

Efficient and Fault-Tolerant Memristive Neural Networks with *In-Situ* Training

Santlal Prajapati, Manobendra Nath Mondal, and Susmita Sur-Kolay

Abstract—Neuromorphic architectures, which incorporate parallel and in-memory processing, are crucial for accelerating artificial neural network (ANN) computations. This work presents a novel memristor-based multi-layer neural network (memristive MLNN) architecture and an efficient in-situ training algorithm. The proposed design performs matrix-vector multiplications, outer products, and weight updates in constant time $\mathcal{O}(1)$, leveraging the inherent parallelism of memristive crossbars. Each synapse is realized using a single memristor, eliminating the need for transistors, and offering enhanced area and energy efficiency. The architecture is evaluated through LTspice simulations on the IRIS, NASA Asteroid, and Breast Cancer Wisconsin datasets, achieving classification accuracies of 98.22%, 90.43%, and 98.59%, respectively. Robustness is assessed by introducing stuck-at-conducting-state faults in randomly selected memristors. The effects of nonlinearity in memristor conductance and a 10% device variation are also analyzed. The simulation results establish that the network's performance is not affected significantly by faulty memristors, non-linearity, and device variation.

Index Terms—Neuromorphic computing, memristor, memristive synapse, memristive neural network, *in-situ* training.

I. INTRODUCTION

A memristive neuromorphic computing system excels in real-time data processing due to its capabilities of parallel and in-memory computation, which overcome the memory-wall bottleneck of von Neumann architectures [2], [3]. Memristors, as non-linear passive circuit elements, are promising synaptic devices in neuromorphic circuits because of their tunable nonvolatile conductance states, similar to biological synapses. The conductance of a memristor in a memristive synapse represents synaptic weight [4]. Various memristor-based artificial synapses have been proposed, including memristor bridges [5], [6], two transistors and one memristor (2T1M) [7], one transistor and one memristor (1T1M) [8], two memristors (2M) [9], [10], [1], and one memristor (1M) [11], [12].

Memristive multi-layer neural networks (memristive MLNNs) have been constructed by using these synapses on memristive crossbars (MCBs) [13]. These MCBs act as memory by storing synaptic weight matrices, and also as in-memory processors by performing matrix-vector

multiplication [14], vector-vector outer product [7], and weight matrix updates [7] on the MCBs.

Memristive MLNNs [15], [1], [16], [17], [18], [12], [7], [10], [8], [19], [9] are trained either *ex-situ* or *in-situ*. In *ex-situ* training, mirror neural networks are trained on external circuits or software, and the final weights are set in the MCBs of memristive MLNNs [15], [20], [19]. The *in-situ* training is performed directly on the neuromorphic hardware, which is more energy-efficient and faster due to reduced communication with external hardware [20]. Moreover, it overcomes hardware imperfections [9] and is thus a suitable alternative to address memristor imperfections and training latency.

For memristive MLNNs, the key factors include the size of artificial synapses and the time required to modify weights during in-situ training. The power consumption and physical area of transistors are much greater than those of memristors [21] which encourages avoiding using transistors in synapses. In [22], a perceptron neural network is proposed where each synapse consists of four memristors. The studies in [7], [23] use two transistors per synapse (2T1M type) but update the conductances of all the memristors in an MCB in constant time. The works in [8], [24] have reduced the MCB size of MLNNs to 1T1M synapses, achieving a constant-update time of MCB during in-situ training. For space and energy efficiency, the 2M type synapses without transistors have been proposed in [9], [10]. But during training, the time taken to update the memristors in the crossbar is longer, namely one crossbar line at a time in [9] and one memristor at a time in [10]. The memristive MLNN [1] with 2M synapse achieved constant time update of the weight matrix during in-situ training. However, to reduce the size of the MCB (1M type), Zhang *et al.* [12] proposed an artificial synapse with only one memristor, but their training method updates the conductances of the memristors of an MCB one by one, making the training time-inefficient.

Our proposed memristive MLNN architecture with *in-situ* training method focuses on reducing the memristive crossbar area (1M type) as well as improving the crossbar update time to $\mathcal{O}(1)$. Additionally, this work scrutinizes the fault tolerance, the effects of device variation, and the non-linearity present in the conductance of a memristor on the performance of memristive MLNNs.

Table I qualitatively compares memristive MLNN and in-situ training with earlier works and shows that the proposed work reduces the size of an artificial synapse to one memristor as well as the time to update all the memristors in a crossbar to $\mathcal{O}(1)$.

Our main contributions are as follows:

S. Prajapati, School of Engineering, Sister Nivedita University, DG 1/2 New Town, Action Area 1, Kolkata - 700156, E-mail: santlal.p@snu.ac.in

S. Sur-Kolay, ACM Unit, Indian Statistical Institute, Kolkata-700108, India, ssk@isical.ac.in

M. N. Mondal is now with School of Computer Science, UPES, Dehradun, India, E-mail: manobmondal@gmail.com

All three authors were in Indian Statistical Institute, Kolkata-700108, India, when this work was done. An earlier version of this paper appeared in the Proc. of SOCC 2022 [1]

Table I: Qualitative comparison of our proposed work with earlier in-situ training methods.

BAM: bidirectional associative memory, MLNN: multi-layered neural network, BP: back propagation

Work	Synapse type	Type of network and its input	Training	Update on weight	Memristor's weight update voltage	Non-linearity analysis	Fault analysis
Zhang <i>et al.</i> [12]	1M	With digital input, MLNN	Adaptive BP	One memristor at a time	Constant voltage amplitude, weight adjustment encode by duration	No	No
Li <i>et al.</i> [9]	2M	Analog input MLNN	General BP	One memristor line at a time	Fixed duration time, weight adjustment, encode by voltage amplitude	No	Yes
Krestinskaya <i>et al.</i> [10]	2M	ANN with digital or analog input	General BP	One memristor at a time	-do-	No	No
Shi <i>et al.</i> [8]	1T1M	BAM with digital and analog input	Network specific algorithm [25]	All memristors at a time	Fixed voltage amplitude with feedback duration time control	No	No
Feng <i>et al.</i> [24]	1T1M	MLNN with analog input	General BP	All memristors at a time	-do-	No	No
Soudry <i>et al.</i> [7]	2T1M	Analog input MLNN	General BP	All memristors at a time	-do-	No	No
Yan <i>et al.</i> [23]	2T1M	Analog input MLNN	Momentum and adaptive learning	All memristors at a time	-do-	No	No
Prajapati <i>et al.</i> [1]	2M	Analog input MLNN	General BP	All memristors at a time	Variable voltage amplitude with feedback duration time control	No	No
Proposed work	1M	Analog input MLNN	General BP	All memristors at a time	-do-	Yes	Yes

- 1) a memristive feed-forward multi-layer neural network architecture, demonstrating resource and space efficiency by utilizing only one memristor for each synapse without any transistors,
- 2) an innovative *in-situ* training algorithm based on online gradient descent backpropagation, enabling constant-time weight-matrix update without any additional storage,
- 3) a scheme for encoding the input and control signals for the inference and training of the memristive MLNNs,
- 4) validation of this memristive MLNN architecture and its *in-situ* training method, and the study of its robustness in the presence of (i) device variations, (ii) stuck-at faults, (iii) non-linearity in the conductance of memristors (iv) sneak path issues, and (v) scalability of our architecture.

In Section II, the preliminaries of the online gradient descent backpropagation algorithm are outlined. The details of the proposed synapse and neuron in a memristive crossbar and the required peripherals, i.e., encoder and control signals for switches, along with a complete memristive MLNN, are presented in Section III. The principles of inference and update operations on a single MCB and then the relevant algorithms for a memristive MLNN are explained in Section IV. Simulation results are given in Section V, with concluding remarks in Section VI. All boldface lowercase and uppercase letters represent vectors and matrices, respectively.

II. BACKGROUND

Multi-layer neural networks (MLNNs) are trained to learn the data environment with data points. Supervised online gradient descent backpropagation learning algorithm [26] trains MLNNs where learning is guided by a loss function. An MLNN comprises an input and an output layer of neurons with one or more hidden layers of neurons. The training process involves two stages: forward or feed forward and backpropagation. During forward propagation, the input data is propagated layer by layer from the input layer to the sequence of hidden layers, and finally to the output layer. Next, the gradient or error computed from the loss function

is back propagated from the output layer through the hidden layers to the input layer. Each neuron has two states, namely pre-activation and activation. In the pre-activation state, the weighted sum of the inputs is calculated while in the activation state, a non-linear function is applied to this sum.

For clarity, let us consider a single-layer neural network (SLNN) (Fig. 1(a)) with n input neurons and m output neurons coupled by an $m \times n$ synaptic weight matrix W ; there are no hidden layers. An input feature vector $\mathbf{x} \in \mathbb{R}^n$ is fed to the SLNN to produce the desired output $\mathbf{d} \in \mathbb{R}^m$, where \mathbf{d} is a one-hot vector (for classification problems but may be different for other types). For a given input \mathbf{x} and its true class label l , all the elements of \mathbf{d} are zero except at position l . At the pre-activation state, matrix-vector multiplication is performed to obtain the pre-activation vector $\mathbf{r} \in \mathbb{R}^m$ where $\mathbf{r} = \mathbf{W}\mathbf{x}$, having $r_j = \sum_{i=1}^n w_{j,i}x_i$ for an output neuron n_j [26].

In order to obtain the forward inferred vector $\mathbf{o} \in \mathbb{R}^m$ at the output layer, a function $f: \mathbb{R}^m \rightarrow \mathbb{R}^m$ is applied to r_j of the output neuron n_j whose inferred value is the output o_j . For classification, f is typically *soft-max*, so $\forall j = 1$ to m $o_j = f(r_j) = \frac{e^{r_j}}{\sum_{i=1}^m e^{r_i}}$. The loss/cost function L defined as $L \triangleq -\sum_{j=1}^m d_j \log(o_j) = -\log(o_l)$ is *cross entropy* where l is the true class label. Other cost functions, for example, mean squared error (MSE) may also be used.

