# Impact of On-Chip Interconnect on In-Memory Acceleration of Deep Neural Networks

GOKUL KRISHNAN*, Arizona State University, USA
SUMIT K. MANDAL*, University of Wisconsin-Madison, USA
CHAITALI CHAKRABARTI, Arizona State University, USA
JAE-SUN SEO, Arizona State University, USA
UMIT Y. OGRAS, University of Wisconsin-Madison, USA
YU CAO, Arizona State University, USA

With the widespread use of Deep Neural Networks (DNNs), machine learning algorithms have evolved in two diverse directions – one with ever-increasing connection density for better accuracy and the other with more compact sizing for energy efficiency. The increase in connection density increases on-chip data movement, which makes efficient on-chip communication a critical function of the DNN accelerator. The contribution of this work is threefold. First, we illustrate that the point-to-point (P2P)-based interconnect is incapable of handling a high volume of on-chip data movement for DNNs. Second, we evaluate P2P and network-on-chip (NoC) interconnect (with regular topology) for SRAM- and ReRAM-based in-memory computing (IMC) architectures for a range of DNNs. This analysis shows the necessity for the optimal interconnect choice for an IMC DNN accelerator. Finally, we perform an experimental evaluation for different DNNs to empirically obtain the performance of the IMC architecture with both NoC-tree and NoC-mesh. We conclude that, at the tile-level, NoC-tree is appropriate for compact DNNs employed at the edge, and NoC-mesh is necessary to accelerate DNNs with high connection density. Furthermore, we propose a technique to determine the optimal choice of interconnect for any given DNN. In this technique, we use analytical models of NoC to evaluate end-to-end communication latency of any given DNN. We demonstrate that the interconnect optimization in the IMC architecture results in up to 6× improvement in energy-delay-area product for VGG-19 inference compared to the state-of-the-art ReRAM-based IMC architectures.

CCS Concepts: • **Hardware** → **Emerging architectures**; **Interconnect**; **Network on chip**; • **Computing methodologies** → **Machine learning**; **Artificial intelligence**.

---

*Authors have equal contributions.
Authors' addresses: Gokul Krishnan*, gkrish19@asu.edu, Arizona State University, School of Electrical, Computer, and Energy Engineering, Tempe, AZ, 85287, USA; Sumit K. Mandal*, skmandal@wisc.edu, University of Wisconsin-Madison, Department of Electrical and Computer Engineering, Madison, WI, 53706, USA; Chaitali Chakrabarti, chaitali@asu.edu, Arizona State University, School of Electrical, Computer, and Energy Engineering, Tempe, AZ, 85287, USA; Jae-sun Seo, jseo28@asu.edu, Arizona State University, School of Electrical, Computer, and Energy Engineering, Tempe, AZ, 85287, USA; Umit Y. Ogras, uogras@wisc.edu, University of Wisconsin-Madison, Department of Electrical and Computer Engineering, Madison, WI, 53706, USA; Yu Cao, Yu.Cao@asu.edu, Arizona State University, School of Electrical, Computer, and Energy Engineering, Tempe, AZ, 85287, USA.
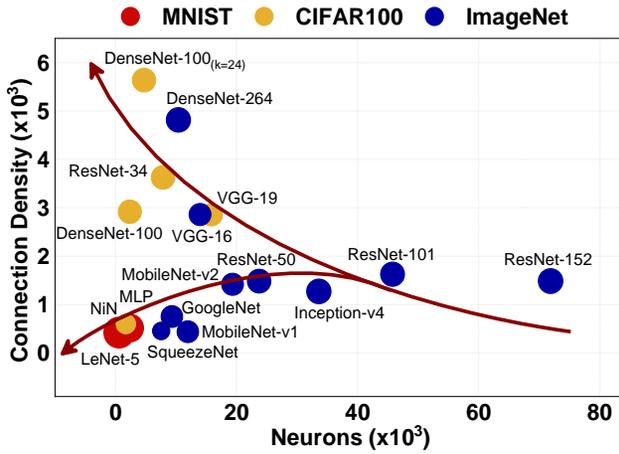
---

Fig. 1. Connection density of different DNNs for three different datasets. Each output feature map (convolution layer) and neural unit (FC layer) represent a neuron. Larger markers represent higher accuracy.

## 1 INTRODUCTION

DNNs have achieved high accuracy that exceeds human-level perception for a variety of applications such as computer vision, natural language processing, and medical imaging [4, 15, 19]. The DNNs that achieve higher accuracy tend to consist of deeper and denser network structures. On the other hand, DNNs for edge devices tend to use smaller and shallower networks.

Figure 1 shows the trend in connection density for various DNNs in the literature, where *connection density* is defined as the average number of connections per neuron in DNNs. In the context of DNNs, a neuron is defined as an output feature of a convolution layer and every neural unit of the fully-connected (FC) layer. Three representative DNN structures and connection patterns are illustrated in Figure 2. Linear structures such as LeNet-5 [17] and VGG-19 [31] have a connection density of one owing to one connection per neuron. Since residual networks such as ResNet [6] have residual skips, it has more connections than the number of neurons resulting in a connection density higher than one. Dense structures like DenseNet [8] have multiple connections from each neuron, resulting in a higher connection density.

We observe two main trends by analyzing the connection density for different DNNs in Figure 1. First, increasing connection density provides higher accuracy, which is essential for cloud-based computing platforms. Second, lower connection density is observed for compact models, which is necessary for edge computing hardware. Both hardware platforms require the processing of large amounts of data with corresponding power and performance constraints. Hence, there is a need to design optimal hardware architectures with low power and high performance for DNNs with different connection densities.

With limited on-chip memory, conventional DNN architectures inevitably involve a significant amount of communication with off-chip memory resulting in increased energy consumption [3]. However, it has been reported that the energy consumption of off-chip communication is 1,000× higher than the energy required to perform the computations [7]. Dense structures like DenseNet perform approximately $2.7 \times 10^7$ off-chip memory accesses to process a frame of an image [8]. As a result, off-chip memory access becomes the energy bottleneck for hardware architectures of dense structures. Employing dense embedded non-volatile memory (NVM) such as ReRAM for in-memory computing (IMC) substantially reduces off-chip memory accesses [30, 32].
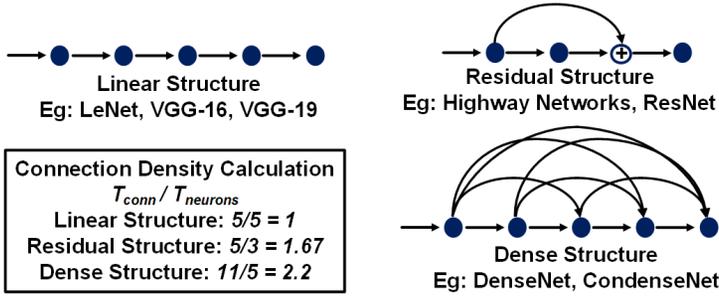
Fig. 2. Different types of DNN structures and their representative connection density.

On-chip interconnect is an integral part of hardware architectures that incorporate in-memory acceleration. Both point-to-point (P2P) interconnect [16, 33] and NoC-based interconnect [3, 14, 30] are used for on-chip communication in state-of-the-art DNN accelerators. Shafiee *et al.* [30] utilizes a concentrated mesh for the interconnect, while Chen *et al.* [3] employs three different NoCs that are used for on-chip data movement in the architecture. In contrast, Krishnan *et al.* [14] utilizes a custom mesh-NoC for on-chip communication. The custom NoC derives the structure based on the on-chip traffic between different IMC processing elements (PEs), where each PE denotes the SRAM- or ReRAM-based IMC crossbar. A technique to construct custom NoC which provides minimum communication latency for a given DNN is proposed in [22]. Since custom NoC requires alteration in hardware for different DNNs, our studies focus on regular NoC topologies. A more detailed survey on work which design efficient interconnect for DNN accelerators can be found in [25].

To better understand the need for an NoC-based on-chip interconnect, we analyze the scalability of P2P interconnect in in-memory computing (IMC) architectures by evaluating the contribution of routing latency to end-to-end latency for different DNNs, as shown in Figure 3. The contribution of routing latency increases up to 94% with increasing connection density. The high routing latency is attributed to the increased connection density, which correlates to more on-chip data movement. VGG-19 shows a reduced contribution compared to lower connection density DNNs due to the high utilization of the IMC PEs or crossbars resulting in reduced on-chip data movement. Hence, P2P networks do not provide a scalable solution for high connection density DNNs. At the same time, NoC-based interconnects require higher area and energy for operation and can result in a significant overhead for low connection density DNNs. Furthermore, different NoC topologies,
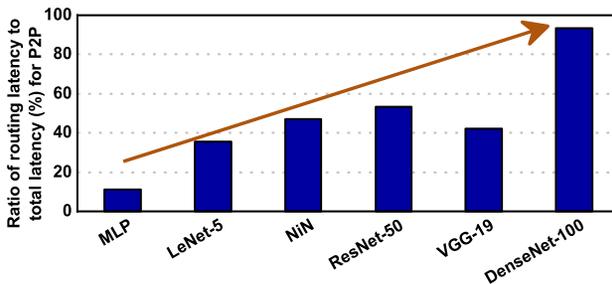


Fig. 3. Contribution of routing latency to total latency for different DNNs for a P2P-based IMC architecture [2]. With increase in connection density, routing latency becomes the bottleneck for performance.

