.

First, energy is lower for smaller  $P$  and  $Q$  values. This is because by reducing  $P$  and  $Q$ , the size of the collapsed region of a crossbar reduces. Therefore, there are more memory cells in the far region that can be power-gated to lower energy.

Second, latency also reduces with a reduction in  $P$  and  $Q$  values (until  $P = Q = 80$ ). This is due to shorter bitlines and wordlines of the collapsed region. However, with  $P = Q = 64$  or  $72$ , more clusters of ResNet need crossbars in the expanded mode of operation. This is because synapses in these clusters can no longer fit onto the reduced dimension of a collapsed crossbar. This increases latency due to isolation transistors on current paths. For ResNet,  $P = Q = 80$  is the tradeoff point. The tradeoff point is different for different applications. To select a single crossbar configuration that gives good results for all applications, we perform similar analysis for all evaluated applications (see Section 6.3). Based on such analysis,  $P = Q = 96$  is the selected configuration for the  $128 \times 128$  crossbar at 16nm technology node.

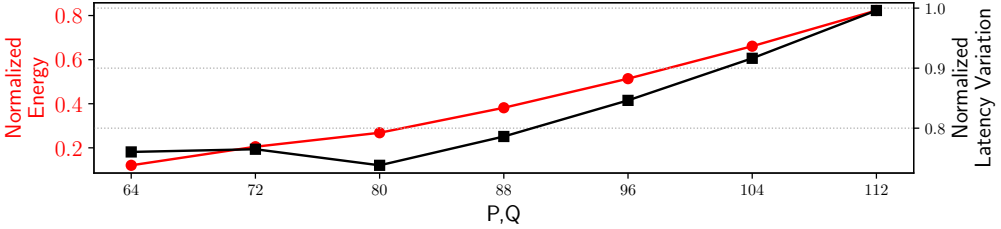


Fig. 12. Selecting  $P$  and  $Q$  values for the ResNet application.

## 5 EXPLOITING TECHNOLOGICAL AND ARCHITECTURAL IMPROVEMENTS VIA THE SYSTEM SOFTWARE

To describe the system software, the left subfigure of Figure 13 shows the final crossbar design with isolation transistors that allow each neuromorphic PE to operate in a collapsed or expanded mode. The right subfigure shows control signals for these transistors generated from a centralized controller implemented inside the system software.

Without loss of generality, Figure 14 shows modifications to the baseline system software [60] to exploit the proposed design changes. A trained machine learning model is first partitioned to generate clusters, where each cluster can fit onto a crossbar. These clusters are stored in a cluster queue (c1Q). In the baseline design, each cluster from the c1Q is mapped to an  $N \times N$  array (exactly replicating the crossbar dimension of the hardware). The mapping is programmed to the hardware using the cluster placement block. In the proposed design, each cluster of c1Q is mapped on four separate arrays – a  $P \times Q$  array, a  $N \times Q$  array, a  $P \times N$  array, and a  $N \times N$  array. These mappings go to a configuration selection block, which selects the final mapping for the cluster and the configuration of the corresponding PE based on energy-latency tradeoffs. The configuration is programmed to the hardware by configuring the two control signals `wl_iso_ctrl` and `bl_iso_ctrl`. This allows to power-gate the far region of the crossbar. It is important to note that since we power-gate unused resources of a crossbar only at design-time when admitting an application, we minimize the switching overhead. In the future, we will extend this work to also consider dynamic power management by dynamically controlling the isolation transistors.

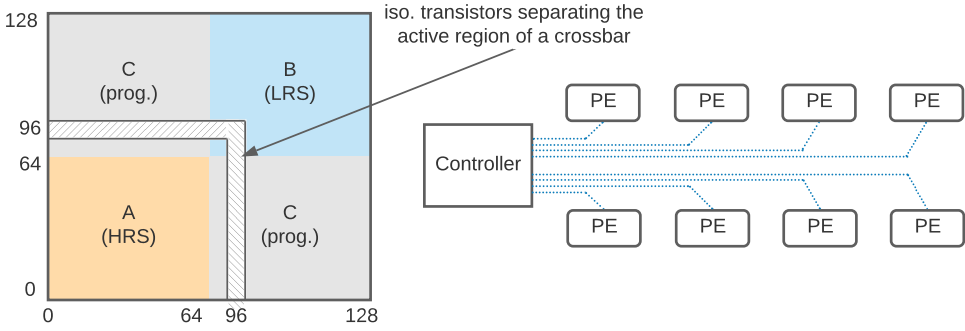


Fig. 13. Final crossbar design using the isolation transistors. The right subfigure shows the control signals generated from the controller when using the proposed partitioned PE architecture in a neuromorphic system.

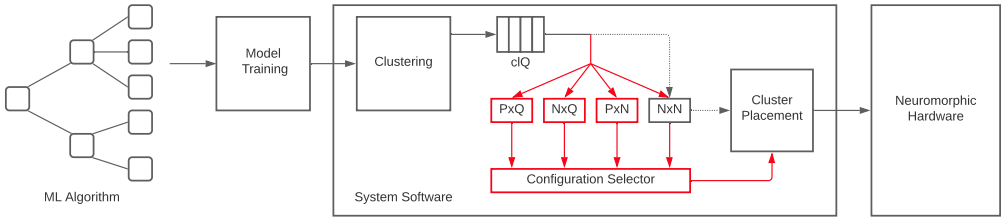


Fig. 14. Proposed system software. All changes are indicated in red.

In selecting the final mapping, the configuration selector first checks to see if a cluster can be mapped to a  $P \times Q$  array. If this is possible, then the mapping to the  $P \times Q$  array is selected as the final mapping for the cluster, and the corresponding PE is set to operate in configuration ‘00’ (collapsed mode). Otherwise, the configuration selector checks to see if the cluster can be mapped to  $N \times Q$  or  $P \times N$  array. If so, the corresponding mapping is selected, and the PE is set to operate in configurations ‘01’ or ‘10’, respectively. If the cluster cannot be mapped to either  $N \times Q$  or  $P \times N$  arrays, the mapping to  $N \times N$  array is selected as the final mapping of the cluster with the PE set to operate in configuration ‘11’ (expanded mode). In this way, the proposed system software uses expanded mode only when it is absolutely necessary to do so. Otherwise, it selects the collapsed region to map synapses, improving both latency and energy.

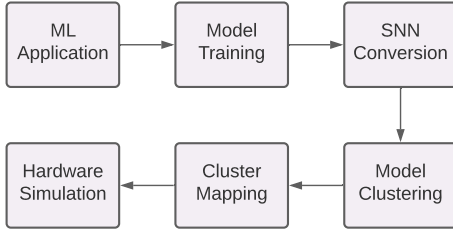
## 6 EVALUATION METHODOLOGY

### 6.1 Simulation Framework

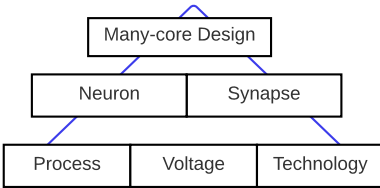
We evaluate the proposed design-technology co-optimization approach for OxRRAM-based neuromorphic PEs. Our simulation framework includes NeuroXplorer [12], a cycle-level in-house neuromorphic simulator [12] with programmable crossbar parameters. We configure this framework to simulate crossbars with parameters listed in Table 2. Circuit-level simulations are performed with technology parameters from the predictive technology model (PTM) [10] and OxRRAM-specific parameters from [17]. We note that, comparing different chip technologies or recommending one technology node over another is not the focus of this work. Instead, we show that for a given process technology node, design optimizations can reduce energy and latency variations. Furthermore, the proposed design-technology co-optimization methodology can be used by system designers to choose the best technology node for their neuromorphic designs by exploring the energy-performance tradeoffs.