The goal of training is to minimize L by updating \mathbf{W} which may be chosen randomly at the beginning. The gradient of L w.r.t \mathbf{r} , given by the vector $\nabla_{\mathbf{r}}^L = -(\mathbf{d} - \mathbf{o})$, is calculated and back propagated to update W [26]. By defining $\mathbf{y}^{(\text{output})} \triangleq \mathbf{d} - \mathbf{o}$ as an error vector at the output layer, $\nabla_{\mathbf{r}}^L$ is simplified to $\nabla_{\mathbf{r}}^L = -\mathbf{y}^{(\text{output})}$. The j^{th} element y_j of $\mathbf{y}^{(\text{output})}$ denotes the error at neuron n_j . The gradient of L with respect to the weight matrix \mathbf{W} , denoted as $\nabla_{\mathbf{W}}^L$, is written in terms of the error vector $\mathbf{y}^{(\text{output})}$ and input \mathbf{x} as the matrix $\nabla_{\mathbf{W}}^L = -(\mathbf{d} - \mathbf{o})\mathbf{x}^T = -\mathbf{y}^{(\text{output})}\mathbf{x}^T$ [26]. The weight matrix \mathbf{W} with learning rate η is updated as

$$\mathbf{W}(\text{new}) = \mathbf{W}(\text{old}) + \Delta \mathbf{W}, \text{ where}$$

$$\Delta \mathbf{W} \propto -\nabla_{\mathbf{W}}^L = -\eta \nabla_{\mathbf{W}}^L = \eta \mathbf{y}^{(\text{output})} \mathbf{x}^T \quad (\text{II.1})$$

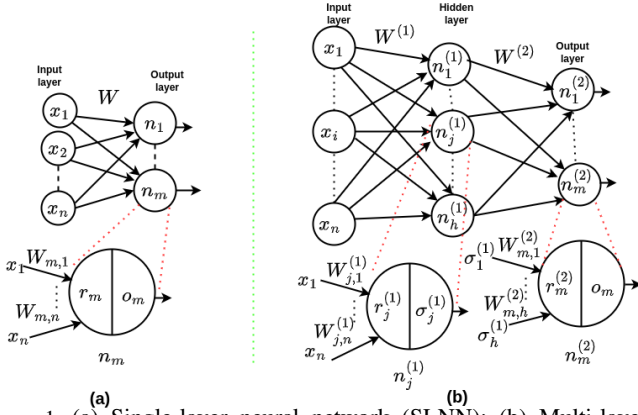


Figure 1: (a) Single-layer neural network (SLNN); (b) Multi-layer neural network (MLNN) with one hidden layer.

Hence, the change of weight matrix $\Delta \mathbf{W}$ is proportional to the outer product of the vectors $\mathbf{y}^{(\text{output})}$ and \mathbf{x} . For a single neuron n_j , $\Delta w_{j,i} = -\eta \nabla_{w_{j,i}}^L = \eta y_j^{(\text{output})} x_i$.

Now, consider an MLNN as in Fig. 1(b) with one hidden layer that contains h hidden neurons. An $h \times n$ $\mathbf{W}^{(1)}$ connecting the input-hidden layers and an $m \times h$ matrix $\mathbf{W}^{(2)}$ connecting the hidden-output layers are the two weight matrices. In order to get the activation vector $\sigma^{(1)}$ in the activation state at the first/hidden layer, a nonlinear activation function σ is applied on the pre-activation vector $\mathbf{r}^{(1)}$ of the 1st layer to obtain $\sigma^{(1)}$, which is then fed to the output neurons through $\mathbf{W}^{(2)}$. The superscripts in the notation denote layer name/number. The output layer (in this case, the 2nd layer) considers the *soft-max* function f as an output function for classification. The four operations $\mathbf{r}^{(1)} = \mathbf{W}^{(1)}\mathbf{x}$, $\sigma^{(1)} = \sigma(\mathbf{r}^{(1)})$, $\mathbf{r}^{(2)} = \mathbf{W}^{(2)}\sigma^{(1)}$, and $\mathbf{o} = \mathbf{f}(\mathbf{r}^{(2)})$ are performed in a sequence to infer the output vector $\mathbf{o} \in \mathbb{R}^m$ of the MLNN for a given input $\mathbf{x} \in \mathbb{R}^n$.

The gradient of L with respect to $\mathbf{r}^{(2)}$, i.e., the vector $\nabla_{\mathbf{r}^{(2)}}^L = -\mathbf{y}^{(\text{output})}$, is similar to that for SLNN where $\mathbf{y}^{(\text{output})}$ is the error vector at the output layer. The change in weights of the matrix $\mathbf{W}^{(2)}$ is given by the matrix $\Delta \mathbf{W}^{(2)} = -\eta \nabla_{\mathbf{W}^{(2)}}^L = \eta \mathbf{y}^{(\text{output})} \sigma^{(1)\text{T}}$. It is to be noted that the error vectors at the hidden and output layers are different. The gradient of L w.r.t $\sigma^{(1)}$ is the vector $\nabla_{\sigma^{(1)}}^L = -(\mathbf{W}^{(2)})^T \mathbf{y}^{(\text{output})} = -\delta^{(2)}$ [26]; and that w.r.t $\mathbf{r}^{(1)}$ is the vector $\nabla_{\mathbf{r}^{(1)}}^L = -\delta^{(2)} \odot \sigma'^{(1)}$ [26] where $\sigma'(\mathbf{x})_i = d\sigma(x_i)/dx_i$ and \odot denotes the element-wise product of two equal-sized vectors. The error vector $\mathbf{y}^{(1)}$ at hidden layer is defined as $\mathbf{y}^{(1)} = \delta^{(2)} \odot \sigma'^{(1)}$, and the gradient of L w.r.t $\mathbf{r}^{(1)}$ in terms of $\mathbf{y}^{(1)}$ is given by $\nabla_{\mathbf{r}^{(1)}}^L = -\mathbf{y}^{(1)}$. The gradient of loss function L w.r.t $\mathbf{W}^{(1)}$ is the matrix $\nabla_{\mathbf{W}^{(1)}}^L = -\mathbf{y}^{(1)} \mathbf{x}^T$. Hence, the final update equations for matrices $\mathbf{W}^{(2)}$ and $\mathbf{W}^{(1)}$ are given by

$$\mathbf{W}^{(2)}(\text{new}) = \mathbf{W}^{(2)}(\text{old}) + \eta \mathbf{y}^{(\text{output})} \sigma^{(1)\text{T}} \quad (\text{II.2})$$

$$\mathbf{W}^{(1)}(\text{new}) = \mathbf{W}^{(1)}(\text{old}) + \eta \mathbf{y}^{(1)} \mathbf{x}^T \quad (\text{II.3})$$

III. MEMRISTIVE NEURAL NETWORK ARCHITECTURE

An overview of memristive SLNN architecture is shown in Fig. 2(a). Fig. 2(b) illustrates the memristive MLNN with

one hidden layer which is obtained by cascading memristive SLNNs. The entire computation mainly consists of forward inference, backward inference, and weight update. The computation of vectors either $\sigma^{(k)}$ at k^{th} hidden layer or \mathbf{o} at the output layer is called forward inference. During backward inference at layer k , the vector $\delta^{(k)} = (\mathbf{W}^{(k)})^T \mathbf{y}^{(k)}$ as defined in Section II is calculated while the weight matrix, i.e., the conductance of the memristors in the MCB are updated.

In forward inference, during the period 0 to T_{rd} , the input feature vector \mathbf{x} is encoded into appropriate voltage waveform by the input encoder and fed to the MCB via left switches. This encoded input vector is multiplied with the matrix of the conductance of the memristors in the MCB (in-memory computation) and the result is passed onto the column output decoder via bottom switches. The column output decoder produces either \mathbf{o} or σ based on layer types.

The vector $\delta^{(k)}$ (Fig. 2(b)) is computed during the time period T_{rd} to $2T_{rd}$ only in MLNNs. The error vector $\mathbf{y}^{(k)}$ is encoded into proper voltages by the error encoder and fed back to the k^{th} MCB via top switches. Then in the MCB, the multiplication of the transposed weight matrix and $\mathbf{y}^{(k)}$ is performed and the result is passed, via right switches, onto the row output decoder to produce vector $\delta^{(k)}$. The error vector $\mathbf{y}^{(k-1)}$ for $(k-1)^{\text{th}}$ layer is obtained by element-wise multiplication (denoted by \odot) of $\delta^{(k)}$ and $\sigma'^{(k-1)}$. While training during the period $2T_{rd}$ to $2T_{rd} + T_{wr}$, the conductances of the memristors in the MCB are updated.

The arrows in Fig. 2 indicate the flow of information. The information flow and update operation are controlled by voltage-controlled switches around the MCB. At k^{th} layer, the switch control interface with $\mathbf{y}^{(k)}$ produces controlling voltage signals to control them. The top and right switches are optional in SLNN and the first layer of MLNN. The details of our proposed memristive neural network architecture are explained in the subsequent sections:

- synapses and neurons in a memristive crossbar,
- input encoder,
- switch control signals for an MCB,
- memristive multi-layer neural network.

A. Synapses and neurons in a memristive crossbar

A memristor is a two-terminal circuit element having the variable conductance states which depend on a state variable or a set of state variables [27]. In this work, only those memristors are considered that show voltage threshold behavior, i.e., changes in their conductances are negligible until they experience voltages greater than the positive threshold voltage (v_{th}^+), or less than the negative threshold voltage (v_{th}^-). In this work, the generalized memristor spice model [28] is used. In Fig. 3, the characteristics of a silver chalcogenide-based memristive device [29] having positive and negative thresholds of +0.16V and -0.15V respectively, are shown. Fig. 3(a) shows the variation of conductance of the memristor when voltage pulses, shown in Fig. 3(b), of amplitude greater than v_{th}^+ and less than v_{th}^- are applied across it. The I-V curve of the memristor is shown in Fig. 3(c).

The structure, where memristors are arranged in a 2D grid as shown in the red dotted box in Fig. 4(a), is the memristive

¹All detailed derivations are in Appendix B (Supplementary Material).