Table 1. Summary of notations

| Symbol | Definition | Symbol | Definition |
|--------|-----------|--------|-----------|
| $N_L$ | Number of Layers | $x_i, y_i$ | Input image size in $i^{\text{th}}$ layer |
| $T_i$ | Number of tiles in $i^{\text{th}}$ layer | $C_i$ | Input channels of $i^{\text{th}}$ layer |
| $A_i$ | Number of activations in $i^{\text{th}}$ layer | $\lambda_{i,j,k}$ | Injection rate from $j^{\text{th}}$ tile of $(i-1)^{\text{th}}$ layer to $k^{\text{th}}$ tile of $i^{\text{th}}$ layer |
| $d_i$ | Input activations in $i^{\text{th}}$ layer | $L_{comm}$ | Total communication latency |

mesh, or tree, are appropriate for DNNs with varying connection densities. Therefore, a connection density-aware interconnect solution is critical to DNN acceleration.

In this work, we first perform an in-depth performance analysis of P2P interconnect-based in-memory computing (IMC) architectures [32]. Through this analysis, we establish that P2P-based interconnects are incapable of handling data communication for dense DNNs and that NoC-based interconnect is needed for IMC architectures. Next, we evaluate P2P-based and NoC-based SRAM and ReRAM IMC architectures for a range of DNNs. Further, we evaluate NoC-tree, NoC-mesh, and c-mesh topologies for the IMC architectures. A c-mesh NoC is used in [30] at the tile-level to connect different tiles. C-mesh uses more number of links and routers, providing better performance in terms of communication latency. However, interconnect area and energy becomes exorbitantly high for c-mesh NoC. Therefore, the energy-delay-area product (EDAP) of c-mesh is higher than NoC-mesh. Hence, we restrict the detailed evaluations to NoC-mesh and NoC-tree. In these evaluations, we perform cycle-accurate NoC simulations through Booksim [10]. However, cycle-accurate NoC simulations are very time consuming and consequently slow down the overall performance analysis of IMC architectures. Our experiment with different DNNs (the simulation framework is described in more detail in Section 3) shows that cycle-accurate NoC simulation takes up to 80% of the total simulation time for high connection density DNNs.

To accelerate the overall performance analysis of the IMC architecture, we propose analytical models to estimate the NoC performance of a given DNN. Specifically, we incorporate the analytical router modeling technique presented in [26] to obtain the performance model for an NoC router. Then we extend the existing analytical model to get an estimation of end-to-end communication latency for NoC-tree and NoC-mesh for any given DNN as a function of the number of neurons and connection density. Through the analytical latency model, the variable communication patterns of different DNNs are incorporated using connection density and number of neurons. Leveraging this analysis and the analytical model, we conclude the importance of the optimal choice of interconnect at different hierarchies of the IMC architecture. Utilizing the same analysis, we provide guidance for the optimal choice of interconnect for IMC architectures. At the tile-level, NoC-mesh for high connection density DNNs and an NoC-tree for low connection density DNNs provide low power and high performance for IMC-based architectures. Leveraging this observation, we propose an NoC-based heterogeneous interconnect IMC architecture for DNN acceleration. We demonstrate that the NoC-based heterogeneous interconnect IMC architecture (ReRAM) achieves up to 6× improvement in the energy-delay-area product (EDAP) for inference of VGG-19 when compared to state-of-the-art implementations. The following are key contributions of this work:
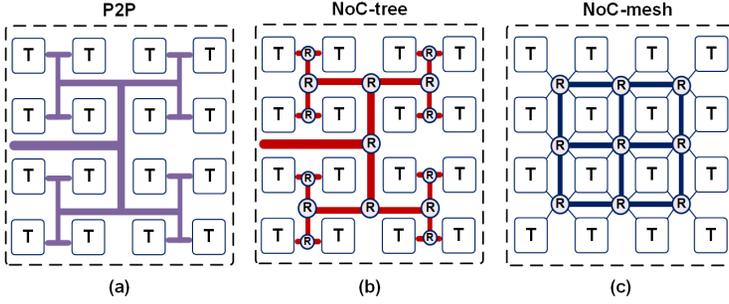
Fig. 4. Multi-tiled IMC architecture with routing architectures based on (a) P2P network, (b) NoC-tree, (c) NoC-mesh. NoC-tree is a P2P network with routers at junctions.

- An in-depth analysis of the shortcomings of P2P-based interconnect and the need for NoC in IMC architectures.
- Analytical and empirical analysis to guide the choice of optimal NoC topology for an NoC-based heterogeneous interconnect.
- The proposed heterogeneous interconnect IMC architecture achieves 6× improvement in EDAP with respect to state-of-the-art ReRAM-based IMC accelerators.

The rest of the paper is organized as follows. Section 2 introduces the background and motivation, Section 3 discusses the simulation framework used in this work, and Section 4 presents the analytical performance modeling-based technique to obtain the optimal choice of NoC for any given DNN. Section 5 presents the in-memory architecture with heterogeneous interconnect. Section 6 discusses the experimental results, and Section 7 concludes the paper.

## 2 MOTIVATION AND RELATED WORK

### 2.1 Deep Neural Networks

We categorize DNNs into three main classes, linear [31], residual [6], and dense [34], as shown in Figure 2. DNN structures include convolution layers stacked on top of each other for feature extraction and a set of classifier layers at the end to classify based on the features. A data point $d(x, y, c)$ in layer $i+1$ can be expressed using the notation summarized in Table 1 as follows:

$$d_{i+1}(x, y, c) = \sum_{c_i=0}^{C_i-1} \sum_{kx_i=0}^{Kx_i-1} \sum_{ky_i=0}^{Ky_i-1} K_i[kx_i, ky_i, c_i, c] \times d_i[(x + kx_i), (y + ky_i), c_i], \qquad (1)$$

where $K_i$ is the kernel, $Kx_i$ is the number of rows in the kernel, and $Ky_i$ is the number of columns in the kernel of the convolution layer i. To implement (1) on hardware, $x_{i+1} \times y_{i+1} \times C_{i+1} \times x_i \times y_i \times C_i$ number of multiplications and $x_{i+1} \times y_{i+1} \times C_{i+1} \times C_i \times Kx_i \times (Ky_i - 1)$ number of additions need to be performed. In addition to convolution and FC layers, pooling and non-linear activation layers such as rectified linear unit (ReLU) are present in the DNN algorithms.

### 2.2 In-Memory Computing with Crossbars

DNNs with a large number of weights requires a considerable amount of computations. Conventional architectures separate access of data from the memory and computation in the computing unit. This results in increased computation and data movement, reducing both the throughput and energy-efficiency for DNN inference. In contrast, in-memory computing (IMC) seamlessly integrates computation and memory access in a single unit such as the crossbar [14, 30, 32]. Through this,
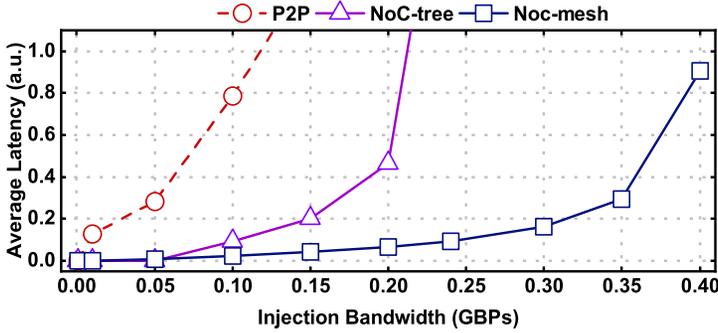
Fig. 5. Comparison of average latency among P2P, NoC-tree, and NoC-mesh interconnect for different injection bandwidth [10]. NoC topologies show better scalability than P2P interconnect.

IMC achieves higher energy efficiency and throughput as compared to conventional von-Neumann architectures.

The IMC technique localizes computation and data memory in a more compact design and enhances parallelism with multiple-row access, resulting in improved performance [11, 12]. The data accumulation is achieved through either current or charge accumulation. The size of the IMC subarray usually varies from 64×64 to 512×512. Along with the computing unit, peripheral circuits such as sample and hold circuit, analog-to-digital converter (ADC), and shift-and-add circuits are used to obtain each DNN layer's result. In this work, we focus on IMC designs based on both SRAM [11, 12, 35] and ReRAM [14, 23, 29, 32] crossbars.

## 2.3 Interconnect Network

As discussed in Section 1, the on-chip interconnect is critical to the accelerator performance for DNN acceleration. There are multiple topologies for Network-on-Chip (NoC). The well-known topologies are mesh, tree, torus, hypercube, and concentrated mesh (c-mesh). NoC with torus topology shows better performance than mesh due to long links between the nodes located at the edges. However, the power consumption by torus is significantly higher than mesh, as shown in [24]. Hypercube and c-mesh have a similar disadvantage as a torus. Therefore, only NoC-tree and NoC-mesh are considered in this work. Also, they are the industrial standard for SoCs used in heavy workloads [9].