Table 2. Major simulation parameters extracted from [28].

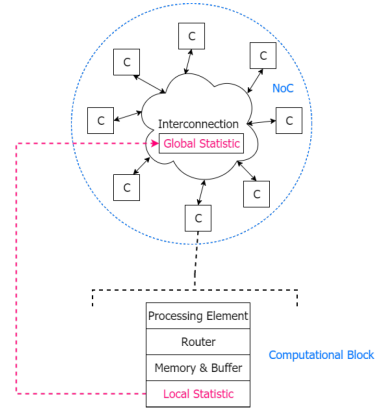
Neuron technology	16nm CMOS (original design is at 14nm FinFET)
Synapse technology	HfO <sub>2</sub> -based OxRRAM [64]
Supply voltage	1.0V
Energy per spike	23.6pJ at 30Hz spike frequency
Energy per routing	3pJ
Switch bandwidth	3.44 G. Events/s



(a) Design Pipeline



(b) Modeling Hierarchy



(c) Statistics Collection Framework

Fig. 15. Design pipeline using NeuroXplorer.

Neuromorphic simulations are performed on a Lambda workstation, which has AMD Threadripper 3960X with 24 cores, 128 MB cache, 128 GB RAM, and 2 RTX3090 GPUs. Figure 15(a) shows the design pipeline implemented using NeuroXplorer. A machine learning model is first trained using frameworks such as Keras and PyTorch. Subsequently, the trained model is converted into SNN using [5, 82]. The trained model is also simulated using an SNN simulator such as CARLsim [21]. NeuroXplorer integrates PyCARL [8], which allows the SNN model to be simulated using other SNN simulators such as Nengo [13], Neuron [43], and Brian [39]. Keras [41] and CARLsim [21] both use the two GPUs to accelerate model training and SNN functional simulation, respectively.

The SNN simulated model is clustered using the best-effort technique of [11], which maximizes cluster utilization. Clusters of the SNN are mapped to the hardware using the SpineMap technique [9]. Finally, we perform cycle-accurate simulation of the clusters using NeuroXplorer [12].

Figure 15(b) shows the modeling hierarchy of the simulator. At the highest level is the many-core design, which is a tile-based architecture, similar to Loihi [28]. Each PE consists of a crossbar, which is an organization of neurons and synapses. A neuron is modeled using [48] and a synaptic circuit using [64]. At the lowest level are the technology models (see Table 2).

Finally, Figure 15(c) shows the statistics collection framework in NeuroXplorer. It facilitates global statistics collection, where spike arrival times are recorded for each PE (shown as C in the figure). These spike times are then used to compute the ISI distortion (see Appendix B).

## 6.2 Power Consideration for Isolation Transistors

The additional power required to control the isolation transistors when accessing the RRAM cells in the far region is approximately 3x that of raising a wordline, since raising a wordline requires driving one access transistor per bitline, while accessing the RRAM cells in the far region requires driving two isolation and one access transistor per bitline. The power overhead for accessing RRAM cells in the collapsed mode ‘01’ and ‘10’ is approximately 2x (one isolation and one access transistor) [56, 79, 83]. The energy numbers reported in Section 7.1 incorporates these overheads.

## 6.3 Evaluated Workloads

We select 10 machine learning inference programs that are representative of three most commonly-used neural network classes: convolutional neural network (CNN), multi-layer perceptron (MLP), and recurrent neural network (RNN). Table 3 summarizes the topology, number of neurons and synapses, number of spikes per image, and baseline quality of these applications on hardware.

Table 3. Applications used to evaluate the proposed approach.

Class	Applications	Dataset				Baseline	Obtained
			Neurons	Synapses	Avg. Spikes/Frame	Quality	Quality
CNN	LeNet	CIFAR-10	80,271	275,110	724,565	86.3%	87.1%
	AlexNet	CIFAR-10	127,894	3,873,222	7,055,109	66.4%	66.9%
	ResNet	CIFAR-10	266,799	5,391,616	7,339,322	57.4%	58.0%
	DenseNet	CIFAR-10	365,200	11,198,470	1,250,976	46.3%	46.5%
	VGG	CIFAR-10	448,484	22,215,209	12,826,673	81.4 %	81.6%
	HeartClass [26]	Physionet	170,292	1,049,249	2,771,634	63.7%	63.9%
MLP	MLPDigit	MNIST	894	79,400	26,563	91.6%	96.4%
	EdgeDet [21]	CARLsim	7,268	114,057	248,603	SSIM = 0.89	0.99
	ImgSmooth [21]	CARLsim	5,120	9,025	174,872	PSNR = 19	22.2
RNN	RNNDigit [30]	MNIST	1,191	11,442	30,508	83.6%	83.7%

## 6.4 Evaluated Approaches

We evaluate the following techniques.

- *Baseline* [9]. The Baseline approach first clusters a machine-learning inference model to minimize the inter-cluster spike communication. Clusters are then mapped to neuromorphic PEs of the hardware with synapses of each cluster implemented on memory cells of a crossbar without incorporating latency variation. Neuromorphic PEs are not optimized to reduce latency variation, i.e., any resistance state (LRS or HRS) can be programmed on any current path (long or short). Unused crossbars are power-gated to reduce energy consumption. This is the coarse-grained power management technique implemented in many state-of-the-art many-core neuromorphic designs such as Loihi [28], DYNAPs [66], and  $\mu$ Brain [93].
- *Baseline + Design Changes*. This is the Baseline mapping approach implemented on the proposed latency-optimized partitioned neuromorphic PE design. In the proposed design, HRS state, which takes long time to sense, is used only on shorter current paths, ones that have lower parasitic delays. Similarly, LRS state is used only on loner current paths. In addition to coarse-grained power management, we facilitate power gating at a finer granularity in the proposed design. Specifically, by controlling the isolation transistors, we power-gate unused resources within each crossbar.



- *Proposed*. This is the proposed solution where the system software is optimized to exploit the design changes.

## 7 RESULTS AND DISCUSSIONS

### 7.1 Energy Efficiency

Figure 16 plots the energy efficiency of the evaluated techniques normalized to Baseline. We make the following two key observations.

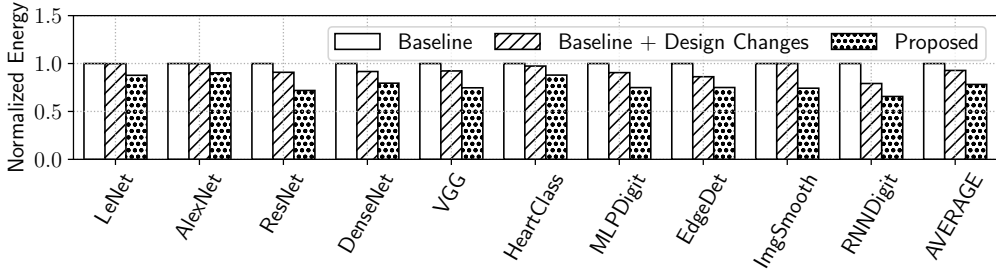


Fig. 16. Energy consumption normalized to Baseline.

First, with the proposed design changes, energy reduces by only 7% compared to Baseline. This is because, both in Baseline and Baseline with the proposed design changes, synapses of a cluster are implemented randomly on NVM cells of a crossbar causing them to be distributed across the crossbar dimension. Therefore, there remains a limited scope to collapse the crossbar and use power-gating to save energy. Second, the proposed design-technology co-optimization approach has the lowest energy (22% lower than Baseline and 16% lower than Baseline with the proposed design changes). This improvement is due to the proposed system software, which exploits the design changes in implementing machine learning inference on neuromorphic PEs. In particular, synapses are implemented to maximize the utilization of the collapsed region in each crossbar of the hardware. If all of a cluster's synapses fit into the collapsed region, then the far region can be isolated from the collapsed region using isolation transistors and power-gated to save energy.