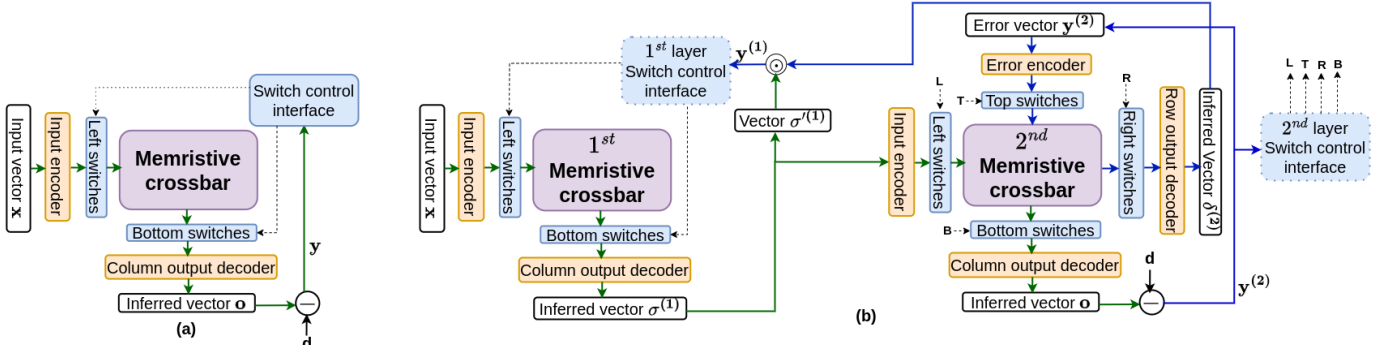


Figure 2: (a) Overview of the proposed memristive SLNN. (b) memristive MLNN with one hidden layer. The operator \odot denotes element-wise multiplication.

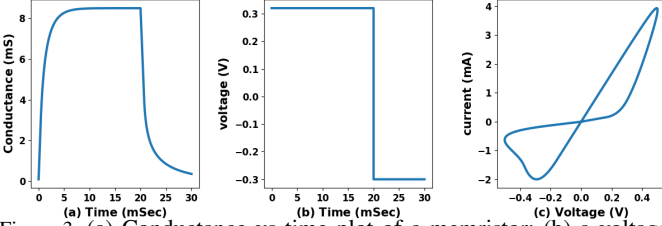


Figure 3: (a) Conductance vs time plot of a memristor; (b) a voltage pulse across a memristor with amplitudes greater than the thresholds; and (c) I-V characteristic curve of a memristor.

Table II: The specifics of memristive SLNN architecture

SLNN Term	Realisation on an $n \times m$ memristive crossbar
Input neuron n_i	Row i
Output neuron n_j	Column j
Synapse $s_{j,i}$	Memristor $M_{j,i}$ connecting input neuron n_i and output neuron n_j .
Weight $w_{j,i}$ of synapse $s_{j,i}$	$w_{j,i} = a \cdot R_0(G - G_{j,i})$, includes conductances G and $G_{j,i}$ of resistor R , and memristor $M_{j,i}$ respectively.

crossbar (MCB). In order to implement a memristive SLNN of n input and m output neurons, an $n \times m$ MCB, $(2n + 2m)$ resistors, and $(n + m + 2)$ op-amps along with a feedback resistor for each, are required. As a well-known practice, a resistor is realized with a pass transistor. The memristive SLNN architecture is shown in Fig. 4(a) and the specifics of its realization on an MCB are in Table II.

In Fig. 4(a), for all $i = 1$ to n , x_i is the i^{th} feature input to the memristive SLNN. The row input interface RII_i encodes x_i into voltage v_{x_i} (see Section III-B below) and feeds it to the i^{th} row via switch S_{li} , indicated on the left of the MCB. The inputs to op-amps OPM_1 and OPM_2 at the top are given via switches S_{l0} and \bar{S}_{l0} , and their outputs are v_f and v_b respectively. The synapse $s_{j,i}$ between input neuron i and output neuron j is represented by memristor $M_{j,i}$.

For all $j = 1$ to m , y_j s are the errors to be backpropagated through the crossbar. The column input interface CII_j encodes y_j into v_{y_j} which (only during backpropagation from T_{rd} to $2T_{rd}$, for T_{rd} period) is fed to the j^{th} column via switch S_{tj} at the top of the crossbar. The switch S_j is connected at the bottom of the j^{th} column of the MCB. The switch \bar{S}_{li} is placed on the right side of the i^{th} row of the MCB.

The row output interface (ROI_i) and column output interface (COI_j) collect the outputs δ_j and σ_i at the i^{th} row and j^{th} column during backward and forward inferences respectively.

During forward inference and update operations, the input voltages are given at the left side of the rows of MCB, whereas the v_{y_j} are fed at the top of columns for backward inference. It is to be noted that RIIs, ROIs, CIIs, and COIs form an input encoder, row output decoder, error encoder, and column output decoder, in Fig. 2, respectively.

B. Input encoder

The input feature and error at neurons are converted into proper voltage waveforms and fed to the memristive crossbar for inferences and updates of the conductances of the memristors in the crossbar during training. In Fig. 4(a), the functional peripheral circuits, namely RII_i and CII_j , encode the input x_i and error y_j into appropriate voltages respectively, as shown in Fig. 5. For forward inference from 0 to T_{rd} , RII_i encodes the input x_i as $v_{x_i} = ax_i$ where a is a positive constant, and $v_{th}^- < v_{x_i} < v_{th}^+$, as in Fig. 5.

Similarly, CII_j encodes error y_j at neuron n_j into v_{y_j} , $1 \leq j \leq m$, and $v_{th}^- < v_{y_j} < v_{th}^+$ to perform backward inference from T_{rd} to $2T_{rd}$. As v_{x_i} and v_{y_j} are within the thresholds, there is no change in weight during inferences.

During updates, the conductance of the memristor changes depending on the signs and magnitude of x_i and y_j as shown in Table III. The x_i part is taken care of by input update voltage. The input voltages for the update are encoded according to x_i and exceed the threshold voltages of memristors. For an update operation during the period $2T_{rd}$ to $2T_{rd} + T_{wr}$, the encoding has four possibilities depending on the signs of x and y , and the interval T_{wr} has four equal sub-periods with the input update voltage as follows:

- $x_i \geq 0$ (Fig. 5(a)): in the first and the second quarters the input update voltages are $v_{x_i} + v_{th}^+$ and $-v_{x_i} + v_{th}^-$, while in the third and the fourth, these are v_{th}^- and v_{th}^+ respectively;
- $x_i < 0$ (Fig. 5(b)): in the first and the second quarters the input update voltages are v_{th}^+ and v_{th}^- , while in the third and the fourth, these are $v_{x_i} + v_{th}^-$ and $-v_{x_i} + v_{th}^+$ respectively.

As $\Delta \mathbf{W}$ depends on input \mathbf{x} by Equation II.1, the voltages during an update operation are also proportional to \mathbf{x} .

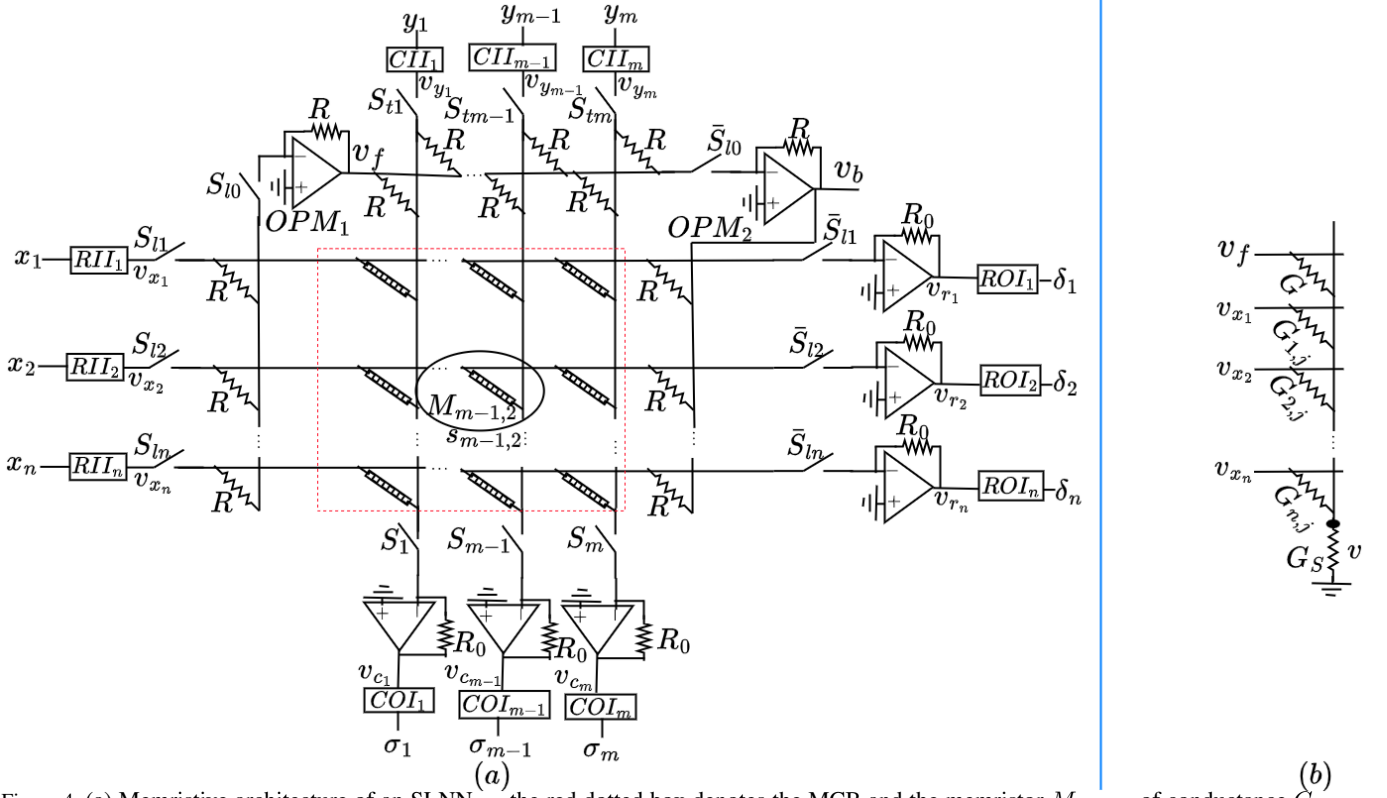


Figure 4: (a) Memristive architecture of an SLNN — the red dotted box denotes the MCB and the memristor $M_{m-1,2}$ of conductance $G_{m-1,2}$ is the synapse $s_{m-1,2}$ of weight $w_{m-1,2}$, (b) for column j denoting neuron n_j , $G_{j,i}$, G , and G_s respectively are the conductances of its n memristors, the resistor R , and the switch S_j ; v is the voltage across switch S_j for n_j .