Figure 4 illustrates representative interconnect schemes of P2P, NoC-tree, and NoC-mesh for multi-tiled IMC architectures. Each tile consists of several crossbar sub-arrays which perform the IMC operation. Existing implementations of DNN accelerators use both P2P-based [16, 33] and NoC-based [14, 30, 36] interconnect for on-chip communication. To better understand the performance of different interconnect architectures, we plot the average interconnect latency for a P2P network with 64 nodes, NoC-tree with 64 nodes, and an 8×8 NoC-mesh with X–Y routing as shown in Figure 5. The NoC utilizes one virtual channel, a buffer size (all input and output buffers) of eight, and three router pipeline stages. We observe that for lower injection rates, the performance is comparable for all topologies, while for higher injection rates, NoC performs better in terms of latency. Hence, NoC provides better scalability and performance compared to P2P interconnects. Moreover, with increasing connection density, injection bandwidth between layers increase due to increased on-chip data movement. Therefore, P2P interconnect performs poorly for DNNs with high connection density. Hence, there is a need for systematic guidance for choosing the optimal
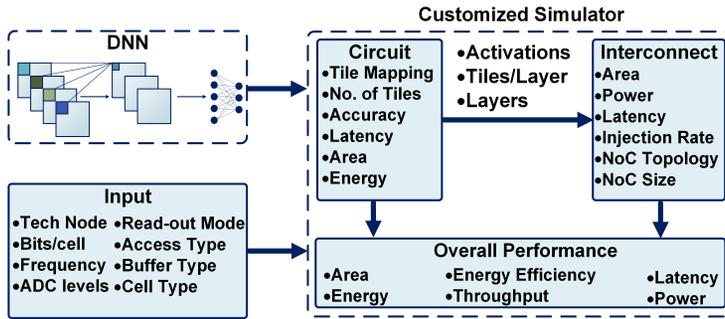
Fig. 6. Block-level representation of the proposed architecture simulator.

interconnect for in-memory acceleration of DNNs. Other works such as [3] utilizes three separate NoC for weights, activations and partial sums. Such a design choice results in increased area and energy cost for the interconnect fabric. Furthermore, the three NoCs are under-utilized, resulting in a sub-optimal design choice for acceleration of DNNs.

## 2.4 Analytical Modeling of NoCs

Till date, multiple NoC performance analysis techniques have been proposed for SoCs [13, 20, 21, 26, 28]. The analytical performance model for an NoC router assumes that the probability distribution of input traffic is in a continuous time domain [26]. However, all transactions in an IMC architecture happen in a discrete clock cycle. An analytical performance modeling technique for NoCs in the discrete-time domain is proposed in [21]. In this work, we estimate end-to-end communication latency for different DNNs as a function of connection density and the number of neurons of the DNN. Specifically, we utilize the analytical model for NoC router presented in [26] with the modifications for discrete time input [21] and extend the model to obtain end-to-end communication latency for NoC-tree and NoC-mesh.

## 3 SIMULATION FRAMEWORK

There exist multiple simulators that evaluate the performance of DNNs on different hardware platforms [2, 5]. These simulators consider different technologies, platforms, and peripheral circuit modeling while providing less consideration to interconnect. With the advent of dense DNN structures [34], the importance of interconnect cost is higher, as discussed in Section 1. In this work, we develop an in-house simulator, where a circuit-level performance estimator of the computing fabric is combined with a cycle-accurate simulator for the interconnect. The simulator also aims at being versatile by supporting multiple DNN algorithms across different datasets, and various interconnect schemes.

Figure 6 shows a block-level representation of the simulator. The inputs of the simulator primarily include the DNN structure, technology node, and frequency of operation. In the proposed simulation framework, any circuit-level performance estimator [2, 5] and any interconnect simulator [1, 10] can be plugged in to extract performance metrics such as area, energy, and latency, proving a common platform for system-level evaluation. In this work, we use customized versions of NeuroSim [2] for circuit simulation and BookSim [10] for cycle-accurate NoC simulation.

## 3.1 Circuit-level Simulator: Customized NeuroSim

The inputs to NeuroSim include the DNN structure indicating the layer size and layer count along with technology node, the number of bits per in-memory compute cell, frequency of operation, read-out mode, etc. The simulator performs the mapping of the entire DNN to a multi-tiled cross-bar

architecture by estimating the number of cross-bar arrays and the number of tiles per layer. Based on the size of the cross-bar $PE_x$ and $PE_y$, the number of cross-bar arrays is determined by (2).

$$\text{No. crossbars} = \sum_{i=1}^{N_L} \left\lceil \frac{(Kx_i \times Ky_i \times C_i)}{(PE_x)_i} \right\rceil \times \left\lceil \frac{(C_{i+1}) \times N_{bits}}{(PE_y)_i} \right\rceil, \tag{2}$$

where $N_{bits}$ is the precision of the weights. The total number of tiles is calculated as the ratio of the total number of crossbar arrays to the number of crossbar arrays per tile. Furthermore, the peripheral circuits are laid out, and the complete tile architecture is determined. The peripheral circuits include an ADC, sample and hold circuit, shift and add circuit, and a multiplexer circuit. However, NeuroSim lacks an accurate estimation of the interconnect cost in latency, energy, and area. Therefore, we replace the interconnect part of NeuroSim with customized BookSim. We also extract the performance metrics for tile-to-tile interconnect in NeuroSim and replace it with the BookSim tile-to-tile interconnect. With this customization, our circuit simulator only reports performance metrics, such as area, energy, and latency of the computing logic. It provides the number of tiles per layer, activations, and the number of layers to the interconnect simulator.

### 3.2 Interconnect Simulator: Customized BookSim

DNNs have varying structures resulting in different traffic loads and data-patterns between the IMC PEs. To accurately capture the NoC traffic of a given DNN configuration, we customize BookSim to evaluate the area, energy, and latency for interconnect, as shown in Figure 6. In the customized version of the BookSim, we enable simulation with non-uniform injection rate. We compute the injection rates for each source-destination pair in the multi-tiled architecture. The placement of tiles and routers in the IMC architecture has a direct impact on the interconnect performance. In this work, we incorporate the impact of mapping into the injection matrix calculation. The mapping of the DNN is performed such that each tile can have at least one layer while no layer is divided between two tiles. Figure 7 shows a sixteen tile IMC architecture with the tiles numbered. The red arrows show the data flow in the IMC architecture. Next, while evaluating the interconnect

---

**Algorithm 1:** Evaluation of interconnect latency through simulation

---

1  **Input:** Number of layers ($N_L$), Number of tiles in each layer ($T_i$), FPS ($F$), Number of activation in each layer ($A_i$), interconnected topology
2  **Output:** End-to-end interconnect latency ($L_{routing}$)
3  **for** *each layer i* **do**
4      **for** *each tile j in layer i − 1* **do**
5          **for** *each tile k in layer i* **do**
6              **if** *i > 0* **then**
7                  Compute $\lambda_{i,j,k}$ following Equation 3.
8              **end**
9          **end**
10     **end**
11     Simulate with interconnect topology and $\lambda_{i,j,k}$
12     Obtain $(l_i)_{sim}$ from the simulator.
13     Calculate $l_i$ following Equation 4.
14 **end**
15 Calculate $L_{comm}^{sim}$ : $L_{comm}^{sim} = \sum_{i=1}^{N_L} l_i$.
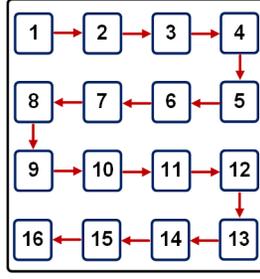
---

Fig. 7. Tile numbering and placement while mapping the DNN to the IMC architecture. The red arrows show the flow of the data across the tiles.

latency, we create an injection matrix that incorporates the position of the tile into the calculation by calculating the number of hops for each source-destination pair. Hence, the injection matrix incorporates the tile placement into the NoC latency calculation. Overall, the proposed approach can be generalized to any tile placement. Algorithm 1 describes the steps performed to compute injection rates and obtain the interconnect latency. Without loss of generality, we assume that the number of nodes required in the interconnect is equal to the total number of tiles across all layers.

The injection rate calculation is shown in lines 5–11 of Algorithm 1. The injection rate is expressed in (3) from each source to each destination in each layer.

$$\lambda_{i,j,k} = \frac{A_i \times N_{bits} \times FPS}{T_i \times T_{i-1} \times W \times freq} \tag{3}$$

where $N_{bits}$, $W$, and $FPS$ represent data precision and bus width, and frames-per-second throughput, respectively. In the numerator of (3), we multiply the number of input activations ($A_i$) for $i^{\text{th}}$ layer by $N_{bits}$ to obtain the total number of bits to be transferred from $(i-1)^{\text{th}}$ layer to $i^{\text{th}}$ layer for one frame of an image. We further multiply this term with FPS to obtain the total number of bits transferred between layers per second. Then, we divide this term by the operating frequency ($freq$) to obtain the total number of bits transferred between layers per cycle. We assume an equal injection rate between all tiles in two consecutive layers. Therefore, to get the number of bits transferred from one tile to another in two consecutive layers, the denominator in (3) includes a multiplication between $T_i$ and $T_{i-1}$. Thus, we divide the expression obtained so far by $W$ to obtain the injection rate ($\lambda_{i,j,k}$). The injection rate from every source to every destination is the input to the interconnect simulator. The interconnect simulator then provides average latency to complete all transactions from $(i-1)^{\text{th}}$ layer to $i^{\text{th}}$ layer ($(l_i)_{sim}$ cycles). Next, we multiply this latency with the number of bits from one tile to the next tile to get the total number of cycles required to transfer all data between two consecutive layers. Then, the latency from one layer to the next layer ($l_i$) is given by:

$$l_i = \frac{(l_i)_{sim} \times A_i \times N_{bits} \times FPS}{freq} \tag{4}$$

Finally, we accumulate the latency of all layers to compute the end-to-end interconnect latency as

$$L_{comm}^{sim} = \sum_{i=1}^{N_L} l_i \tag{5}$$

## 4 ANALYTICAL PERFORMANCE MODELS FOR NOCS IN IMC ARCHITECTURE

In this section, we discuss an analytical approach to estimate NoC performance for IMC architecture. The analytical performance model of NoCs is primarily useful to overcome longer simulation time