### 7.2 Latency Variation

Figure 17 plots the latency variation normalized to Baseline. We make the following three key observations.

First, with the proposed design changes, latency variation increases compared to Baseline by average 1%. This is because of the increase in latency associated with the delay of isolation transistors on current paths. Second, the latency variation using the proposed approach is 30% lower than Baseline and 32% lower than Baseline with the proposed design changes. The reason for these improvements is three fold – 1) optimizing NVM resistance states in a crossbar such that the state that takes the longest time to sense is programmed on current paths that have the least propagation delay, 2) isolating the collapsed region of a crossbar from the far region, to reduce current propagation delay, and 3) exploiting these changes during the implementation of a machine learning inference using the proposed system software, which uses the far region of a crossbar only when it is absolutely necessary to do so. Otherwise, it improves both latency and energy by operating the crossbar in the collapsed mode.

Finally, the latency variation using the proposed approach varies across different applications. This is because the proposed approach exploits the latency and energy tradeoffs differently for

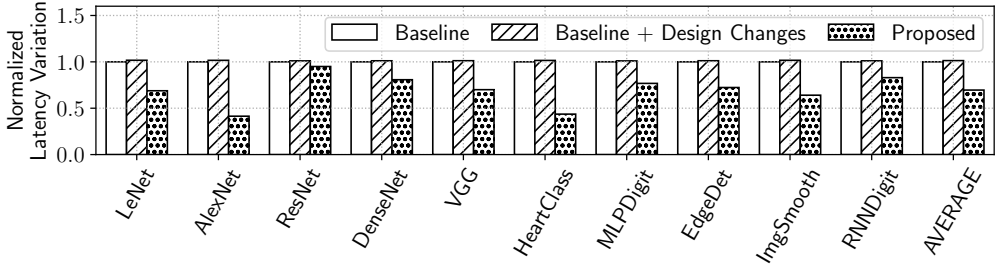


Fig. 17. Latency variation normalized to Baseline.

different applications. The latency variation is similar to the Baseline for ResNet, while it is significantly lower than the Baseline for HeartClass.

Using the results from Sections 7.1 and 7.2, we conclude that the proposed approach introduces maximum gain for applications where the latency and energy tradeoffs can be better exploited. For all other applications, it either minimizes energy or minimizes latency variation.

### 7.3 Real-time Performance

One of the key hardware performance metrics for neuromorphic computing is real-time performance, which is a function of the crossbar latency. To evaluate real-time performance, Figure 18 plots the crossbar latency of the proposed approach and the Baseline for the evaluated applications. Results are normalized to the Baseline.

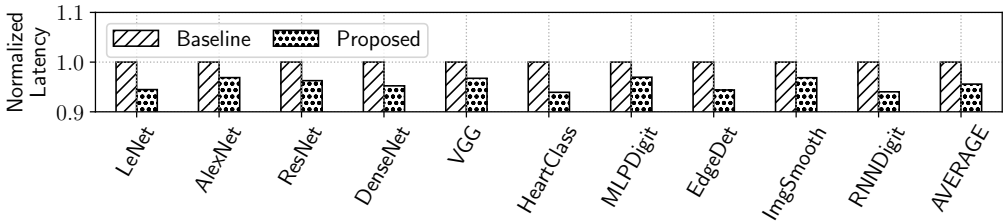


Fig. 18. Crossbar latency normalized to Baseline.

We observe that the crossbar latency using the proposed approach is on average 4.5% lower than the Baseline. This reduction is because the proposed approach places synapses with the HRS state on shorter current paths, which lowers the overall spike latency on those synapses. We have elaborated this in Section 3.2.

### 7.4 Inference Quality

Figure 19 shows the improvement in inference quality using the proposed approach, normalized to Baseline. We observe that the image quality improves by an average of 4%. This is due to the reduction in ISI distortion caused by a reduction of the latency variation in neuromorphic PEs using the proposed changes, which we have analyzed in Section 7.2. In addition, the improvement of inference quality with PSNR and SSIM metrics for EdgeDet and ImgSmooth is higher than other inference tasks with accuracy metrics. This is because PSNR and SSIM metrics are computed on

individual images where we see a large improvement in quality. For accuracy-based tasks, we observe that feature representation in hidden layers of these models changes due to ISI distortion, but not all such changes lead to misclassification. So the accuracy of these inference tasks is comparable to Baseline.

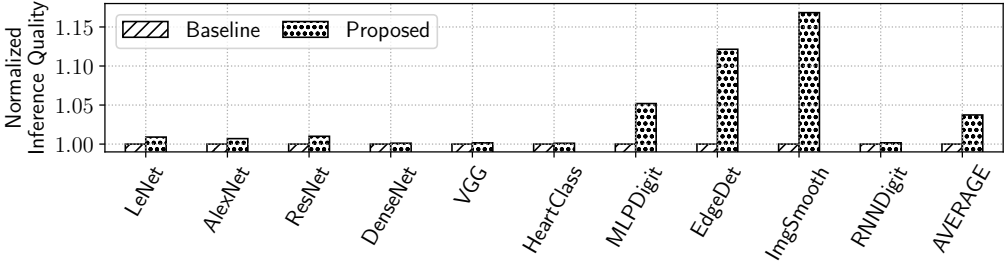


Fig. 19. Inference quality normalized to Baseline.

### 7.5 Single vs. Double Control Design

Figure 20 plots the energy efficiency of the proposed design with single control signal and the default, which uses two control signals for each PE. We observe that using single control, energy reduces by only 2% compared to Baseline. This is because most crossbars are operated in the expanded mode due to limited scope to collapse the crossbar. Our default design leads to 14.4% lower energy than with single control. This is because in the default design, a crossbar can be collapsed along X- and Y- dimensions independently, leading to three collapsed array configurations. Therefore, the system software has a higher probability to use the collapsed mode, leading to a reduction in energy.

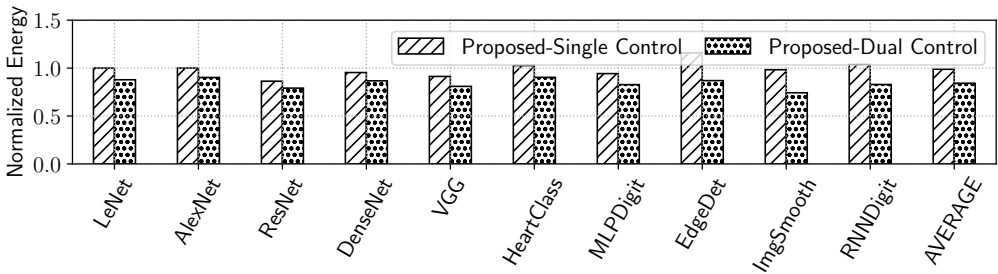


Fig. 20. Partitioned PE architecture with single and double control.

### 7.6 Die Area Analysis

Adding an isolation transistor to the bitline increases the height of the crossbar, whereas that on the wordline increases the width. Without the isolation transistors, the height of a baseline crossbar is equal to the sum of height of the memory cells and the sense-amplifier, while the width is equal to the sum of the width of the memory cells. For RRAM-based neuromorphic PEs, a sense amplifier in the peripheral circuit and a isolation transistor is approximately 384x and 9.6x taller

than an individual RRAM cell, respectively [18, 64, 96]. In terms of width, an isolation transistor is only 1.3x wider than an RRAM cell. Therefore, for a crossbar with 128 RRAM cells per bitline and wordline (i.e.,  $128 \times 128$  array), the overhead along the height of the crossbar is  $\frac{9.6}{384+128} = 1.83\%$ , and the overhead along the width of the crossbar is  $\frac{1.3}{128} = 1.01\%$ .