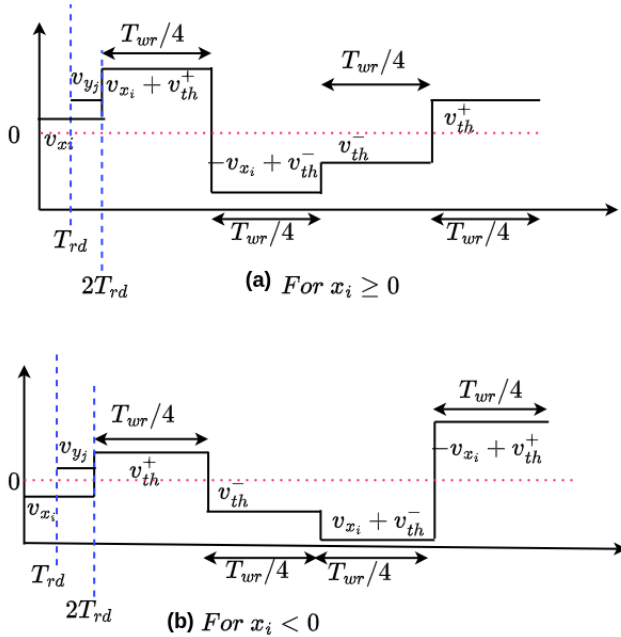


Figure 5: Encoding of input feature x_i and error y_j into voltage waveforms for forward and backward propagation, also for updating weights — (a) $x_i \geq 0$, (b) $x_i < 0$. For simulation, both rise time and fall times are considered as 1 ns.

C. Switch control signals for an MCB

In Fig. 4(a), the signals flow from the left of the rows to the bottom of the columns, and from the top of the columns to the right of the rows of an MCB during forward and backward inferences respectively. Additionally, the change in the conductance of its memristors is controlled by the bottom switches in Fig. 4(a). The directions of the signals and conductance update are guided by controlling the ON/OFF states of the voltage-controlled switches in memristive SLNN (Fig 4(a)). The functional block called the switch control interface shown in Fig. 2, produces voltages to control the ON/OFF timings of the switches. The switches S_{li} , S_j , S_{tj} and \bar{S}_{li} (for $0 \leq i \leq n$ and $1 \leq j \leq m$) are turned ON/OFF to perform the inferences and update operations as follows:

- forward inference from 0 to T_{rd} time units: *only* the switch S_{l0} , switches S_{li} ($1 \leq i \leq n$) on the left, and switches S_j ($1 \leq j \leq m$) at the bottom are ON;
- backward inference (while finding $\delta^{(k)} = W^{(k)T}y^{(k)}$ during backpropagation from T_{rd} to $2T_{rd}$ time units): *only* switch \bar{S}_{l0} , switches S_{tj} ($1 \leq j \leq m$) at the top, and \bar{S}_{li} ($1 \leq i \leq n$) on the right are ON;
- during update from $2T_{rd}$ to $2T_{rd} + T_{wr}$ time units: switch S_{l0} , switch \bar{S}_{l0} , switches S_{tj} ($1 \leq j \leq m$) at the top, and \bar{S}_{li} ($1 \leq i \leq n$) on the right are OFF. The switches S_{li} ($1 \leq i \leq n$) on the left are ON but switches S_j ($1 \leq j \leq m$) at the bottom may be ON or OFF depending on error y_j at column j during update (Fig. 6).

For synapse $s_{j,i}$, the $\Delta w_{j,i} = \eta y_j x_i$ depends on input x_i and

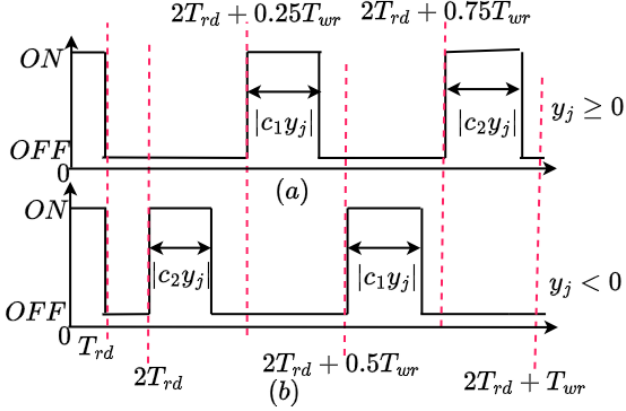


Figure 6: The states of switch S_j at the bottom of column j in the MCB of Fig. 4(a) — (a) for $y_j \geq 0$; (b) for $y_j < 0$. In the simulation, both rise and fall times are considered as 1 ns.

error y_j . The y_j part is taken care of by the states of switch S_j . The total update interval T_{wr} is divided into four equal sub-periods. The rules for turning switches S_j ON/OFF depending on the sign as well as the magnitude of the error y_j at neuron n_j are as follows (vide Fig. 6) with c_1 and c_2 being the respective slopes of linearly decreasing and increasing regions of the memristor's conductance (Fig. 3(a)):

- $y_j \geq 0$: S_j ON for periods proportional to $|c_2 y_j|$ and $|c_1 y_j|$ in the 2nd and 4th quarters of T_{wr} respectively (Fig. 6(a));
- $y_j < 0$: switch S_j ON for periods proportional to $|c_1 y_j|$ and $|c_2 y_j|$ in the 1st and 3rd quarters of T_{wr} respectively (Fig. 6(b)).

D. Memristive multi-layer neural network

A memristive MLNN of N hidden layers is implemented by cascading $N + 1$ memristive crossbars (MCBs) as shown in Fig. 7, where the output of MCB in layer l is fed to the subsequent MCB for layer $l + 1$. In the memristive MLNN, the input encoder and the error encoder encode input vector \mathbf{x} and error vector \mathbf{y} into \mathbf{v}_x and \mathbf{v}_y respectively. The input \mathbf{v}_x propagates layer by layer from MCB_1 to MCB_{N+1} while \mathbf{v}_y propagates from MCB_{N+1} to MCB_1 . The desired output vector \mathbf{d} is given at the output layer to get the error vector $\mathbf{y}^{(\text{output})}$ only during the training. The $\delta^{(k)}$ at layer k is rescaled (optional) using \tanh function. The arrows in Fig. 7 indicate the flow of data.

IV. IN-SITU TRAINING ON A MEMRISTIVE CROSSBAR

A. Inference and update operations on an MCB

In this work, memristors are operated in the linearly increasing and decreasing portions of their conductance vs time plot (Fig. 3(a)) to maintain all the in-memory activities in the linear domain. Let G_{max} and G_{min} denote the maximum and minimum conductance in the linear region. For neuron n_j , the conductance of memristor $M_{j,i}$ and switch S_j denoted as $G_{j,i}$ and G_s respectively typically have the following properties:

- 1) $G_s \gg G_{j,i}$ when switch S_j is ON;

- 2) $G_s \ll G_{j,i}$ when switch S_j is OFF.

With these premises, the inference and update on an MCB are illustrated below.

1) *Inference*: During forward inference to get the vector σ , the states of the switches in Fig. 4(a) in anti-clockwise order are $S_{li} : \text{ON}$, $S_j : \text{ON}$, $\bar{S}_{li} : \text{OFF}$ and $S_{tj} : \text{OFF}$; for $0 \leq i \leq n$ and $1 \leq j \leq m$ from 0 to T_{rd} time units. Fig. 4(b) shows the conductance equivalent circuit of the neuron n_j , in which the memristor $M_{j,i}$, resistor R of OPM_1 , and switch S_j of column j have been replaced by the corresponding conductance $G_{j,i}$, G , and G_s respectively, where $1 \leq i \leq n$. It is noted that resistors and switches can be realized with transistors. We compute the voltages across memristors of the neuron n_j during inference and update to observe their effects on conductance $G_{j,i}$ of $M_{j,i}$. Recall that $R II_i$ encodes x_i as $v_{x_i} = ax_i$ (refer Fig. 5) such that $v_{th}^- < v_{x_i} < v_{th}^+$ for inference. It is necessary to know the voltages v across switch S_j and $(v_{x_i} - v)$ across memristor $M_{j,i}$ (refer Fig. 4(b)) to perform the inference and update on memristive MCB. With $G = \frac{1}{R} = \frac{G_{max} + G_{min}}{2}$, we have

$$(v_f - v)G + \sum (v_{x_i} - v)G_{j,i} = vG_s$$

$$v = \frac{v_f G + \sum_{i=1}^n v_{x_i} G_{j,i}}{G_s + G + \sum_{i=1}^n G_{j,i}} \quad (\text{IV.1})$$

As S_j is ON, $G_s \gg G$ and $G_s \gg G_{j,i}$, for $i = 1$ to n . Thus, if $G_s \gg \sum_{i=1}^n G_{j,i}$, from Equation IV.1 we get

$$v \approx \frac{v_f G + \sum_{i=1}^n v_{x_i} G_{j,i}}{G_s} = v_f \frac{G}{G_s} + \sum_{i=1}^n v_{x_i} \frac{G_{j,i}}{G_s} \quad (\text{IV.2})$$

Therefore, $M_{j,i}$ experiences $v_{x_i} - v \approx v_{x_i}$, the input voltage to i^{th} row which lies in the threshold range of $M_{j,i}$. This implies no change in $G_{j,i}$ of $M_{j,i}$, and in the weight $w_{j,i}$ ($= aR_0(G - G_{j,i})$) of synapse $s_{j,i}$. Again, referring to Fig. 4(a), the voltage $v_f = -\sum_{\forall i} \frac{R}{R} v_{x_i} = -\sum_{\forall i} v_{x_i}$ and the output v_{c_j} of op-amp at bottom of j^{th} column is

$$v_{c_j} = - \left[\sum_{i=1}^n (R_0 G_{j,i} \times v_{x_i}) + R_0 G \times v_f \right]$$

$$v_{c_j} = - \sum_{i=1}^n R_0 (G_{j,i} - G) v_{x_i} = \sum_{i=1}^n a R_0 (G - G_{j,i}) \frac{v_{x_i}}{a}$$

Defining the weights $w_{j,i} = aR_0(G - G_{j,i})$, we have $v_{c_j} = \sum_{i=1}^n w_{j,i} x_i$ where $x_i = \frac{v_{x_i}}{a}$. The weight $w_{j,i}$ is positive if $G > G_{j,i}$ else non-positive. It is to be noted that for $j = 1$ to m , all v_{c_j} together form a pre-activated vector $\mathbf{r} = \mathbf{W}\mathbf{x}$, i.e., the weight matrix \mathbf{W} with $w_{j,i}$ as its entries multiplied by the input vector \mathbf{x} . Therefore, the matrix-vector multiplication is performed within the period T_{rd} and hence $\mathcal{O}(1)$ time. The COI_j produces activated output σ_j over r_j of n_j .