---

**Algorithm 2:** End-to-end latency computation through analytical models

---

1 **Input:** Input activation, Number of routers in each layer $l$ ($R_l$), Number of layers ($N_L$)

2 **Output:** End-to-end communication latency ($L_{comm}$)

3 **for** $l$ = 1: $N_L$-1 **do**

4     **for** $r$ = 1: $R_l$ **do**

        `/* Computing injection rate matrix */`

5         Compute $A_p^r$

6         Compute $\lambda_p^r$ using (6)

7         Construct $\Lambda^r$

        `/* Computing contention matrix */`

8         Compute forwarding probability matrix ($F^r$)

9         Compute contention matrix ($C^r$)

        `/* Computing average waiting time */`

10         Compute average queue length ($N^r$) using (8)

11         Compute average waiting time ($W_{avg}^r$) using (9)

12     **end**

13     Compute average latency for the layer ($L_{avg}^l$) using (10)

14 **end**

15 $L_{comm}^{ana} = \sum_{l=1}^{N_L} L_{avg}^l$

---

incurred by cycle-accurate NoC simulators. Specifically, we utilize analytical performance models for NoCs to compare the performance of NoC-tree and NoC-mesh for a given DNN. The analytical model of an NoC router is adopted from the work proposed in [26]. We extend this router model for NoC-tree and NoC-mesh to obtain end-to-end communication latency for different DNNs. Algorithm 2 describes the technique to evaluate the communication latency through analytical models. There are two major steps involved in analyzing the performance of an NoC: 1) Computing injection rate and 2) Computing contention probability matrix.

**Computing injection rate matrix ($\Lambda$):** First, the injection rate from each source to each destination ($\lambda_{sd}$) for each layer of the DNN is computed through (3). We note that the injection rate calculation incorporates the tile placement as detailed in Section 3.2. Each NoC router has five ports: North ($N$), South ($S$), East ($E$), West ($W$), and Self ($Se$). The injection rate at each port $p$ of every router $r$ ($\lambda_p^r, p \in \{N, S, E, W, Se\}$) is computed as:

$$\lambda_p^r = \frac{A_p^r \times N_{bits} \times FPS}{T_l \times T_{l+1} \times W \times freq} \tag{6}$$

where $T_l$ denotes the number of tiles in the $l^{\text{th}}$ layer. $\lambda_p^r$ is a function of the number of activations through each port $p$ of router $r$ ($A_p^r$). From $\lambda_p^r$, the injection rate matrix for router $r$ ($\Lambda^r$) is computed (as shown in line 5–7 of Algorithm 2), where $\Lambda^r = \{\lambda_{ij}^r\}, 1 \leq i \leq 5, 1 \leq j \leq 5, \lambda_{ij}^r = 0 \; \forall i \neq j$.

**Computing contention matrix ($C$):** Each element of the contention matrix $C$ ($c_{ij}$) denotes the contention between port $i$ and port $j$. To compute the contention matrix of router $r$ ($C^r = \{c_{ij}^r\}$), we first compute forwarding probability matrix $F^r = \{f_{ij}^r\}$. $f_{ij}^r$ denotes the probability of a packet that arrived at the port $i$ of the router $r$ to be forwarded to the port $j$, and is computed as shown
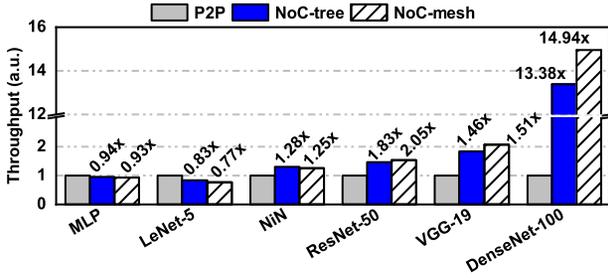
Fig. 8. Throughput comparison for three interconnect topologies (P2P, NoC-tree, and NoC-mesh) for SRAM-based IMC architecture, normalized to P2P, for different DNNs. NoC shows superior performance and scalability than P2P-based network.

in (7) [26].

$$f^r_{ij} = \frac{\lambda^r_{ij}}{\sum^5_{k=1} \lambda^r_{jk}} \tag{7}$$

The contention probability between port $i$ and port $j$ of the router $r$ is computed as $c^r_{ij} = \sum^5_{k=1} f^r_{ik} f^r_{jk}$. Line 10-11 of Algorithm 2 shows the computation of the contention matrix.

Next, the average queue length of each port of the router $r$ ($N^r$) is computed through the technique described in [26].

$$N^r = (I - t\Lambda^r C^r)^{-1} \Lambda^r R, \tag{8}$$

where $t$ is the service time of the router, and we assume $t = 1$ for our evaluation. $R$ is the average residual time and is calculated assuming that the packets arrive in discrete clock cycles [21]. Waiting time of the packets at each port of the router $r$ is computed as $W^r = N^r (\Lambda^r)^{-1}$. End-to-end average latency for each layer $l$ ($L^l_{avg}$) is obtained by averaging the waiting time through all 5 ports ($W^r_{avg}$) of router $r$ and then adding across all routers, as shown in (9) and (10).

$$W^r_{avg} = \frac{1}{5} \sum^5_{p=1} W^r_p \tag{9}$$

$$L^l_{avg} = \sum^R_{r=1} W^r_{avg} \tag{10}$$

Finally, total communication latency ($L^{ana}_{comm}$) is obtained by adding end-to-end average latency for each layer $l$ as:

$$L^{ana}_{comm} = \sum^{N_L}_{l=1} L^l_{avg} \tag{11}$$

## 5 CONNECTION-CENTRIC ARCHITECTURE

In this section, we first discuss a multi-tiled SRAM-based IMC architecture with three different interconnect topologies, namely, P2P, NoC-tree, and NoC-mesh at the tile level. We perform a comprehensive analysis of these three interconnect-based SRAM IMC architectures for different DNNs using the simulation framework described in Section 3. Based on the analysis, we show the need for an NoC-based heterogeneous interconnect IMC architecture for efficient DNN acceleration. We assume all weights are stored on-chip to avoid any DRAM access. The weights are loaded pre-execution and stored on-chip. The inputs are then loaded, and the computation is performed.
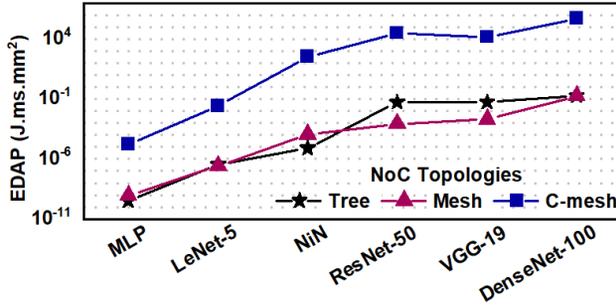
Fig. 9. Comparison of energy-delay-area product (EDAP) of NoC-tree, NoC-mesh, and c-mesh for different DNNs.

There is no re-loading of intermediate results or weights from the off-chip memory during the execution of the DNN. The SRAM buffer is designed large enough to hold the intermediate results on-chip rather than moving them off-chip. Multiple inferences of the images can be performed using one pre-execution loading of the weights. Hence, we do not consider the initial loading of the weights into the energy calculation, consistent with prior work [30, 32] compared in the manuscript. In addition, we adhere to layer-by-layer design instead of a layer-pipelined design, since a pipelined design introduces pipeline bubbles in the execution flow and complicates the control logic [29].

## 5.1 Design Space Exploration

We evaluate different performance metrics for a wide range of DNNs with P2P, NoC-tree, and NoC-mesh-based interconnect for SRAM-based IMC architectures. We consider routers with five ports, one virtual channel for NoCs and X–Y routing for NoC-mesh for this evaluation. To facilitate fair comparison, we normalize the throughput of the hardware architectures with three interconnect topologies to that of P2P interconnect.

Figure 8 shows the throughput comparison for different DNNs. For low connection density DNNs such as MLP and LeNet-5 [17], the choice of interconnect does not make a significant difference to the throughput, due to low data movement between different tiles of the IMC architecture. However, P2P interconnect results in 1.25× and 2× higher area cost than NoC-tree for MLP and LeNet-5, respectively. Hence, NoC-tree provides better overall performance than P2P for both MLP and LeNet-5. We further analyze dense DNNs such as NiN [18], VGG-19, ResNet-50 [6] and DenseNet-100 [8]. The performance comparison shows that the NoC-tree and NoC-mesh-based IMC architectures perform better than the P2P-based architectures (up to 15× for DenseNet-100). Since higher connection density of the DNNs results in increased on-chip data movement, the routing latency dominates the end-to-end latency. We see a similar trend with ReRAM-based IMC architectures with similar throughput for MLP and 15× improvement in throughput for DenseNet-100. Through this, we establish that the performance of the P2P-based IMC architecture (SRAM- or ReRAM-based) diminishes with increasing connection density. In contrast, the performance of the NoC-based (tree, mesh) IMC architecture scales better (Figure 8).