## 8 CONCLUSIONS

We present a design-technology co-optimization approach to implement energy-efficient machine-learning inference on NVM-based neuromorphic processing elements (PEs). First, we optimize the NVM resistance state such that the state that takes the longest time to sense is placed on current paths with fewer parasitics, and hence incurs lower propagation delay, and vice versa. Second, we use isolation transistors to partition a PE into collapsed and far regions such that the NVM cells of the far region can be opportunistically power-gated to save both energy and latency. Finally, we use the system software to exploit the design changes, maximizing the utilization of the collapsed region of each PE in the hardware. Our system software uses the far region only when it is absolutely necessary to do so, otherwise it improves both latency and energy by operating the PE in the collapsed mode. We evaluate our design-technology co-optimization approach for a state-of-the-art neuromorphic architecture. Evaluations with different machine-learning inference tasks show that the proposed approach improves both latency and energy without incurring significant cost-per-bit.

## ACKNOWLEDGMENTS

This work is supported by 1) U.S. Department of Energy under Award Number DE-SC0022014, 2) the National Science Foundation Award CCF-1937419 (RTML: Small: Design of System Software to Facilitate Real-Time Neuromorphic Computing) and 3) the National Science Foundation Faculty Early Career Development Award CCF-1942697 (CAREER: Facilitating Dependable Neuromorphic Computing: Vision, Architecture, and Impact on Programmability).

## REFERENCES

- [1] A. Amir, P. Datta, W. P. Risk, A. S. Cassidy, J. A. Kusnitz, S. K. Esser, A. Andreopoulos, T. M. Wong, M. Flickner, R. Alvarez-Icaza *et al.*, "Cognitive computing programming paradigm: a corelet language for composing networks of neurosynaptic cores," in *IJCNN*, 2013.
- [2] A. Ankit, A. Sengupta, and K. Roy, "TraNNsformer: Neural network transformation for memristive crossbar based neuromorphic system design," in *ICCAD*, 2017.
- [3] A. Balaji and A. Das, "Compiling spiking neural networks to mitigate neuromorphic hardware constraints," in *IGSC Workshops*, 2020.
- [4] A. Balaji and A. Das, "A framework for the analysis of throughput-constraints of SNNs on neuromorphic hardware," in *ISVLSI*, 2019.
- [5] A. Balaji, F. Corradi, A. Das, S. Pande, S. Schaafsma, and F. Catthoor, "Power-accuracy trade-offs for heartbeat classification on neural networks hardware," *JOLPE*, 2018.
- [6] A. Balaji, S. Song, A. Das, N. Dutt, J. Krichmar, N. Kandasamy, and F. Catthoor, "A framework to explore workload-specific performance and lifetime trade-offs in neuromorphic computing," *CAL*, 2019.
- [7] A. Balaji, Y. Wu, A. Das, F. Catthoor, and S. Schaafsma, "Exploration of segmented bus as scalable global interconnect for neuromorphic computing," in *GLSVLSI*, 2019.
- [8] A. Balaji, P. Adiraju, H. J. Kashyap, A. Das, J. L. Krichmar, N. D. Dutt, and F. Catthoor, "PyCARL: A PyNN interface for hardware-software co-simulation of spiking neural network," in *IJCNN*, 2020.
- [9] A. Balaji, A. Das, Y. Wu, K. Huynh, F. G. Dell'anna, G. Indiveri, J. L. Krichmar, N. D. Dutt, S. Schaafsma, and F. Catthoor, "Mapping spiking neural networks to neuromorphic hardware," *TVLSI*, 2020.
- [10] A. Balaji, T. Marty, A. Das, and F. Catthoor, "Run-time mapping of spiking neural networks to neuromorphic hardware," *JSPS*, 2020.
- [11] A. Balaji, S. Song, A. Das, J. Krichmar, N. Dutt, J. Shackleford, N. Kandasamy, and F. Catthoor, "Enabling resource-aware mapping of spiking neural networks via spatial decomposition," *ESL*, 2020.

- [12] A. Balaji, S. Song, T. Titirsha, A. Das, J. Krichmar, N. Dutt, J. Shackelford, N. Kandasamy, and F. Catthoor, "NeuroXplorer 1.0: An extensible framework for architectural exploration with spiking neural networks," in *ICONS*, 2021.
- [13] T. Bekolay, J. Bergstra, E. Hunsberger, T. DeWolf, T. C. Stewart, D. Rasmussen, X. Choo, A. Voelker, and C. Eliasmith, "Nengo: a Python tool for building large-scale functional brain models," *Frontiers in Neuroinformatics*, 2014.
- [14] S. Bose, J. Acharya, and A. Basu, "Is my neural network neuromorphic? taxonomy, recent trends and future directions in neuromorphic engineering," *ACSSC*, 2019.
- [15] G. W. Burr, R. M. Shelby, A. Sebastian, S. Kim, S. Kim, S. Sidler, K. Virwani, M. Ishii, P. Narayanan, A. Fumarola, L. L. Sanches, I. Boybat, M. Le Gallo, K. Moon, J. Woo, H. Hwang, and Y. Leblebici, "Neuromorphic computing using non-volatile memory," *Advances in Physics: X*, 2017.
- [16] F. Catthoor, S. Mitra, A. Das, and S. Schaafsma, "Very large-scale neuromorphic systems for biological signal processing," in *CMOS Circuits for Biological Sensing and Processing*, 2018.
- [17] P.-Y. Chen and S. Yu, "Compact modeling of RRAM devices and its applications in 1T1R and 1S1R array design," *TED*, 2015.
- [18] P.-Y. Chen, Z. Li, and S. Yu, "Design tradeoffs of vertical RRAM-based 3-D cross-point array," *TVLSI*, 2016.
- [19] Y. Chen, "ReRAM: History, status, and future," *TED*, 2020.
- [20] Y.-H. Chiu, Y.-B. Liao, M.-H. Chiang, C.-L. Lin, W.-C. Hsu, P.-C. Chiang, Y.-Y. Hsu, W.-H. Liu, S.-S. Sheu, K.-L. Su *et al.*, "Impact of resistance drift on multilevel PCM design," in *ICDT*, 2010.
- [21] T. Chou, H. Kashyap, J. Xing, S. Listopad, E. Rounds, M. Beyeler, N. Dutt, and J. Krichmar, "CARLsim 4: An open source library for large scale, biologically detailed spiking neural network simulation using heterogeneous clusters," in *IJCNN*, 2018.
- [22] D. V. Christensen, R. Dittmann, B. Linares-Barranco, A. Sebastian, M. L. Gallo, A. Redaelli, S. Slesazek, T. Mikolajick, S. Spiga, S. Menzel *et al.*, "2021 roadmap on neuromorphic computing and engineering," *arXiv*, 2021.
- [23] S. Curzel, N. B. Agostini, S. Song, I. Dagli, A. Limaye, C. Tan, M. Minutoli, V. G. Castellana, V. Amatya, J. Manzano *et al.*, "Automated generation of integrated digital and spiking neuromorphic machine learning accelerators," in *ICCAD*, 2021.
- [24] A. Das, P. Pradhapan, W. Groenendaal, P. Adiraju, R. Rajan, F. Catthoor, S. Schaafsma, J. Krichmar, N. Dutt, and C. Van Hoof, "Unsupervised heart-rate estimation in wearables with Liquid states and a probabilistic readout," *Neural Networks*, 2018.
- [25] A. Das and A. Kumar, "Dataflow-based mapping of spiking neural networks on neuromorphic hardware," in *GLSVLSI*, 2018.
- [26] A. Das, F. Catthoor, and S. Schaafsma, "Heartbeat classification in wearables using multi-layer perceptron and time-frequency joint distribution of ECG," in *CHASE*, 2018.
- [27] A. Das, Y. Wu, K. Huynh, F. Dell'Anna, F. Catthoor, and S. Schaafsma, "Mapping of local and global synapses on spiking neuromorphic hardware," in *DATE*, 2018.
- [28] M. Davies, N. Srinivasa, T. H. Lin, G. China, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain, Y. Liao, C. K. Lin, A. Lines, R. Liu, D. Mathaikutty, S. McCoy, A. Paul, J. Tse, G. Venkataramanan, Y. H. Weng, A. Wild, Y. Yang, and H. Wang, "Loihi: A neuromorphic manycore processor with on-chip learning," *IEEE Micro*, 2018.
- [29] M. V. Debole, B. Taba, A. Amir, F. Akopyan, A. Andreopoulos, W. P. Risk, J. Kusnitz, C. O. Otero, T. K. Nayak, R. Appuswamy, P. J. Carlson, A. S. Cassidy, P. Datta, S. K. Esser, G. J. Garreau, K. L. Holland, S. Lekuch, M. Mastro, J. McKinstry, C. Di Nolfo, J. Sawada, B. Paulovicks, K. Schleupen, B. G. Shaw, J. L. Klamo, M. D. Flickner, J. V. Arthur, and D. S. Modha, "TrueNorth: Accelerating from zero to 64 million neurons in 10 years," *Computer*, 2019.
- [30] P. Diehl and M. Cook, "Unsupervised learning of digit recognition using spike-timing-dependent plasticity," *Front. in Comp. Neuroscience*, 2015.
- [31] J. Doevenspeck, R. Degraeve, A. Fantini, S. Cosemans, A. Mallik, P. Debacker, D. Verkest, R. Lauwereins, and W. Dehaene, "OxRRAM-Based Analog in-Memory Computing for Deep Neural Network Inference: A Conductance Variability Study," *TED*, 2021.
- [32] B. R. Fernando, Y. Qi, C. Yakopcic, and T. M. Taha, "3D memristor crossbar architecture for a multicore neuromorphic system," in *IJCNN*, 2020.
- [33] M. E. Fouda, A. M. Eltawil, and F. Kurdahi, "Modeling and analysis of passive switching crossbar arrays," *TCAS I*, 2017.
- [34] M. E. Fouda, J. Lee, A. M. Eltawil, and F. Kurdahi, "Overcoming crossbar nonidealities in binary neural networks through learning," in *NANOARCH*, 2018.
- [35] M. E. Fouda, E. Neftci, A. Eltawil, and F. Kurdahi, "Effect of asymmetric nonlinearity dynamics in rrams on spiking neural network performance," in *ACSSC*, 2019.
- [36] M. E. Fouda, S. Lee, J. Lee, G. H. Kim, F. Kurdahi, and A. Eltawil, "IR-QNN framework: An IR drop-aware offline training of quantized crossbar arrays," *IEEE Access*, 2020.
- [37] C. Frenkel, "Bottom-up and top-down neuromorphic processor design: Unveiling roads to embedded cognition," Ph.D. dissertation, PhD thesis, UCL-Université Catholique de Louvain, 2020.