During backpropagation, to get vector δ consisting of δ_i , for $i = 1$ to n , needs the states of the switches (in Fig. 4(a)) to be $S_{li} : \text{OFF}$, $S_j : \text{OFF}$, $\bar{S}_{li} : \text{ON}$ and $S_{tj} : \text{ON}$; $0 \leq i \leq n$ and $1 \leq j \leq m$ from T_{rd} to $2T_{rd}$ time units.

This is simply an inference operation where the inputs v_{y_j} for $j : 1 \leq j \leq m$ are fed at the top of the MCB, and the vector δ are inferred by ROIs interfaces. This operation is also performed within T_{rd} time period, i.e., $\mathcal{O}(1)$.

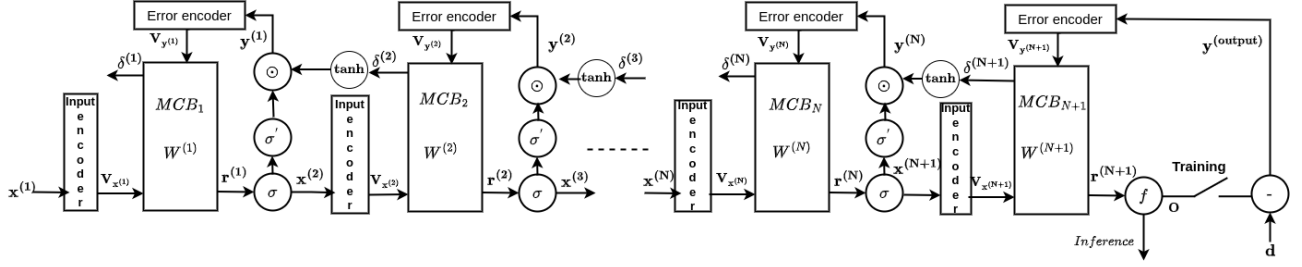


Figure 7: A schematic of the architecture and training process for the memristive MLNN. The memristive crossbar MCB_k with synaptic weight matrix $W^{(k)}$ connects layers k and $k+1$. The f is the output function at the output layer. All RIs together form an Input encoder to encode the input vector $\mathbf{x}^{(k)}$ into a proper voltage form using the encoding scheme presented in Fig. 5. The neural activation function is σ , and its derivative on the input is σ' . Element wise product of σ' of $(k-1)^{th}$ layer and $\delta^{(k)}$ is denoted by \odot , where $\delta^{(k)} = W^{(k)T} \mathbf{y}^{(k)}$. The $\delta^{(1)}$ is neglected. $\mathbf{y}^{(k)}$ is the error vector at layer k .

2) *Update*: By Equation II.1, the change in weight of synapse $s_{j,i}$ is given by $\Delta w_{j,i} = \eta y_j x_i$, which depends on the signs and values of the input x_i and the error y_j in four possible ways, as shown in Table III.

Table III: The rules to update the weight $w_{j,i} = a \cdot R_0(G - G_{j,i})$ — change in conductance $G_{j,i}$ of memristor $M_{j,i}$ with changes in input x_i and error y_j .

Input x_i	Error y_j	Weight change $\Delta w_{j,i}$	Conductance change $G_{j,i}$	Updating period decided by v_{x_i} (Fig. 5) and S_j (Fig. 6)
+ve	+ve	+ve	↓	$2T_{rd} + 0.25T_{wr}$ to $2T_{rd} + 0.5T_{wr}$
+ve	-ve	-ve	↑	$2T_{rd}$ to $2T_{rd} + 0.25T_{wr}$
-ve	+ve	-ve	↑	$2T_{rd} + 0.75T_{wr}$ to $2T_{rd} + T_{wr}$
-ve	-ve	+ve	↓	$2T_{rd} + 0.5T_{wr}$ to $2T_{rd} + 0.75T_{wr}$

The weight $w_{j,i} (= a \cdot R_0(G - G_{j,i}))$ of synapse $s_{j,i}$ increases or decreases if the value of $G_{j,i}$ of $M_{j,i}$ decreases or increases respectively. During the update from $2T_{rd}$ to $2T_{rd} + T_{wr}$ time units, the switch states are:

S_{l0} and $\bar{S}_{l0} : OFF$, $S_{li} : ON$, $\bar{S}_{li} : OFF$ and $S_{tj} : OFF$, $1 \leq i \leq n$ and $1 \leq j \leq m$.

Therefore, $v_f = 0$. At n_j , assuming the input feature $x_i \geq 0$, then input update voltage $v_{x_i}^*$ at row_i has to be $v_{th}^+ < v_{x_i}^* < 2v_{th}^+$ and $2v_{th}^- < v_{x_i}^* < v_{th}^-$ in the first and second quarters of T_{wr} as per input encoding in Fig. 5. Referring to Fig. 6 and Section III-C, if error $y_j \geq 0$ at n_j , then switch S_j is *ON* up to the period proportional to $|y_j|$ in the second and the fourth quarters of T_{wr} . Under this condition, upto a time period of $|c_1 y_j|$ in the second quarter of T_{wr} , the memristor $M_{j,i}$ experiences a negative voltage $v_{x_i}^* (= -v_{x_i} + v_{th}^-)$ and $G_{j,i}$ decreases ($\Delta G_{j,i} \propto -x_i y_j$) that results in the increase of $w_{j,i} (\propto x_i y_j)$. In the fourth quarter of T_{wr} (refer Fig. 6 and Fig. 5), $M_{j,i}$ experiences a positive voltage but within the threshold range, so no change in $G_{j,i}$ occurs.

However, the salient question is that if S_j of neuron n_j is *OFF* during T_{wr} , will $G_{j,i}$ change? The answer is NO. The reason is as follows.

If S_j is *OFF*, then $G_s \ll G$ and $G_s \ll G_{j,i}$ for $i = 1, 2, \dots, n$ and $v_f = 0$, then from Equation IV.1

$$v = \frac{\sum_{i=1}^n v_{x_i}^* G_{j,i}}{G + \sum_{i=1}^n G_{j,i}}$$

Suppose we get $v = v^*$. The input update voltage $v_{x_i}^*$ at i^{th} row is either positive or negative but not both in any quarter of T_{wr} as shown in Fig. 5. So the v^* will be either $2v_{th}^- < v^* < v_{th}^-$, or $v_{th}^+ < v^* < 2v_{th}^+$ in any quarter of T_{wr} . Therefore, each $M_{j,i}$ of column j experiences a voltage across it in the range $v_{th}^- < v_{x_i}^* - v = v_{x_i}^* - v^* < v_{th}^+$. Therefore, no change in

$G_{j,i}$ results in no update in $w_{j,i}$. Similarly, for the other three combinations of signs of the input x_i and error y_j as given in Table III, the $G_{j,i}$ of $M_{j,i}$ is updated accordingly resulting $w_{j,i}$ is modified correctly. Therefore, there is no conductance alteration in column j if S_j is *OFF*.

In the linear region of the conductance $G_{j,i}$ of $M_{j,i}$ (Fig. 3(a)), $\Delta G_{j,i}$ is proportional to the product of the update voltages and the time duration of this voltage. During the update, the change $\Delta G_{j,i}$ is jointly taken care of by input voltage $v_{x_i}^*$ and the time duration ($\propto |y_j|$) for which the bottom switches S_j are *ON*. Since $\Delta G_{j,i}$ is proportional to the product of input x_i and time $|y_j|$ ($\propto |y_j| \cdot x_i$), the weight $w_{j,i}$ is updated as per Equation II.1.

Since each $M_{j,i}$ is updating its $G_{j,i}$ independently for an input x_i and error y_j , therefore all $G_{j,i}$ s of a crossbar are updated in $O(1)$ time. Therefore, the weight updates that involve the outer product of the vectors and addition operations (Equation II.1), are carried out in constant time with an MCB. With these inferences and update operations on an MCB, the *in-situ* training of the memristive MLNN is presented next.

B. Training of a memristive MLNN

The proposed *in-situ* supervised algorithms for the training of a memristive MLNN are based on gradient descent backpropagation. Fig. 7 is a flow diagram of the training. The data set is divided into training and test sets.

1) *Inference on memristive MLNN*: This operation is performed during both training and testing. The inference vector \mathbf{o} at the output layer from the input \mathbf{x} is obtained with Algorithm 1. The accuracy is calculated on the test set.

2) *Update on memristive MLNN*: During training, the weight matrix in the MCB is updated. The error $\mathbf{y}^{(\text{output})}$ at the output layer is the difference between the forward inferred vector \mathbf{o} and the labeled output vector \mathbf{d} . The $\mathbf{y}^{(\text{output})}$ is back propagated and the synaptic weights are adjusted layer by layer as described in Algorithm 2. During backpropagation, the crossbar MCB_i is used to get the vectors $\delta^{(i)}$ at the i^{th} layer in $O(1)$ time as described in Algorithm 2 (line 5), which accelerates the backpropagation.

From Algorithms 1 and 2, it is to be noted that for a layer, the three operations of matrix-vector multiplication, $\delta^{(2)}$ operation, and MCB update are performed within periods of T_{rd} , T_{rd} , and T_{wr} respectively, all in $O(1)$ time.

Algorithm 1: Inference operation on a memristive MLNN with N hidden layers

Data: Input vector \mathbf{x} from training data set;

Result: Anticipated output vector \mathbf{o} ;

```

1  $\mathbf{x}^{(1)} = \mathbf{x}$ ;
2 for  $k = 1$  to  $N + 1$  do
3   Apply Input ENCODER comprising row input interfaces
   (RIIs), to encode  $\mathbf{x}^{(k)}$  into equivalent voltages  $\mathbf{v}_{\mathbf{x}^{(k)}}$ ;
4   Apply  $\mathbf{v}_{\mathbf{x}^{(k)}}$  to  $MCB_k$ ;
5   for  $p = 0$  to  $n$ , keep all switches  $S_{lp}$  (left) ON and  $\bar{S}_{lp}$ 
   (right) OFF in  $MCB_k$  for inference duration of  $T_{rd}$ ;
6   for  $q = 1$  to  $m$ , keep all switches  $S_q$  (bottom) ON and
    $\bar{S}_{tq}$  (top) OFF in  $MCB_k$  for inference duration of
    $T_{rd}$ ;
7   Assign voltage  $\mathbf{v}_c^{(k)}$  of  $k^{th}$   $MCB_k$  to  $\mathbf{r}^{(k)}$ ;
8   if  $k < N + 1$  then
9      $\mathbf{x}^{(k+1)} = \sigma(\mathbf{r}^{(k)})$ , the input for  $MCB_{k+1}$  by
     applying the activation function to  $\mathbf{r}^{(k)}$ ;
10  else
11     $\mathbf{o} = \mathbf{f}(\mathbf{r}^{(k)})$ , the predicted values at the output layer.