**Exploration of other NoC topologies:** Apart from tree and mesh, the other commonly known NoC topologies include c-mesh, hypercube, and torus. These topologies utilize more resources in terms of routers and links to reduce communication latency. However, the usage of more resources increases power consumption and the area of the NoC. For example, we performed experiments with c-mesh topology for different DNNs. Figure 9 compares energy-delay-area product (EDAP) of mesh-, tree- and c-mesh-based NoC for different NoC. We observe that while mesh- and tree-NoC
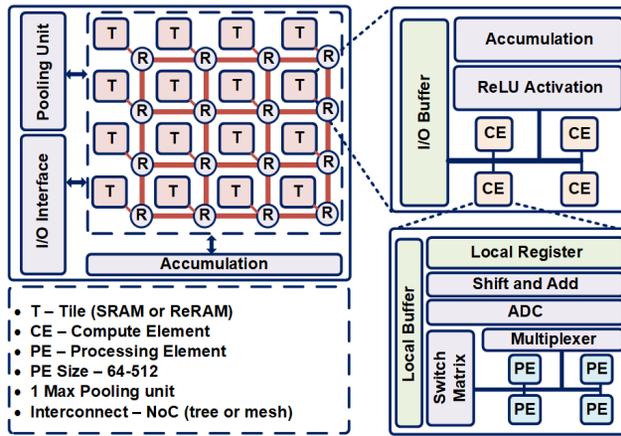
Fig. 10. NoC-based heterogeneous interconnect IMC architecture. A three-level interconnect scheme consisting of NoC (tree or mesh) between tiles, P2P network between CEs, and bus between PEs.

provides comparable EDAP, the same for c-mesh is a minimum of five orders of magnitude higher than mesh- and tree-NoC.

## 5.2 Hardware Architecture

Based on the conclusions from Section 5.1, we derive an NoC-based heterogeneous interconnect IMC architecture for DNN acceleration. Figure 10 shows the hardware architecture which employs the heterogeneous interconnect system.

The proposed architecture is divided into a number of tiles, with each tile having a set of computing elements (CE). The tile architecture includes non-linear activation units, I/O buffer, and accumulators to manage data transfer efficiently. Each CE further consists of multiple processing elements (PE) or crossbar arrays, multiplexers, buffers, a sense amplifier, and flash ADCs. The ADC precision is set to four bits such that there is minimum or no accuracy degradation for DNNs. In addition, the architecture does not utilize a digital-to-analog (DAC) converter; instead, it uses sequential signaling to represent multi-bit inputs [27]. The proposed heterogeneous tile architecture can be used for both SRAM and ReRAM (1T1R) technologies. However, the peripheral circuits change based on the technology. In this work, we choose a homogeneous tile design consisting of four CEs and a CE structure consisting of four PEs. We evaluate both SRAM- and ReRAM-based IMC architectures for PE sizes varying from 64×64 to 512×512. We sample 8 DNNs (LeNet, NiN, SqueezeNet, ResNet-152, ResNet-50, VGG-16, VGG-19, and DenseNet-100) and a crossbar size of 256×256 provides the lowest EDAP for 75% of the DNNs. Hence, in this work, we choose 256×256 as the crossbar size for both SRAM- and ReRAM-based IMC architectures. To maximize performance, the architecture uses heterogeneous interconnects. It employs the NoC-based interconnect on the global tile-level with a P2P interconnect (H-Tree) at the CE-level and bus at the PE-level due to significantly lower data volume. For low data volume, the NoC-based interconnect provides marginal performance gain while increasing energy consumption.

## 6 EXPERIMENTS AND RESULTS

### 6.1 Experimental Setup

We consider an IMC architecture (Figure 10) with a homogeneous tile structure (SRAM, ReRAM) and one NoC router per tile. Table 2 summarizes the design parameters considered. We report
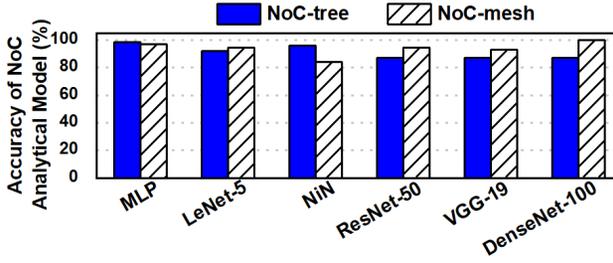
Fig. 11. Accuracy of NoC analytical model for NoC-mesh and NoC-tree with respect to cycle-accurate simulator [10].

the end-to-end latency, chip area, and total energy obtained for a PE size of 256×256 for each of the DNNs using the simulation framework discussed in Section 3. We incorporate conventional mapping [30], IMC SRAM bitcell/array design from [12] and 1T1R ReRAM bitcell/array properties from [2]. The IMC compute fabric utilizes a parallel read-out method. We utilize the same crossbar array size of 256×256 for both SRAM and ReRAM-based IMC architectures. All rows of the IMC crossbar are asserted together, analog MAC computation is performed along the bitline, and the analog voltage/current is digitized with a 4-bit flash ADC at the column periphery. We perform an extensive evaluation of the IMC architecture with both SRAM-based and ReRAM-based PE arrays for both NoC-tree and NoC-mesh. Unless specified, the NoC utilizes one virtual channel, a buffer size (all input and output buffers) of eight, and three router pipeline stages.

Table 2. Summary of design parameters

| PE array size | 256×256 | Read-out Method | Parallel |
|---|---|---|---|
| Technology node | 32nm | Flash ADC resolution | 4 bits |
| Cell levels | 1 bit/cell | Operating frequency | 1 GHz |
| Data precision | 8 bits | NoC bus width | 32 |

## 6.2 Evaluation of NoC Analytical Model

Figure 11 shows the accuracy of the analytical model (presented in Algorithm 2 in Section 4) to estimate the end-to-end communication latency with both NoC-tree and NoC-mesh. We observe that the accuracy is always more than 85% for different DNNs. On an average, the NoC analytical model achieves 93% accuracy with respect to cycle-accurate NoC simulation [10]. Moreover, we achieve 100×-2000× speed-up with the NoC analytical model with respect to cycle-accurate NoC simulation. Figure 12 shows the speed-up for different DNNs with mesh-NoC. This speed-up is useful to perform design space exploration by considering various sizes of PE arrays and other NoC topologies. Due to the high speed-up in NoC performance analysis , we achieve 8× speed-up in overall performance analysis with respect to the framework which uses cycle-accurate NoC simulation.
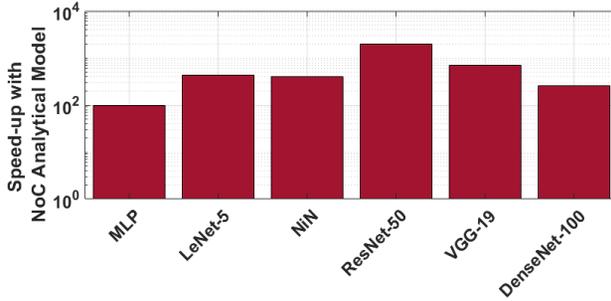
Fig. 12. Speed-up (in NoC simulation) with NoC analytical models with respect to cycle-accurate NoC simulation for different DNNs with mesh-NoC.
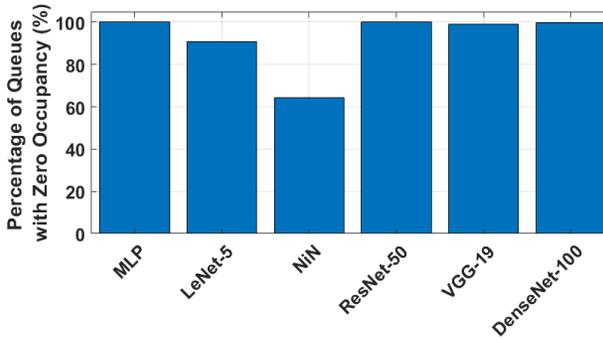


Fig. 13. Percentage of queues with zero occupancy when a new flit arrives.

## 6.3 Analysis on Traffic Congestion in NoC

In this section, we present an analysis on traffic congestion in NoC for various DNNs. To this end we discuss about average queue length of different buffers in the NoC and worst case communication latency.

**Analysis of the average queue length:** Furthermore, we investigated the average queue length at different ports of different routers in the NoC through a cycle-accurate NoC simulator. We performed this experiment with mesh-NoC considering the configuration parameters shown in Table 2. Figure 13 shows that 64%-100% of the queues contain no flit when a new flit arrives for different DNNs. The percentage of queues with zero occupancy for LeNet-5 and NiN is 91% and 65%, respectively. These two DNNs utilize fewer number of routers, which results in less parallelism in data communication. However, we note that determining the optimal number of routers for a given DNN is not a scope of this work.

Figure 14 shows the average queue length for NiN and VGG-19 for the queues with non-zero length when a new flit arrives to the queues. We observe that the average queue length varies from 0.004-0.5 for these DNNs. Average queue length is very low in these cases since the injection rate to the queues are less, and NoC introduces a high degree of parallelism in data transmission between routers.