- [38] F. Galluppi, S. Davies, A. Rast, T. Sharp, L. A. Plana, and S. Furber, "A hierarchical configuration system for a massively parallel neural hardware platform," in *Computing Frontiers*, 2012, pp. 183–192.
- [39] D. F. Goodman and R. Brette, "The brian simulator," *Frontiers in Neuroscience*, 2009.
- [40] R. Gopalakrishnan, Y. Chua, P. Sun, A. J. S. Kumar, and A. Basu, "HFNet: A CNN architecture co-designed for neuromorphic hardware with a crossbar array of synapses," *Frontiers in Neuroscience*, 2020.
- [41] A. Gulli and S. Pal, *Deep learning with Keras*, 2017.
- [42] Y. He, Y. Wang, X. Zhao, H. Li, and X. Li, "Towards state-aware computation in ReRAM neural networks," in *DAC*, 2020.
- [43] M. L. Hines and N. T. Carnevale, "The NEURON simulation environment," *Neural Computation*, 1997.
- [44] A. Hore and D. Ziou, "Image quality metrics: PSNR vs. SSIM," in *ICPR*, 2010.
- [45] M. Hu, H. Li, Y. Chen, Q. Wu, G. S. Rose, and R. W. Linderman, "Memristor crossbar-based neuromorphic computing system: A case study," *TNNLS*, 2014.
- [46] M. Hu, J. P. Strachan, Z. Li, E. M. Grafals, N. Davila, C. Graves, S. Lam, N. Ge, J. J. Yang, and R. S. Williams, "Dot-product engine for neuromorphic computing: Programming 1T1M crossbar to accelerate matrix-vector multiplication," in *DAC*, 2016.
- [47] P. K. Huynh, M. L. Varshika, A. Paul, M. Isik, A. Balaji, and A. Das, "Implementing spiking neural networks on neuromorphic architectures: A review," *arXiv*, 2022.
- [48] G. Indiveri, "A low-power adaptive integrate-and-fire neuron circuit," in *ISCAS*, 2003.
- [49] Y. Jeong, M. A. Zidan, and W. D. Lu, "Parasitic effect analysis in memristor-array-based neuromorphic systems," *TNANO*, 2017.
- [50] Y. Ji, Y. Zhang, S. Li, P. Chi, C. Jiang, P. Qu, Y. Xie, and W. Chen, "NEUTRAMS: Neural network transformation and co-design under neuromorphic hardware constraints," in *MICRO*, 2016.
- [51] Y. Kim, Y. Zhang, and P. Li, "A digital neuromorphic VLSI architecture with memristor crossbar synaptic array for machine learning," in *SOCC*, 2012.
- [52] Y. Kim, Y. Zhang, and P. Li, "A reconfigurable digital neuromorphic processor with memristive synaptic crossbar for cognitive computing," *JETC*, 2015.
- [53] O. Krestinskaya, A. Irmanova, and A. P. James, "Memristive non-idealities: Is there any practical implications for designing neural network chips?" in *ISCAS*, 2019.
- [54] S. Kundu, K. Basu, M. Sadi, T. Titirsha, S. Song, A. Das, and U. Guin, "Special session: Reliability analysis for ML/AI hardware," in *VTS*, 2021.
- [55] M. Le and S. N. Truong, "Memristor crossbar circuits for neuromorphic pattern recognition," in *ISOCC*, 2021.
- [56] D. Lee, Y. Kim, V. Seshadri, J. Liu, L. Subramanian, and O. Mutlu, "Tiered-latency DRAM: A low latency and low cost DRAM architecture," in *HPCA*, 2013.
- [57] T. Li, X. Bi, N. Jing, X. Liang, and L. Jiang, "Sneak-path based test and diagnosis for 1R RRAM crossbar using voltage bias technique," in *DAC*, 2017.
- [58] Y. Li and K.-W. Ang, "Hardware implementation of neuromorphic computing using large-scale memristor crossbar arrays," *Advanced Intelligent Systems*, 2021.
- [59] C.-Y. Liao, K.-Y. Hsiang, F.-C. Hsieh, S.-H. Chiang, S.-H. Chang, J.-H. Liu, C.-F. Lou, C.-Y. Lin, T.-C. Chen, C.-S. Chang *et al.*, "Multibit ferroelectric FET based on nonidentical double  $h\text{fzr}_2$  for high-density nonvolatile memory," *EDL*, 2021.
- [60] C.-K. Lin, A. Wild, G. N. Chinya, T.-H. Lin, M. Davies, and H. Wang, "Mapping spiking neural networks onto a manycore neuromorphic architecture," in *PLDI*, 2018.
- [61] C. Liu, B. Yan, C. Yang, L. Song, Z. Li, B. Liu, Y. Chen, H. Li, Q. Wu, and H. Jiang, "A spiking neuromorphic design with resistive crossbar," in *DAC*, 2015.
- [62] X. Liu, W. Wen, X. Qian, H. Li, and Y. Chen, "Neu-NoC: A high-efficient interconnection network for accelerated neuromorphic systems," in *ASP-DAC*, 2018.
- [63] W. Maass, "Networks of spiking neurons: The third generation of neural network models," *Neural Networks*, 1997.
- [64] A. Mallik, D. Garbin, A. Fantini, D. Rodopoulos, R. Degraeve, J. Stuijt, A. Das, S. Schaafsma, P. Debacker, G. Donadio *et al.*, "Design-technology co-optimization for OxRRAM-based synaptic processing unit," in *VLSIT*, 2017.
- [65] C. Mead, "Neuromorphic electronic systems," *Proc. of the IEEE*, 1990.
- [66] S. Moradi, N. Qiao, F. Stefanini, and G. Indiveri, "A scalable multicore architecture with heterogeneous memory structures for dynamic neuromorphic asynchronous processors (DYNAPs)," *TBCAS*, 2017.
- [67] O. Mutlu, "Memory scaling: A systems architecture perspective," in *IMW*, 2013.
- [68] O. Mutlu and L. Subramanian, "Research problems and opportunities in memory systems," *SUFI*, 2015.
- [69] N. S. Nukala, N. Kulkarni, and S. Vrudhula, "Spintronic threshold logic array (STLA)—A compact, low leakage, non-volatile gate array architecture," *JPDC*, 2014.