```

Algorithm 2: Update operation for the weight matrices of a memristive MLNN with N hidden layers

Data: Input vector \mathbf{x} , anticipated output vector \mathbf{o} of Algorithm 1, labeled output vector \mathbf{d} for \mathbf{x} ;

Result: Proper weight adjustment of each synapse in the memristive MLNN;

```

1  $\mathbf{y}^{(N+1)} = \mathbf{y}^{(output)}$ , the error at the output layer with  $\mathbf{o}$ 
  obtained from Algorithm. 1 and  $\mathbf{d}$ ;
2 for  $k = N + 1$  to 1 do
3   Get error voltages  $\mathbf{v}_{\mathbf{y}^{(k)}}$  from error ENCODER by
   feeding  $\mathbf{y}^{(k)}$  to it;
4   Apply  $\mathbf{v}_{\mathbf{y}^{(k)}}$  to  $MCB_k$ ;
5   for  $p = 0$  to  $n$ ,  $q = 0$  to  $m$  do
6     keep switches  $\bar{S}_{lp}$  and  $S_{tq}$  ON, and  $S_{lp}$  and  $\bar{S}_q$  OFF
     for the time period  $T_{rd}$  to get
      $\delta^{(k)} = (\mathbf{W}^{(k)})^T \mathbf{y}^{(k)}$ ;
7   Update the crossbar  $MCB_k$  using update rules in
   Table III (Section IV-A2);
8   if  $k > 1$  then
9      $\mathbf{y}^{(k-1)} = \tanh(\delta^{(k)}) \odot \sigma'(\mathbf{r}^{(k-1)})$ .

```

V. EXPERIMENTAL RESULTS

The proposed memristive SLNNs and MLNN were simulated using the open source LTSpice XVII simulator, running on an Ubuntu 20.04 LTS environment with an 8-core 1.6GHz Intel Core i5 processor and 8GB RAM. The datasets for classification, network type with the number of features in each layer, number of crossbars and their sizes, and resistor R_0 are given in Table IV. The values of T_{rd} and T_{wr} are chosen as $10\mu s$ and $1ms$ respectively. For all simulations carried out with memristors, memristive crossbars, resistors, Opamps, switches, and functional blocks such as RII_i , ROI_j , etc., the wire resistance has been assumed to be negligible. Further, we simulated a memristive MLNN with MNIST dataset [30] in Python to demonstrate the scalability of the architecture.

For a dataset, if n and m are the number of features and classes, then $n + 1$ (along with a bias) and m are the input

Table IV: Neural network parameters used in our LTSpice simulation

Data set	Network type, # features	#MCBs with size	Memristor model [29]	Memristor model [31]
			R_0 in Ω	R_0 in Ω
NASA Asteroid [32]	Single layer, $20 \rightarrow 1$	One; 21×1	100	2k
Breast Cancer Wisconsin [33]	Single layer, $30 \rightarrow 1$	One; 31×1	100	15k
XOR	Multi-layer, $2 \rightarrow 2 \rightarrow 2$	Two; 3×2 and 3×2	1k	—
IRIS [33]	Multi-layer, $4 \rightarrow 4 \rightarrow 3$	Two; 5×4 and 5×3	1k	15k

and output neurons at the input and output layers respectively. The number of hidden layers and hidden neurons are decided empirically by experiments. A sigmoid function was used at the output layer in SLNNs. The sigmoid and \tanh were used as activation functions at the hidden layers in memristive MLNNs for IRIS and XOR datasets respectively. The output function chosen in MLNNs was *soft-max* while cross entropy was the cost function. The input data \mathbf{x} belongs to that neuron/class which has the highest *soft-max* value.

The Spice models of silver chalcogenide [29] and anodic titania [31] based memristors have been employed. These memristor spice models are in [28]. The values of the parameters² in Table V used in general Spice model [28] are to match the characterizations of silver chalcogenide [29] and anodic titania [31] based memristors. From the simulation, the average energy consumed at a synapse to perform both read and write operations is $1.589\mu J$.

An illustration of the steps: A memristive MLNN on the XOR dataset with network structure mentioned in Table IV was simulated to explain the training steps (Fig. 7) and the results are presented in Fig. 8. The four inputs are $[0, 0]$, $[0, 1]$, $[1, 0]$, and $[1, 1]$ whereas the respective output vectors are $[1, 0]$, $[0, 1]$, $[0, 1]$, and $[1, 0]$. As an example, for input $[0, 1]$ with constant $a = 0.6$, the encoded input voltages $[v_{x_1}^{(1)}, v_{x_2}^{(1)}]$ during forward inference and update are shown Fig. 8(a) and the corresponding desired output vector $[d_1, d_2]$ is in Fig. 8(b). The weights (memristors' conductance) are initialized randomly. The weighted sums $[r_1^{(1)}, r_2^{(1)}]$ (equivalent to v_{c_1} and v_{c_2} in MCB in Fig. 4) during forward inference of the hidden layer is in Fig. 8(c). The output $[\sigma_1^{(1)}, \sigma_2^{(1)}]$ in the hidden layer is in Fig. 8(d). The $[\sigma_1^{(1)}, \sigma_2^{(1)}] \equiv [\mathbf{x}_1^{(2)}, \mathbf{x}_2^{(2)}]$ is encoded for the next layer, as presented in Fig. 8(e). Fig. 8(f) shows the weighted sums $[r_1^{(2)}, r_2^{(2)}]$ of output layer. The inferred output vector $[o_1, o_1]$ is in Fig. 8(g). The error $[y_1^{(2)}, y_2^{(2)}]$ at output layer is shown in Fig. 8(h). The vector $[\delta_1^{(2)}, \delta_2^{(2)}]$ is calculated with MCB_2 during backward inference and rescaled with tanh function as $[\tanh(\delta_1^{(2)}), \tanh(\delta_2^{(2)})]$ and shown in Fig. 8(i). Fig. 8(j) shows the error at the first layer.

Table V: Parameters used in Spice models for memristors in [28]

Memristor Parameter \rightarrow	a_1	a_2	b	a_p	a_n	X_p	X_n	V_p	V_n	α_p	α_n	η
Memristor model [29]	0.17	0.17	0.05	4000	4000	0.3	0.5	0.16	0.15	1	5	1
Memristor model [31]	1.4	1.4	0.05	16	11	0.3	0.5	0.65	0.56	1.1	6.2	-1

Simulation with memristor model in [29] : In the spice model [29] for silver chalcogenide memristors, the threshold values v_{th}^+ and v_{th}^- are $0.16V$ and $-0.15V$ respectively. The minimum and maximum conductance are $0.255mS$ and $8.5mS$ respectively while G_{min} and G_{max} in the linear region of the memristor's conductance are $3.18mS$ and $6.38mS$ respectively. The conductance $G_{j,i}$ of memristor $M_{j,i}$ was initialized in the range $\{4.4mS, 5mS\}$. The memristive SLNNs

²Glossary in Table XI of Appendix A

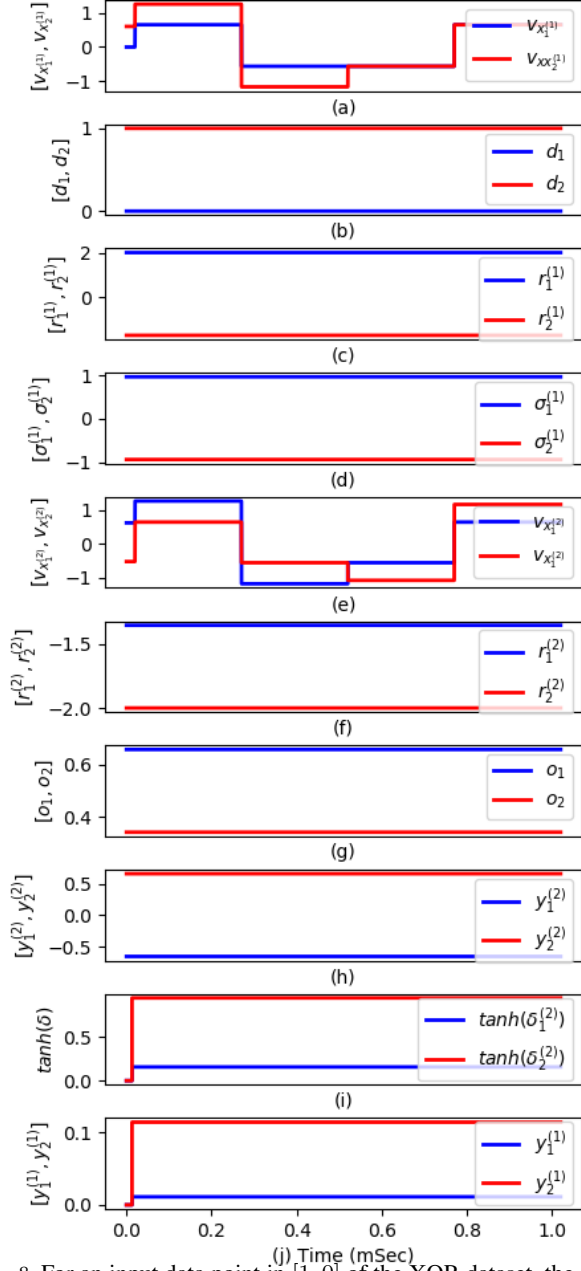


Figure 8: For an input data point in $[1, 0]$ of the XOR dataset, the plots (a) to (j) represent the training steps of our memristive MLNN. In (d), the output $[\sigma_1^{(1)}, \sigma_2^{(1)}]$ of the hidden layer is the input $[\mathbf{x}_1^{(2)}, \mathbf{x}_2^{(2)}]$ to the output layer. In (i), $\tanh(\delta) = [\tanh(\delta_1^{(2)}), \tanh(\delta_2^{(2)})]$.

and MLNNs were trained with algorithms 1 and 2 on NASA Asteroid, Breast Cancer Wisconsin, and IRIS data sets with 90.43%, 98.59%, and 98.22% accuracies respectively. The corresponding cost functions during training are shown in Figures 9(a), (b), and (c) respectively for these three data sets. For each dataset, the number of training epochs as indicated has been chosen empirically to avoid overfitting.