**Analysis of the worst case latency:** Furthermore, we extracted the worst-case latency ($L_{max}$) for different source to destination pairs of different DNNs with mesh-NoC. We compared $L_{max}$ of each source to destination pair with corresponding average latency ($L_{avg}$). Then we compute mean
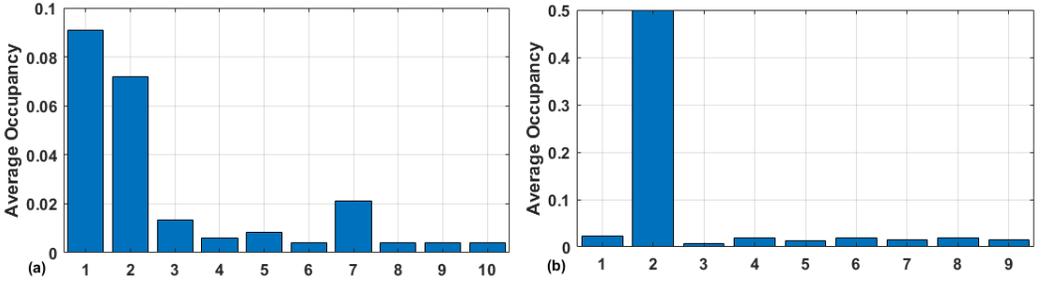
Fig. 14.  Average Occupancy of Queues with non-zero length for (a) NiN, (b) VGG-19.

Table 3.  Mean absolute percentage deviation (*MAPD*) of worst-case NoC latency from average NoC for different DNNs.

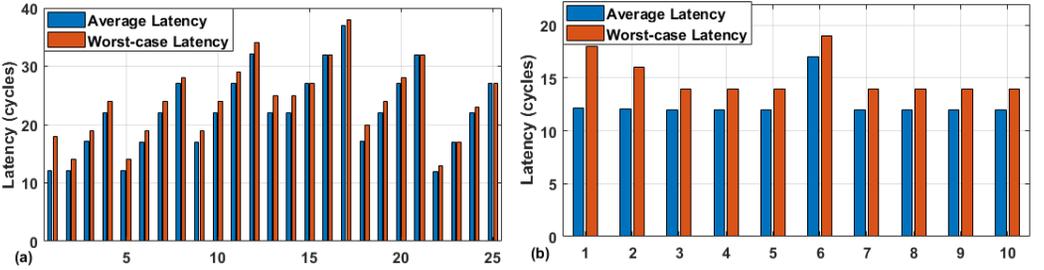| DNNs | MLP | LeNet-5 | NiN | ResNet-50 | VGG-19 | DenseNet-100 |
|---|---|---|---|---|---|---|
| MAPD(%) | 0 | 9.13 | 20.76 | 0 | 0.14 | 0 |



Fig. 15.  Comparison between average latency and worst-case latency for source to destination pairs with non-zero latency for (a) LeNet-5 and (b) NiN.

absolute percentage deviation (MAPD) of $L_{max}$ from $L_{avg}$ as the equation below.

$$MAPD = 100 \times \frac{1}{N} \sum_{i=1}^{N} \frac{(L_{max}^i - L_{avg}^i)}{L_{avg}^i} \qquad (12)$$

Where $N$ is the total number of source to destination pairs with non-zero average latency. $L_{max}^i$ and $L_{avg}^i$ are the worst-case latency and the average latency respectively of $i^{th}$ source to destination pair. Table 3 shows the mean absolute percentage deviation for different DNNs. We observe that the deviation is insignificant, except for LeNet-5 and NiN. The deviations for these two networks are 9.13% and 20.76%, respectively.

Furthermore, in Figure 15 we show the absolute difference between the worst-case latency and the average latency for LeNet-5 and NiN for different source to destination pairs with non-zero latency. The maximum difference is 6 cycles both for LeNet-5 and NiN. This analysis shows that the worst-case latency has very less deviation from the average latency. Therefore, the studies of average queue length and worst-case latency confirm that there is no congestion in the NoC.
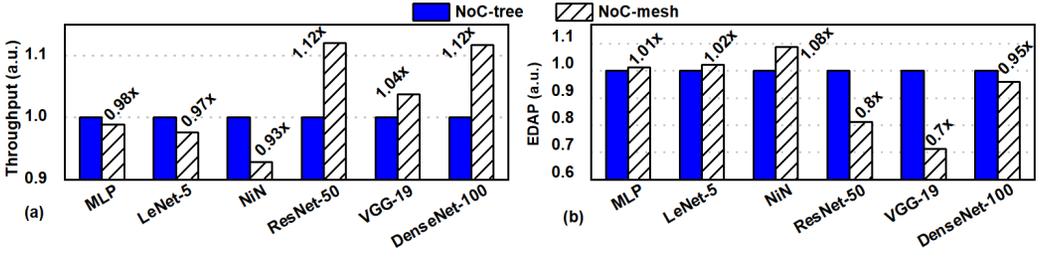
Fig. 16. (a) Normalized throughput and (b) normalized EDAP of NoC-tree and NoC-mesh-based on-chip interconnect for **SRAM-based** IMC architecture for different DNNs. Dense DNNs favor NoC-mesh while NoC-tree is sufficient for shallow DNNs.
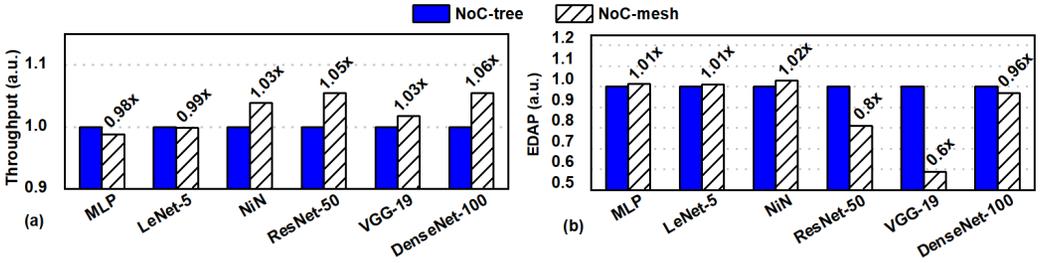


Fig. 17. (a) Normalized throughput and (b) normalized EDAP of NoC-tree and NoC-mesh-based on-chip interconnect for **ReRAM-based** IMC architecture for different DNNs. Dense DNNs favor NoC-mesh while NoC-tree is sufficient for shallow DNNs.

## 6.4 Guidance on Optimal Choice of Interconnect

*6.4.1 Empirical Analysis.* We compare the performance of the IMC architecture using both NoC-tree and NoC-mesh for both SRAM and ReRAM-based technologies. We perform the experiments for representative DNNs. MLP, LeNet-5, and NiN depict low connection density DNNs; ResNet-50, VGG-19, and DenseNet-100 depict high connection density DNNs. We report throughput and the product of energy consumption, end-to-end latency, and area (EDAP) of the IMC architectures. EDAP is used as the metric to guide the optimal choice for the interconnect for IMC architectures.

Figure 16(a) shows the ratio of the throughput of the SRAM-based IMC architecture using NoC-tree and NoC-mesh interconnect. We normalize the throughput values with respect to that of NoC-tree. NoC-tree performs better than the NoC-mesh for DNNs with low connection density. This is because of the reduced injection bandwidth into the interconnect. In addition, while NoC-mesh provides lower interconnect latency than NoC-tree, it comes at an increased area and energy cost. However, NoC-mesh performs better for DNNs with high connection density. The improved performance stems from the reduced interconnect latency for high injection rates of data into the interconnect. The reduction in latency is much higher than the additional overhead due to both area and energy of NoC-mesh.

To better understand the performance, we report the EDAP for the SRAM-based IMC architecture. Figure 16(b) shows normalized EDAP of the NoC-tree and NoC-mesh for both low and high connection density DNNs. DNNs with low connection density have significantly lower EDAP for NoC-tree than that with NoC-mesh. Such an improved EDAP performance for NoC-tree complements the observation for throughput. At the same time, for DNNs with high connection density, the EDAP of NoC-mesh is lower than that of the NoC-tree for IMC architectures. A similar observation is seen

for ReRAM-based IMC architectures as shown in Figure 17(a) and Figure 17(b). In contrast to the SRAM-based IMC architecture, NiN provides better performance in throughput for the NoC-mesh interconnect. At the same time, NoC-tree provides better EDAP compared to NoC-mesh, similar to that of the SRAM-based IMC architecture.

Furthermore, we performed another two sets of experiments with NoC-mesh and NoC-tree by varying the number of virtual channels and bus-width. In this case, we consider ReRAM-based IMC architectures. Figure 18 shows the comparison with different numbers of virtual channels, and Figure 19 shows the comparison with different bus width of the NoC. We observe similar trends for different DNNs with a different NoC configurations.

Since the injection rate to the input buffer of the NoC is always low (less than one packet in 100 cycles), increasing the number of virtual channels does not alter the inference latency significantly. Therefore, throughput remains similar (for all DNNs) both for NoC-tree and NoC-mesh with an increasing number of virtual channels. However, the area and power of both NoC-mesh and NoC-tree increase proportionally with an increasing number of virtual channels. Therefore, the normalized EDAP (EDAP of mesh-NoC divided by the EDAP of tree-NoC) is similar for all DNNs with different numbers of virtual channels.