- [70] A. Paul, S. Song, and A. Das, "Design technology co-optimization for neuromorphic computing," in *IGSC Workshop*, 2021.
- [71] A. Paul, S. Song, T. Titirsha, and A. Das, "On the mitigation of read disturbances in neuromorphic inference hardware," *D&T*, 2022.
- [72] B. Rajendran, A. Sebastian, M. Schmuker, N. Srinivasa, and E. Eleftheriou, "Low-power neuromorphic hardware for signal processing applications: A review of architectural and system-level design approaches," *Signal Processing Magazine*, 2019.
- [73] M. Rakka, M. Fouda, R. Kanj, A. Eltawil, and F. Kurdahi, "Design exploration of sensing techniques in 2T-2R resistive ternary CAMs," *TCAS II: Express Briefs*, 2020.
- [74] W. Shim, Y. Luo, J.-s. Seo, and S. Yu, "Impact of read disturb on multilevel RRAM based inference engine: Experiments and model prediction," in *IRPS*, 2020.
- [75] S. Song and A. Das, "Design methodologies for reliable and energy-efficient PCM systems," in *IGSC Workshops*, 2020.
- [76] S. Song and A. Das, "A case for lifetime reliability-aware neuromorphic computing," in *MWSCAS*, 2020.
- [77] S. Song, A. Das, O. Mutlu, and N. Kandasamy, "Enabling and exploiting partition-level parallelism (PALP) in phase change memories," *TECS*, 2019.
- [78] S. Song, A. Balaji, A. Das, N. Kandasamy, and J. Shackleford, "Compiling spiking neural networks to neuromorphic hardware," in *LCTES*, 2020.
- [79] S. Song, A. Das, and N. Kandasamy, "Exploiting inter- and intra-memory asymmetries for data mapping in hybrid tiered-memories," in *ISMM*, 2020.
- [80] S. Song, A. Das, and N. Kandasamy, "Improving dependability of neuromorphic computing with non-volatile memory," in *EDCC*, 2020.
- [81] S. Song, A. Das, O. Mutlu, and N. Kandasamy, "Improving phase change memory performance with data content aware access," in *ISMM*, 2020.
- [82] S. Song, H. Chong, A. Balaji, A. Das, J. Shackleford, and N. Kandasamy, "DFSynthesizer: Dataflow-based synthesis of spiking neural networks to neuromorphic hardware," *TECS*, 2021.
- [83] S. Song, A. Das, O. Mutlu, and N. Kandasamy, "Aging-aware request scheduling for non-volatile main memory," in *ASP-DAC*, 2021.
- [84] S. Song, J. Hanamshet, A. Balaji, A. Das, J. Krichmar, N. Dutt, N. Kandasamy, and F. Catthoor, "Dynamic reliability management in neuromorphic computing," *JETC*, 2021.
- [85] S. Song, L. V. Mirtinti, A. Das, and N. Kandasamy, "A design flow for mapping spiking neural networks to many-core neuromorphic hardware," in *ICCAD*, 2021.
- [86] S. Song, T. Titirsha, and A. Das, "Improving inference lifetime of neuromorphic systems via intelligent synapse mapping," in *ASAP*, 2021.
- [87] S. A. Thomas, S. K. Vohra, R. Kumar, R. Sharma, and D. M. Das, "Analysis of parasitics on CMOS based memristor crossbar array for neuromorphic systems," in *MWSCAS*, 2021.
- [88] T. Titirsha and A. Das, "Reliability-performance trade-offs in neuromorphic computing," in *IGSC Workshops*, 2020.
- [89] T. Titirsha and A. Das, "Thermal-aware compilation of spiking neural networks to neuromorphic hardware," in *LCPC*, 2020.
- [90] T. Titirsha, S. Song, A. Balaji, and A. Das, "On the role of system software in energy management of neuromorphic computing," in *CF*, 2021.
- [91] T. Titirsha, S. Song, A. Das, J. Krichmar, N. Dutt, N. Kandasamy, and F. Catthoor, "Endurance-aware mapping of spiking neural networks to neuromorphic hardware," *TPDS*, 2021.
- [92] S. Tuli, M. Rios, A. Levisse, and D. A. ESL, "RRAM-VAC: A variability-aware controller for RRAM-based memory architectures," in *ASP-DAC*, 2020.
- [93] M. L. Varshika, A. Balaji, F. Corradi, A. Das, J. Stuijt, and F. Catthoor, "Design of many-core big little  $\mu$ Brains for energy-efficient embedded neuromorphic computing," in *DATE*, 2022.
- [94] Z. Wang, H. Zhang, T. Luo, W.-F. Wong, A. T. Do, P. Vishnu, W. Zhang, and R. S. M. Goh, "NCPower: Power modelling for NVM-based neuromorphic chip," in *ICONS*, 2020.
- [95] P. Wijesinghe, A. Ankit, A. Sengupta, and K. Roy, "An all-memristor deep spiking neural computing system: A step toward realizing the low-power stochastic brain," *TETCI*, 2018.
- [96] C. Xu, X. Dong, N. P. Jouppi, and Y. Xie, "Design implications of memristor-based RRAM cross-point structures," in *DATE*, 2011.
- [97] C.-X. Xue, W.-H. Chen, J.-S. Liu, J.-F. Li, W.-Y. Lin, W.-E. Lin, J.-H. Wang, W.-C. Wei, T.-W. Chang, T.-C. Chang *et al.*, "24.1 a 1Mb multibit ReRAM computing-in-memory macro with 14.6 ns parallel MAC computing time for CNN based AI edge processors," in *ISSCC*, 2019.
- [98] S. R. Young, P. Devineni, M. Parsa, J. T. Johnston, B. Kay, R. M. Patton, C. D. Schuman, D. C. Rose, and T. E. Potok, "Evolving energy efficient convolutional neural networks," in *Big Data*, 2019.