Simulation with memristor model in [31]: Our proposed memristive ANN architecture was also tested with anodic titania-based spice model [31] of a memristor. Its threshold voltages v_{th}^+ and v_{th}^- are 0.65V and -0.56V respectively. The minimum and maximum conductances are 1mS and 70mS, respectively. The linear section of its conductance plot

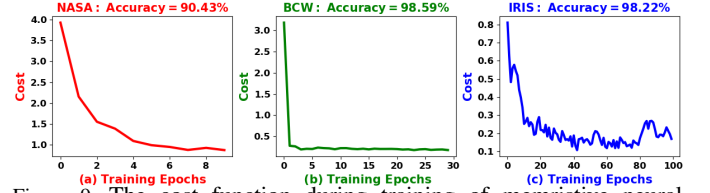


Figure 9: The cost function during training of memristive neural networks constructed with silver chalcogenide based memristors [29] for (a) NASA Asteroid (b) Breast Cancer Wisconsin (c) IRIS datasets.

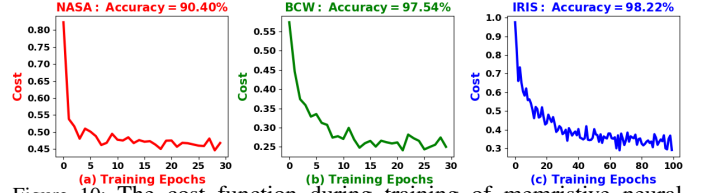


Figure 10: The cost function during training of memristive neural networks constructed with anodic titania-based memristors [31] for (a) NASA Asteroid (b) Breast Cancer Wisconsin (c) IRIS datasets.

has 28mS and 48mS as the values of G_{min} and G_{max} , respectively. The initialization range for the conductance of a memristor was between 35mS and 41mS. The memristive SLNNs and MLNNs were trained using proposed in-situ algorithms on NASA Asteroid, Breast Cancer Wisconsin, and IRIS data sets with 90.40%, 97.54%, and 98.22% accuracies respectively. Figures 10(a), (b), and (c) are the costs during training.

Performance of proposed architecture: The performance is compared with similar works in [1], [7], [34], [23] and the classification accuracy with both memristor models [29] and [31] are presented in Table VI, where the one in [29] gives a slightly better value.

Table VI: Comparison of classification accuracy (%) with prior similar works [1], [7], [34], [23].

Data set	Network structure	[1]	[7]	[34]	[23]	Proposed work with memristor	
						Model [29]	Model [31]
NASA Asteroid	Single-layer	89.04	90.43	90.40
Breast Cancer Wisconsin	Single-layer	90.14	98.7	97	97.72	98.59	97.54
IRIS	Multi-layer with one hidden layer	99.11	97.33	84.33	88.43	98.22	98.22

A. Variability analysis for the proposed architecture

In order to analyze the memristor's device variation issues, we added variations in memristors' spice models [29], [31] and evaluated the performance of memristive MLNN on IRIS dataset. In order to add $\nu\%$ variations in the I-V characteristics of the memristors, the fitting parameters were adjusted until the change in I-V characteristic exceeded $\nu\%$ [28]. The 10% change in I-V characteristic is measured using the total average difference between the altered I-V characteristic with parameters in Table VII and the initial I-V characteristic with parameters in Table V [28].

In order to analyze the effects of variation on the performance of our memristive MLNN, we have chosen randomly 10%, 20%, and 30% of the memristors and added 10% I-V curve variation in these selected memristors. Simulation results in Table VIII demonstrate that the variations in memristors have not affected the performance significantly.

Table VII: Memristor parameters increased or decreased from original values in Table V in order to obtain 10% variation of I-V characteristic in memristor models of [29], [31]; these values are taken from [28].

Memristor Parameter \rightarrow		a_1	a_2	b	a_p	a_n	X_p	X_n	V_p	V_n	α_p	α_n
Decreased values	Model [29]	0.153	0.153	0.045	2680	2680	0.18462	0.3077	0.104848	0.098295	0.145	0.725
	Model [31]	1.26	1.26	0.045	9.888	6.798	0.2139	0.3565	0.594165	0.511896	0	0
Increased values	Model [29]	0.187	0.187	0.055	5924	5924	0.42363	0.70605	0.217696	0.20409	2.115	10.575
	Model [31]	1.54	1.54	0.055	32.16	22.11	0.57903	0.96505	0.6994	0.60256	1.9855	11.191

Table VIII: Accuracy and F1-Score with 10% I-V curve variation added to randomly chosen memristors in the MLNN on IRIS dataset.

% of memristors having variation	Variation added by decreasing memristor parameters				Variation added by increasing memristor parameters			
	Model [29]		Model [31]		Model [29]		Model [31]	
	Accuracy	F1-Score	Accuracy	F1-Score	Accuracy	F1-Score	Accuracy	F1-Score
10	97.33%	96.214%	96.44%	94.667%	98.22%	97.598%	97.33%	96.02%
20	97.33%	96%	96.44%	94.667%	98.22%	97.598%	97.33%	96.209%
30	97.33%	96.214%	96.44%	94.667%	98.22%	97.598%	96.44%	95.03%

B. Robustness of the architecture against stuck-at-a-conductance state

The robustness against faulty memristors is crucial for maintaining performance. To address this, several simulation experiments were carried out when a certain percentage (tested for 1%, 5%, 10%, and 20% of memristors randomly chosen) of the memristors in the crossbars were stuck-at-a-conducting state and unable to update themselves during training. These experiments were done with the spice model [29] to assess the effect in the above circumstance. After training, the accuracy has not been affected much as evident from Table IX and compared with a similar study in [1]. This demonstrates that compared to [1], the architecture and *in-situ* training are more robust to the presence of faulty memristors. The authors [9] performed fault analysis on a memristive MLNN with one hidden layer (2M type synapse and $O(n)$ MCB update time) on the MNIST dataset (8×8 image size, i.e., 64 input neurons) where the accuracy was 91.7% with 11% memristors being stuck.

C. Sneak path immunity

The sneak-path problem is a common challenge in 1M memristive crossbar arrays. Various methods have been proposed to either mitigate sneak-path effects during functional operation [35] or to exploit the sneak-path current for testing purposes [36]. The proposed model exhibits inherent robustness against this issue, as the sneak-path voltages across all unintended memristors remain below v_{th} , resulting in no undesired updates under all conditions.

To demonstrate this behavior, a 2×2 memristive crossbar array is considered, as illustrated in Fig. 11. During the inference phase Fig. 11(a), each memristor is subjected to either v_1 or v_2 , both of which are below the threshold voltage v_{th} . As a result, no change in memristor conductance occurs, effectively preventing sneak path-induced disturbances. During the update phase Fig. 11(b), assume that only neuron n_1 (corresponding to the first column) undergoes synaptic modification, while neuron n_2 does not. During update operation, voltages v'_1 and v'_2 are applied across M_{11} and M_{21} , respectively, where $v'_1 = v_1 + v_{th}$ and $v'_2 = v_2 + v_{th}$. An undesired sneak path voltage $|v'_1 - v'_2| = |v_1 - v_2| \leq v_{th}$ is appeared across M_{12} and M_{21} , respectively, which remains below the threshold. This

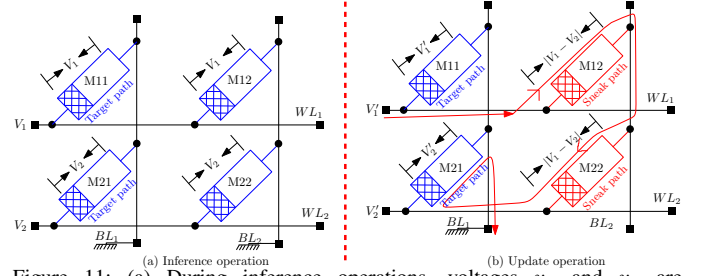


Figure 11: (a) During inference operations, voltages v_1 and v_2 are applied across M_{11}, M_{12} and M_{21}, M_{22} , respectively, where $v_1 \geq v_2$ and $v_1, v_2 \leq v_{th}$. (b) During update operation, voltages v'_1 and v'_2 are applied across M_{11} and M_{21} , respectively, where $v'_1 = v_1 + v_{th}$ and $v'_2 = v_2 + v_{th}$. An undesired sneak path voltage $|v'_1 - v'_2| = |v_1 - v_2| \leq v_{th}$ is appeared across M_{12} and M_{21} , respectively, resulting no effect.

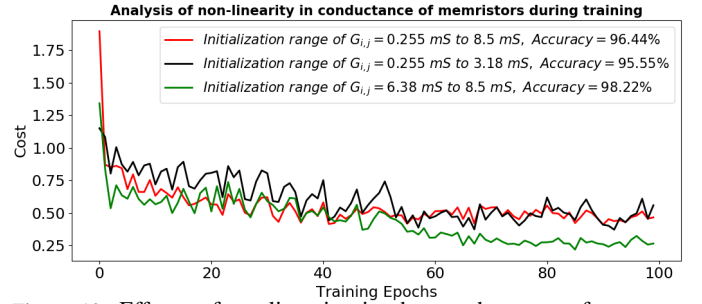


Figure 12: Effects of nonlinearity in the conductance of a memristor (Fig. 3(a)) on the performance for IRIS dataset; the conductance of the memristors [29] of the crossbar are initialized in the range (i) $\{0.225mS, 8.5mS\}$ (red), (iii) only lower non-linear $\{0.225mS, 3.18mS\}$ (black), and (iv) only upper non-linear $\{6.38mS, 8.5mS\}$ (green). For range (i) only linear conductance region, the accuracy shown in Fig. 9(c) is 98.22%.

ensures that no unintended updates occur, thereby effectively mitigating the sneak path issue.

D. Effects of non-linearity on the performance

The memristor's conductance changes non-linearly. In order to examine the effect of the non-linearity in conductance (refer Fig. 3(a)) on performance, four experiments were performed where memristors were initialized (i) only in the linear region, (ii) in both linear and non-linear region, (iii) only in the lower non-linear region, and (iv) only in the upper non-linear region. The analysis was performed with memristor model [29] that was trained and tested with the IRIS data set.