While we change the bus width of the NoC, the latency increases (decreases) with decreasing (increasing) bus width proportionally, i.e., the latency with a bus width of 32 is twice than the latency with bus width of 64. Moreover, the area and power of the NoC increases (decreases) with increasing (decreasing) bus width proportionally. Therefore, the normalized EDAP is similar for all DNNs with different NoC bus widths. Consequently, for all configurations, we obtain exactly the same guidance on the choice of NoC for different DNNs. Therefore, the guidance is consistent across different parameters of NoCs.

*6.4.2 Theoretical Analysis.* We utilize the analytical model in Section 4 and the experimental results described in Figure 16 and Figure 17 to provide guidance on the optimal choice of interconnect for IMC architectures. The injection rate at each port of an NoC router for each layer of the DNN is expressed in (6). The numerator of (6) denotes the total data volume between $i^{th}$ layer and $(i + 1)^{th}$ layer for each port of the router per cycle. This is divided by $(T_i \times T_{i+1} \times W)$ to obtain the injection rate from for each port of every router as detailed in Section 4. For a fixed NoC-based IMC architecture, target throughput ($FPS$), frequency of operation ($freq$), and bus width ($W$) remain constant. Hence, from (6) we obtain,

$$\lambda_i \propto \frac{A_i \times N_{bits}}{T_i \times T_{i+1}} \tag{13}$$

Let the connection density for $i^{th}$ layer be $\rho_i$ and the number of neurons be $\mu_i$. Data volume between $i^{th}$ and $(i + 1)^{th}$ layer is proportional to the product of $\rho_i$ and $\mu_i$, as shown in (14).

$$(A_i \times N_{bits}) \propto (\rho_i \times \mu_i) \tag{14}$$

Additionally, the number of tiles in $i^{th}$ layer is directly proportional to $\mu_i$. Hence, from (13) and (14) we get,

$$\lambda_i \propto \frac{\rho_i \times \mu_i}{\mu_i \times \mu_{i-1}} = \frac{\rho_i}{\mu_{i-1}} \tag{15}$$

Generalising (15), we obtain,

$$\lambda \propto \frac{\rho}{\mu} \tag{16}$$

Therefore, the injection rate is directly proportional to the connection density and inversely proportional to the number of neurons of the DNN. Figure 20 presents the preferred regions for
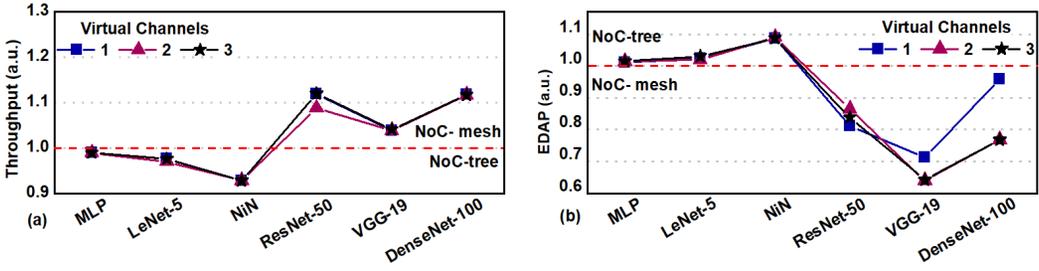
Fig. 18. Assessment of (a) throughput and (b) EDAP between NoC-tree and NoC-mesh with different numbers of virtual channels for different DNNs. The throughput is normalized to that of NoC-tree. The preferred NoC topology for optimal performance is shown for the regions above and below the red line.
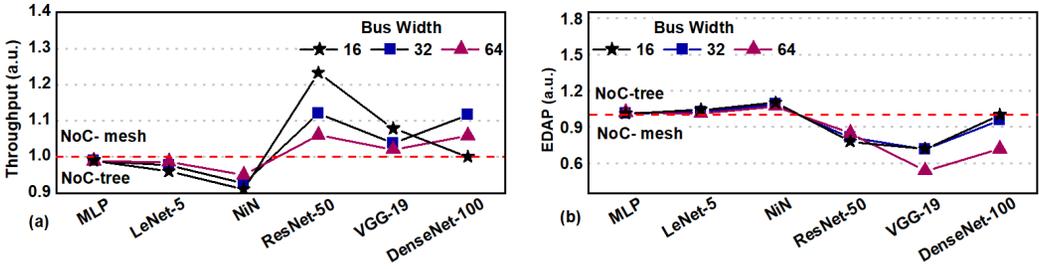


Fig. 19. Assessment of (a) throughput and (b) EDAP between NoC-tree and NoC-mesh with different bus width for different DNNs. The throughput is normalized to that of NoC-tree. The preferred NoC topology for optimal performance is shown for the regions above and below the red line.
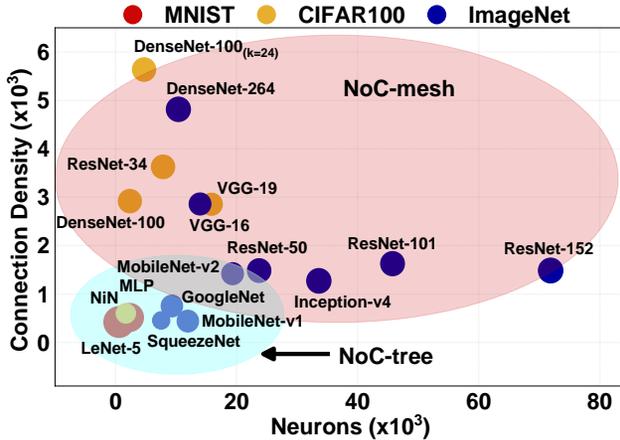


Fig. 20. Optimal NoC topology for IMC architectures for different DNNs.

NoC-tree and NoC-mesh for best throughput for different DNNs with IMC architectures. If the connection density of a DNN is more than $2 \times 10^3$, then NoC-mesh is suitable. If the connection density is less than $1 \times 10^3$, then NoC-tree is appropriate. Both NoC-tree and NoC-mesh are suitable for the DNNs with connection density in the range of $1 \times 10^3$-$2 \times 10^3$ (the region where red and blue ovals overlap in Figure 20).

Table 4. Inference Performance Results for VGG-19. *Reported in [29]

|  | Latency | Power/frame (W/frame) | FPS | EDAP (J.ms.mm$^2$) |
|---|---|---|---|---|
| Proposed-SRAM | 0.68 | 1.96 | 1458 | 0.46 |
| Proposed-ReRAM | 1.49 | 0.43 | 670 | 0.28 |
| AtomLayer [29] | 6.92 | 4.8 | 145 | 1.58 |
| PipeLayer [32] | 2.6* | 168.6 | 385 | 94.17 |
| ISAAC [30] | 8.0* | 65.8 | 125 | 359.64 |

## 6.5 Comparison with state-of-the-art architectures

Table 4 compares the proposed architecture with state-of-the-art DNN accelerators. Prior works show the efficacy of their ReRAM-based IMC architectures for VGG-19 DNN [29, 30, 32]. Hence for comparison, we choose VGG-19 network as the representative DNN. Moreover, we compare the dynamic power consumption of the DNN hardware since prior work utilizes dynamic power in their results, hence making the comparison consistent. The inference latency of the proposed architecture with SRAM arrays is 2.2× lower than the architecture with ReRAM arrays. The proposed ReRAM-based architecture achieves 4.7× improvement in FPS and 6× improvement in EDAP than AtomLayer [29]. The improvement in performance is attributed to the optimal choice of interconnect coupled with the absence of off-chip accesses. The proposed ReRAM-based architecture consumes 400× lower power per frame along with 1.74× improvement in FPS than PipeLayer [32]. Moreover, there is a 5.4× improvement in inference latency compared to ISAAC [30], which is achieved by the heterogeneous interconnect structure.

## 6.6 Connection Density and Hardware Performance

Figure 1 showed a trend of DNNs moving toward a high connection density structure for performance and low connection density structure for compact models. Figure 21 shows the performance for both P2P and NoC-based interconnect at the tile-level for IMC architecture for DNNs with different connection density. We observe a steep increase in total latency with a P2P interconnect. However, the IMC architecture with NoC interconnect shows a stable curve as we move towards high connection density DNNs. With the advent of neural architecture search (NAS) techniques [34, 37], DNNs are moving towards a highly branched structure with very high connection densities. Hence, the NoC-based heterogeneous interconnect architecture provides a scalable and suitable platform for IMC acceleration of DNNs.

## 7 CONCLUSION

The trend of connection density in modern DNNs requires a re-evaluation of the underlying interconnect architecture. Through a comprehensive evaluation, we demonstrate that the P2P-based interconnect is incapable of handling the high volume of on-chip data movement of DNNs. Further, we provide guidance backed by empirical and analytical results to select the appropriate NoC topology as a function of the connection density and the number of neurons. We conclude that NoC-mesh is preferred for DNNs with high connection density, while NoC-tree is suitable for DNNs with low connection density. Finally, we show that the NoC-based heterogeneous interconnect IMC architecture achieves 6× lower EDAP than state-of-the-art ReRAM-based IMC accelerators.
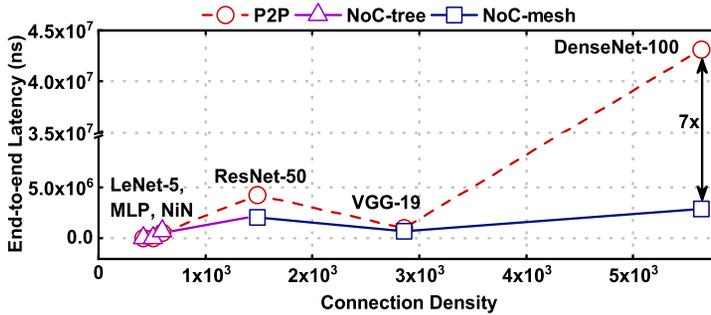
Fig. 21. Appropriate selection of NoC topology significantly improves performance for both SRAM- and ReRAM-based IMC architectures.