- [99] S. Yu, Y. Deng, B. Gao, P. Huang, B. Chen, X. Liu, J. Kang, H.-Y. Chen, Z. Jiang, and H.-S. P. Wong, “Design guidelines for 3D RRAM cross-point architecture,” in *ISCAS*, 2014.
- [100] X. Zhang, A. Huang, Q. Hu, Z. Xiao, and P. K. Chu, “Neuromorphic computing with memristor crossbar,” *Physica Status Solidi (a)*, 2018.
- [101] W. Zhao and Y. Cao, “Predictive technology model for nano-CMOS design exploration,” *JETC*, 2007.
- [102] Z. Zhu, J. Lin, M. Cheng, L. Xia, H. Sun, X. Chen, Y. Wang, and H. Yang, “Mixed size crossbar based RRAM CNN accelerator with overlapped mapping method,” in *ICCAD*, 2018.

## A SPIKING NEURAL NETWORKS

Spiking Neural Networks (SNNs) enable powerful computations due to their spatio-temporal information encoding capabilities [63]. An SNN consists of neurons, which are connected via synapses. A neuron can be implemented as an integrate-and-fire (IF) logic, which is illustrated in Figure 21 (left). Here, an input current  $U(t)$  (i.e., spike from a pre-synaptic neuron) raises the membrane voltage of the neuron. When this voltage crosses a threshold  $V_{th}$ , the IF logic emits an output spike, which propagates to its post-synaptic neuron. Figure 21 (middle) illustrates the membrane voltage of the IF neuron due to an input spike train. The moment of threshold crossing is illustrated in Figure 21 (right). These are the firing times of the output spike train of the neuron.

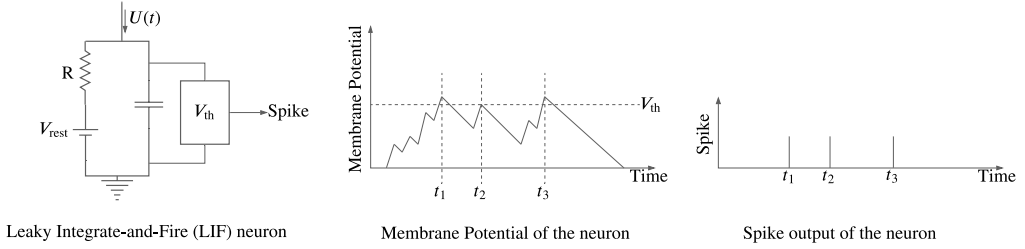


Fig. 21. A leaky integrate-and-fire (LIF) neuron with current input  $U(t)$  (left). The membrane potential over time of the neuron (middle). The spike output of the neuron representing its firing time (right).

SNNs can implement many machine learning approaches such as supervised learning, unsupervised learning, reinforcement learning, and lifelong learning. We focus on supervised machine learning, where an SNN is pre-trained with representative data. Machine learning **inference** refers to feeding live data points to this trained SNN to generate the corresponding output.

## B QUALITY OF INFERENCE

The **quality** of machine learning inference can be expressed in terms of accuracy [5], Mean Square Error (MSE) [24], Peak Signal-to-Noise Ratio (PSNR) [21], and Structural Similarity Index Measure (SSIM) [44]. While accuracy is commonly used for assessing the quality of supervised learning, e.g., using Convolution Neural Networks (CNNs), there are also applications such as edge detection, where the quality is assessed using other metrics such as SSIM. In our prior work [9], we have shown that these quality metrics are a function of the inter-spike interval (ISI) between neurons. Therefore, any deviation of ISI (called ISI distortion) from its trained value may lead to quality loss. To describe ISI, let  $\{t_1, t_2, \dots, t_K\}$  denote a neuron’s firing times in the time interval  $[0, T]$ , the average ISI of this spike train is

$$I = \sum_{i=2}^K (t_i - t_{i-1}) / (K - 1). \quad (7)$$



To illustrate how a change in ISI, called **ISI distortion**, impacts inference quality, we use a small SNN in which three input neurons are connected to an output neuron. Figure 22 illustrates the impact of ISI distortion on the output spike. In the top sub-figure, a spike is generated at the output neuron at  $22\mu\text{s}$  due to spikes from the input neurons. In the bottom sub-figure, the second spike from input 3 is delayed, i.e., it has an ISI distortion. Due to this distortion, there is no output spike generated. Missing spikes can impact inference quality, as spikes encode information in SNNs.

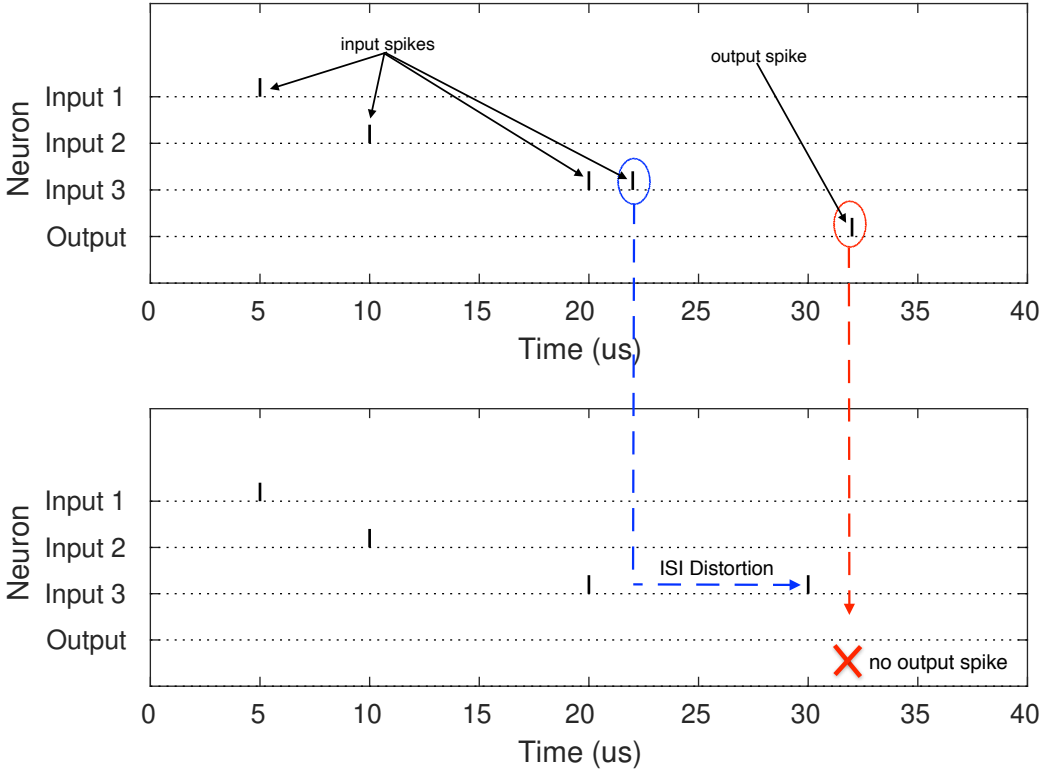


Fig. 22. Impact of ISI distortion on accuracy [8]. Top sub-figure shows a scenario where an output spike is generated based on the spikes received from the three input neurons. Bottom sub-figure shows a scenario where the second spike from neuron 3 is delayed. There are no output spikes generated.

Figure 23 shows the impact of ISI distortion on the quality of image smoothing implemented using an SNN [21]. Figure 23a shows the input image, which is fed to the SNN. Figure 23b shows the output of the image smoothing application with no ISI distortion. PSNR of the output with reference to the input is 20. Figure 23c shows the output with ISI distortion due to variation in latency within neuromorphic PEs of the hardware. PSNR of this output with respect to the input is 19. A reduction in PSNR indicates that the output image quality with ISI distortion is lower than the one without distortion. In fact, image quality deteriorates with increase in ISI distortion. We use ISI distortion as a measure of the quality of machine learning inference [9]. Our aim is to improve this inference quality via technological and architectural enhancements that reduce ISI distortion when the inference task is implemented on neuromorphic PEs of a hardware.

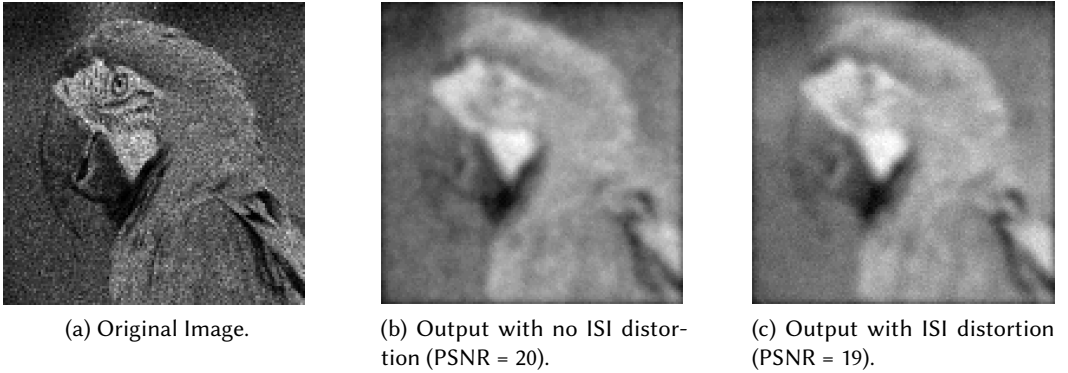


Fig. 23. Impact of ISI distortion on image smoothing.

### C HARDWARE IMPLEMENTATION OF MACHINE LEARNING INFERENCE

Most neuromorphic hardware platforms are implemented as tiled-based architectures [16, 28, 29, 37, 72, 93], where the tiles are interconnected via a shared interconnect such as Network-on-Chip [62] and Segmented Bus [7]. Figure 24 illustrates a tile-based neuromorphic hardware platform, where the tiles can communicate concurrently. Each tile includes 1) a neuromorphic PE, which consists of neuron and synapse circuitries and 2) a network interface, which encodes spikes into Address Event Representation (AER) and communicates these AER packets to the switch for routing to their destination tiles. A common design practice is to use analog crossbars to implement a neuromorphic PE [2, 9, 45, 52, 55, 58, 61, 100]. Within a crossbar, a pre-synaptic neuron circuit acts as a current driver and is placed on a wordline, while a post-synaptic neuron circuit acts as a current sink and is placed on a bitline as illustrated in Figure 1 (left).

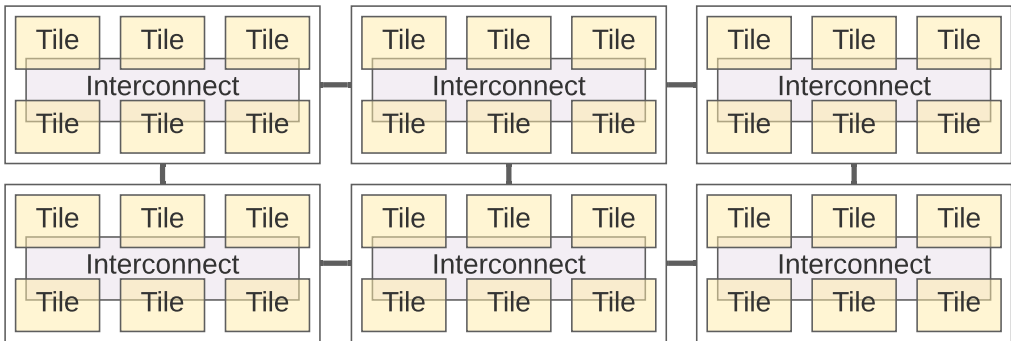


Fig. 24. Tile-based neuromorphic hardware, representative of hardware platforms such as TrueNorth [29], Loihi [28], DYNAPs [66], and  $\mu$ Brain [93].

Since a crossbar can accommodate only a limited number of neurons and synapses, a machine learning model is first partitioned into clusters, where each cluster can be implemented on a crossbar of the hardware. Partitioned clusters are then mapped to different crossbars when admitting the model to the hardware platform. To this end, several heuristic approaches are proposed in literature.

PSOPART [27] minimizes spike latency on the shared interconnect, SpiNeMap [9] minimizes interconnect energy, DFSynthesizer [82] maximizes throughput, DecomposedSNN [11] maximizes crossbar utilization, EaNC [90] minimizes overall energy of a machine learning task by targeting both computation and communication energy, TaNC [89] minimizes the average temperature of each crossbar, eSpine [91] maximizes NVM endurance in a crossbar, RENEU [80] minimizes the circuit aging in a crossbar's peripheral circuits, and NCil [86] reduces read disturb issues in a crossbar, improving the inference lifetime. Beside these techniques, there are also other software frameworks [1, 3, 4, 6, 12, 23, 25, 38, 47, 50, 54, 60, 71, 75, 76, 78, 85, 88] and run-time approaches [10, 84], addressing one or more of these optimization objectives.

We investigate the internal architecture of a crossbar and find that the parasitic components introduce delay in propagating current from a pre-synaptic neuron to a post-synaptic neuron as illustrated in Figure 1 (right). This delay depends on the specific current path used in the mapping.

Higher the number of parasitic components on a current path, larger is its propagation delay. Parasitic components on bitlines and wordlines are a major source of latency at scaled process technology nodes and they create significant **latency variation** in a crossbar. Specifically, the latency of a synaptic connection in an SNN depends precisely on the memory cell in the crossbar that is used to implement it. Such latency variation can introduce ISI distortion (Section B), which may impact the quality of an inference task.

#### D NON-VOLATILE MEMORY TECHNOLOGY

RRAM technology presents an attractive option for implementing memory cells of a crossbar due to its demonstrated potential for low-power multilevel operation and high integration density [64]. An RRAM cell is composed of an insulating film sandwiched between conducting electrodes forming a metal-insulator-metal (MIM) structure (see Figure 25). Recently, conducting filament-based metal-oxide RRAM implemented with transition-metal-oxides such as  $\text{HfO}_2$ ,  $\text{ZrO}_2$ , and  $\text{TiO}_2$  has received considerable attention due to their low-power and CMOS-compatible scaling.

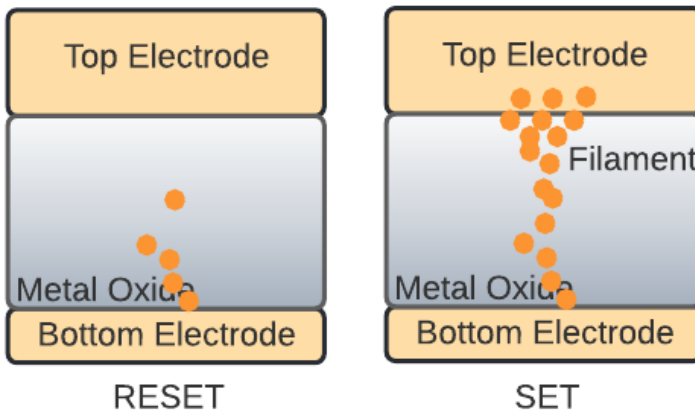


Fig. 25. Operation of an RRAM cell with the  $\text{HfO}_2$  layer sandwiched between the metals Ti (top electrode) and TiN (bottom electrode). The right subfigure shows the formation of LRS/SET state. The left subfigure shows the HRS/RESET state.

Synaptic weights are represented as conductance of the insulating layer within each RRAM cell. To program an RRAM cell, elevated voltages are applied at the top and bottom electrodes, which

re-arranges the atomic structure of the insulating layer. Figure 25 shows the High-Resistance State (HRS) and the Low-Resistance State (LRS) of an RRAM cell. An RRAM cell can also be programmed into intermediate low-resistance states, allowing its multilevel operations [17].