## 8 ACKNOWLEDGEMENT

## REFERENCES

[1] Niket Agarwal, Tushar Krishna, Li-Shiuan Peh, and Niraj K Jha. 2009. GARNET: A Detailed On-Chip Network Model Inside a Full-System Simulator. In *2009 IEEE ISPASS*. IEEE, 33–42.

[2] Pai-Yu Chen, Xiaochen Peng, and Shimeng Yu. 2018. NeuroSim: A Circuit-level Macro Model for Benchmarking Neuro-Inspired Architectures in Online Learning. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 37, 12 (2018), 3067–3080.

[3] Yu-Hsin Chen, Tien-Ju Yang, Joel Emer, and Vivienne Sze. 2019. Eyeriss v2: A Flexible Accelerator for Emerging Deep Neural Networks on Mobile Devices. *IEEE JETCAS* 9, 2 (2019), 292–308.

[4] Li Deng, Geoffrey Hinton, and Brian Kingsbury. 2013. New Types of Deep Neural Network Learning for Speech Recognition and Related Applications: An Overview. In *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 8599–8603.

[5] Xiangyu Dong, Cong Xu, Yuan Xie, and Norman P Jouppi. 2012. Nvsim: A Circuit-level Performance, Energy, and Area Model for Emerging Nonvolatile Memory. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 31, 7 (2012), 994–1007.

[6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 770–778.

[7] Mark Horowitz. 2014. Computing's Energy Problem (and What We Can Do About It). In *IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*. 10–14.

[8] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely Connected Convolutional Networks. In *IEEE Conference on Computer Vision and Pattern Recognition*. 4700–4708.

[9] James Jeffers, James Reinders, and Avinash Sodani. 2016. *Intel Xeon Phi Processor High Performance Programming: Knights Landing Edition*. Morgan Kaufmann.

[10] Nan Jiang, Daniel U Becker, George Michelogiannakis, James Balfour, Brian Towles, David E Shaw, John Kim, and William J Dally. 2013. A Detailed and Flexible Cycle-Accurate Network-on-Chip Simulator. In *IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. 86–96.

[11] Z. Jiang, S. Yin, J. Seo, and M. Seok. 2020. C3SRAM: An In-Memory-Computing SRAM Macro Based on Robust Capacitive Coupling Computing Mechanism. *IEEE Journal of Solid-State Circuits* (2020), 1–1.

[12] Win-San Khwa, Jia-Jing Chen, Jia-Fang Li, Xin Si, En-Yu Yang, Xiaoyu Sun, Rui Liu, Pai-Yu Chen, Qiang Li, Shimeng Yu, et al. 2018. A 65nm 4Kb algorithm-dependent computing-in-memory SRAM unit-macro with 2.3 ns and 55.8 TOPS/W fully parallel product-sum operation for binary DNN edge processors. In *2018 IEEE International Solid-State Circuits Conference-(ISSCC)*. IEEE, 496–498.

[13] Abbas Eslami Kiasari, Zhonghai Lu, and Axel Jantsch. 2012. An Analytical Latency Model for Networks-on-Chip. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 21, 1 (2012), 113–123.

[14] Gokul Krishnan, Sumit K Mandal, Chaitali Chakrabarti, Jae-sun Seo, Umit Y Ogras, and Yu Cao. 2020. Interconnect-aware area and energy optimization for in-memory acceleration of DNNs. *IEEE Design & Test* 37, 6 (2020), 79–87.

[15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*. 1097–1105.

[16] Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna. 2018. Maeri: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects. In *ACM SIGPLAN Notices*, Vol. 53. 461–475.

[17] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. 1998. Gradient-based learning applied to document recognition. *Proc. of the IEEE* 86, 11 (1998), 2278–2324.

[18] Min Lin, Qiang Chen, and Shuicheng Yan. 2013. Network in Network. *arXiv preprint arXiv:1312.4400* (2013).

[19] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I Sánchez. 2017. A Survey on Deep Learning in Medical Image Analysis. *Medical image analysis* 42 (2017), 60–88.

[20] Sumit K Mandal, Raid Ayoub, Michael Kishinevsky, Mohammad M Islam, and Umit Y Ogras. 2020. Analytical Performance Modeling of NoCs under Priority Arbitration and Bursty Traffic. *IEEE Embedded Systems Letters* (2020).

[21] Sumit K Mandal, Raid Ayoub, Michael Kishinevsky, and Umit Y Ogras. 2019. Analytical Performance Models for NoCs with Multiple Priority Traffic Classes. *ACM Transactions on Embedded Computing Systems (TECS)* 18, 5s (2019), 1–21.

[22] Sumit K Mandal, Gokul Krishnan, Chaitali Chakrabarti, Jae-Sun Seo, Yu Cao, and Umit Y Ogras. 2020. A Latency-Optimized Reconfigurable NoC for In-Memory Acceleration of DNNs. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* 10, 3 (2020), 362–375.

[23] Manqing Mao, Xiaochen Peng, Rui Liu, Jingtao Li, Shimeng Yu, and Chaitali Chakrabarti. 2019. MAX2: An ReRAM-based Neural Network Accelerator that Maximizes Data Reuse and Area Utilization. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* (2019).

[24] Mohammad Mirza-Aghatabar, Somayyeh Koohi, Shaahin Hessabi, and Massoud Pedram. [n.d.]. An Empirical Investigation of Mesh and Torus NoC Topologies under Different Routing algorithms and Traffic Models. In *10th Euromicro conference on digital system design architectures, methods and tools (DSD 2007)*. 19–26.

[25] Seyed Morteza Nabavinejad, Mohammad Baharloo, Kun-Chih Chen, Maurizio Palesi, Tim Kogel, and Masoumeh Ebrahimi. 2020. An Overview of Efficient Interconnection Networks for Deep Neural Network Accelerators. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* 10, 3 (2020), 268–282.

[26] Umit Y Ogras, Paul Bogdan, and Radu Marculescu. 2010. An Analytical Approach for Network-on-Chip Performance Analysis. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 29, 12 (2010), 2001–2013.

[27] Xiaochen Peng, Minkyu Kim, Xiaoyu Sun, Shihui Yin, Titash Rakshit, Ryan M Hatcher, Jorge A Kittl, Jae-sun Seo, and Shimeng Yu. 2019. Inference Engine Benchmarking Across Technological Platforms from CMOS to RRAM. In *Proceedings of the International Symposium on Memory Systems*. 471–479.

[28] Zhi-Liang Qian, Da-Cheng Juan, Paul Bogdan, Chi-Ying Tsui, Diana Marculescu, and Radu Marculescu. 2015. A Support Vector Regression (SVR)-based Latency Model for Network-on-Chip (NoC) Architectures. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 35, 3 (2015), 471–484.

[29] Ximing Qiao, Xiong Cao, Huanrui Yang, Linghao Song, and Hai Li. 2018. Atomlayer: A Universal Reram-based CNN Accelerator with Atomic Layer Computation. In *IEEE/ACM DAC*.

[30] Ali Shafiee, Anirban Nag, Naveen Muralimanohar, Rajeev Balasubramonian, John Paul Strachan, Miao Hu, R Stanley Williams, and Vivek Srikumar. 2016. ISAAC: A Convolutional Neural Network Accelerator with in-situ Analog Arithmetic in Crossbars. *Proceedings of the 43rd International Symposium on Computer Architecture* (2016).

[31] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[32] Linghao Song, Xuehai Qian, Hai Li, and Yiran Chen. 2017. Pipelayer: A Pipelined Reram-based Accelerator for Deep Learning. In *IEEE International Symposium on High Performance Computer Architecture (HPCA)*. 541–552.

[33] Swagath Venkataramani, Ashish Ranjan, Subarno Banerjee, Dipankar Das, Sasikanth Avancha, Ashok Jagannathan, Ajaya Durg, Dheemanth Nagaraj, Bharat Kaul, Pradeep Dubey, et al. 2017. Scaledeep: A Scalable Compute Architecture for Learning and Evaluating Deep Networks. *ACM SIGARCH Computer Architecture News* 45, 2 (2017).

[34] Saining Xie, Alexander Kirillov, Ross Girshick, and Kaiming He. 2019. Exploring Randomly Wired Neural Networks for Image Recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 1284–1293.

[35] Shihui Yin, Zhewei Jiang, Minkyu Kim, Tushar Gupta, Mingoo Seok, and Jae-sun Seo. 2019. Vesti: Energy-Efficient In-Memory Computing Accelerator for Deep Neural Networks. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 28, 1 (2019), 48–61.

[36] Zhenhua Zhu, Hanbo Sun, Kaizhong Qiu, Lixue Xia, Gokul Krishnan, Guohao Dai, Dimin Niu, Xiaoming Chen, X Sharon Hu, Yu Cao, et al. 2020. MNSIM 2.0: A Behavior-Level Modeling Tool for Memristor-based Neuromorphic Computing Systems. In *Proceedings of the 2020 on Great Lakes Symposium on VLSI*. 83–88.

[37] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. 2018. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8697–8710.