Table IX: Classification accuracy of our architecture after training with faulty memristors and comparison with [1].

	Memristive neural network architecture in [1]					Proposed memristive neural network architecture				
% of memristors randomly chosen as stuck-at a conductance state	0 %	1 %	5 %	10 %	20 %	0 %	1 %	5 %	10 %	20 %
Classification accuracy (%) for NASA Asteroid data set	89.04	85.68	79.97	82.89	72.83	90.43	90.43	91.86	92.24	89.39
Classification accuracy (%) for Breast Cancer Wisconsin data set	90.14	83.8	82.04	81.34	79.23	98.59	98.59	98.59	99.65	98.24
Classification accuracy (%) for IRIS data set	99.11	99.11	96.44	94.67	80.44	98.22	98.22	97.33	97.33	98.22

For case (i), we get 98.22% accuracy as shown in Fig. 9(c). The red plot in Fig. 12 shows the loss during training with an accuracy of 96.44% when the memristors are initialized in the range of $\{0.225mS, 8.5mS\}$ (linear and non-linear regions of conductance in Fig. 3(a)). The black plot in Fig. 12 shows the loss during training with accuracy of 95.55% when memristors were initialized in the lower non-linear range $\{0.225mS, 3.18mS\}$. When the memristors were randomly initialized only in the upper non-linear range $\{6.38mS, 8.5mS\}$, it gave 98.22% accuracy and the cost is shown in the green plot in Fig. 12. In summary, the results of these experiments establish that the training in the non-linear region on this architecture with its *in-situ* training algorithm does not affect the accuracy drastically compared to that with the linear region only.

E. Scalability of the proposed architecture

We performed digit classification of MNIST [30] data set to test the scalability of the proposed architecture. The training and testing data sets had 60000 and 10000 images of shape 28×28 pixels respectively and each pixel is considered as a feature. The LTspice simulator becomes very slow for large MCBs. Hence the spice netlists of memristive MLNN for the MNIST dataset are not suitable. There is also a limit

of 1024 on the number of nodes in a subcircuit that this simulator can handle. Because of the above limitations, the proposed algorithms for memristive MLNN with two hidden layers ($784 \rightarrow 397 \rightarrow 204 \rightarrow 10$ and total 394887 synapses) were simulated in Python. The cost during training is shown in Fig. 13. The proposed model achieves 91.27% accuracy on MNIST using Python with just 2 hidden layers and 7 epochs. Table X summarizes previous works with varying accuracy with different synapse configurations, update times, and simulation environments.

Discussion: Compared to [1], this architecture uses 50% fewer memristors, employing a 1M MCB without transistors or synapse-controlling devices. It addresses resilience to faults from stuck-at-conductance states, not covered in [1], and, to the best of our knowledge, is the first to explore the impact of non-linearity in memristor conductance. We study performance variations due to memristor imperfections in detail, unlike [1]. While [1] only considered conductance increases, our model supports both increases and decreases based on input and error, utilizing separate encoding and distinct control signals for positive and negative inputs and errors.

This work focuses on memristive MLNN and its related aspects. Compared with MLNNs, the data flow in convolutional neural networks (CNNs) is different. Moreover, major operations like stride, pooling during forward propagation, and full convolution operation during backpropagation are performed in CNNs. The future direction will be to modify the proposed memristive architecture for *in-situ* trainable CNNs that will be efficient in area, energy, and training latency.

VI. CONCLUSION

A novel architecture of memristive multi-layer neural networks with an efficient *in-situ* training algorithm is proposed here. The training algorithm based on gradient descent back propagation updates all the memristors of a crossbar in $\mathcal{O}(1)$ time. Here, one memristor suffices for a single synapse. Furthermore, it is demonstrated that the accuracy of the MLNN is not significantly impacted, even if some of the memristors (tested for 1%, 5%, 10%, and 20% of the memristors) are stuck at a conducting state. Experimental analysis revealed that the variation and non-linearity in the conductance of memristors do not have a notable impact on the accuracy of the MLNN. At a synapse, the average energy required to execute read and write operations is $1.589\mu J$. The comprehensive energy analysis is being investigated and will be provided separately.

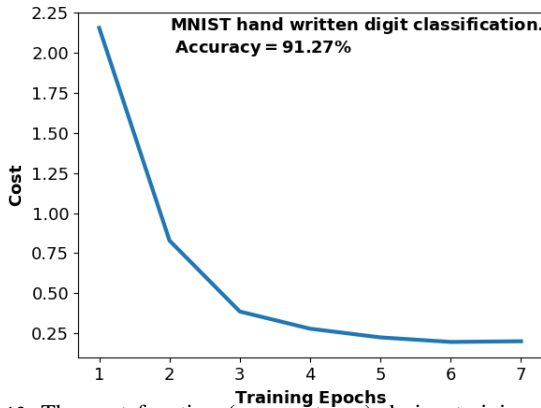


Figure 13: The cost function (cross-entropy) during training of proposed architecture on MNIST data set for digit classification.

Table X: Comparative Performance on MNIST Dataset

Reference	Accuracy (%)	Dataset	Key Details
[37]	97.16	MNIST	4000 epochs, Matlab simulation mathematical model of network
[23]	82.2	MNIST	2T1M synapses, SimElectronics, MATLAB
[38]	92.39	Binarised MNIST (28 × 28)	1M synapses, $\mathcal{O}(n^2)$ update time, Neurosim
[39]	94.0	Binarised MNIST (20 × 20)	1M synapses, $\mathcal{O}(n^2)$ update time, Neurosim
Our Work	91.27	MNIST	1M synapse, $\mathcal{O}(1)$ update time, Python

REFERENCES

- [1] S. Prajapati, M. N. Mondal, and S. Sur-Kolay, "Memristive neural network with efficient in-situ supervised training," in *2022 IEEE 35th International System-on-Chip Conference (SOCC)*, 2022, pp. 1–6.
- [2] W. A. Wulf and S. A. McKee, "Hitting the memory wall: Implications of the obvious," *ACM SIGARCH computer architecture news*, vol. 23, no. 1, pp. 20–24, 1995, <https://doi.org/10.1145/216585.216588>.
- [3] C. Mead, "Neuromorphic electronic systems," *Proceedings of the IEEE*, no. 10, pp. 1629–1636, 1990.
- [4] S. H. Jo, T. Chang *et al.*, "Nanoscale memristor device as synapse in neuromorphic systems," *Nano letters*, no. 4, pp. 1297–1301, 2010.
- [5] S. P. Adhikari, H. Kim, R. K. Budhathoki *et al.*, "A circuit-based learning architecture for multilayer neural networks with memristor bridge synapses," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 62, no. 1, pp. 215–223, 2014.
- [6] M. S. Tarkov, "Oscillatory neural associative memories with synapses based on memristor bridges," *Optical Memory and Neural Networks*, vol. 25, pp. 219–227, 2016.
- [7] D. Soudry, D. Di Castro, A. Gal, A. Kolodny, and S. Kvatinisky, "Memristor-based multilayer neural networks with online gradient descent training," *IEEE transactions on neural networks and learning systems*, vol. 26, no. 10, pp. 2408–2421, 2015.
- [8] J. Shi and Z. Zeng, "Design of in-situ learning bidirectional associative memory neural network circuit with memristor synapse," *IEEE TETCI*, vol. 5, no. 5, pp. 743–754, 2021.
- [9] C. Li, D. Belkin, Y. Li *et al.*, "Efficient and self-adaptive in-situ learning in multilayer memristor neural networks," *Nature Communications*, vol. 9, no. 1, p. 2385, 2018.
- [10] O. Krestinskaya, K. N. Salama, and A. P. James, "Learning in memristive neural network architectures using analog backpropagation circuits," *IEEE TCAS*, vol. 66, pp. 719–732, 2019.
- [11] P. M. Sheridan, F. Cai, C. Du *et al.*, "Sparse coding with memristor networks," *Nature Nanotechnology*, vol. 12, no. 8, pp. 784–789, 2017.
- [12] Y. Zhang, X. Wang, and E. G. Friedman, "Memristor-based circuit design for multilayer neural networks," *IEEE TCAS I: Regular Papers*, vol. 65, no. 2, pp. 677–686, 2018.
- [13] E. Linn, "Memristive nano-crossbar arrays enabling novel computing paradigms," in *2014 IEEE International Symposium on Circuits and Systems (ISCAS)*, 2014, pp. 2596–2599.
- [14] C. Li, D. Belkin, Y. Li, P. Yan, M. Hu, N. Ge, H. Jiang, E. Montgomery, P. Lin, Z. Wang, J. P. Strachan, M. Barnell, Q. Wu, R. S. Williams, J. J. Yang, and Q. Xia, "In-memory computing with memristor arrays," in *2018 IEEE International Memory Workshop (IMW)*, 2018, pp. 1–4.
- [15] X.-B. Tian and H. Xu, "The design and simulation of a titanium oxide memristor-based programmable analog filter in a simulation program with integrated circuit emphasis," *Chinese Physics B*, vol. 22, no. 8, p. 088501, 2013.
- [16] P. Yao, H. Wu, B. Gao *et al.*, "Face classification using electronic synapses," *Nature Communications*, vol. 8, pp. 1–8, 2017.
- [17] M. Prezioso, F. Merrih-Bayat, B. Hoskins *et al.*, "Training and operation of an integrated neuromorphic network based on metal-oxide memristors," *Nature*, vol. 521, no. 7550, pp. 61–64, 2015.
- [18] Y.-P. Lin, C. H. Bennett, T. Cabaret *et al.*, "Physical realization of a supervised learning system built with organic memristive synapses," *Scientific Reports*, vol. 6, no. 1, pp. 1–12, 2016.
- [19] C. Yakopcic, R. Hasan, and T. M. Taha, "Memristor based neuromorphic circuit for ex-situ training of multi-layer neural network algorithms," in *IJCNN*, 2015, pp. 1–7.
- [20] F. Alibart, E. Zamanidoost, and D. B. Strukov, "Pattern classification by memristive crossbar circuits using ex situ and in situ training," *Nature Communications*, vol. 4, no. 1, pp. 1–7, 2013.
- [21] M. Elshamy, H. Mostafa, Y. H. Ghallab, and M. S. Said, "A novel nondestructive read/write circuit for memristor-based memory arrays," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 23, no. 11, pp. 2648–2656, 2015.
- [22] C. Wen, J. Zha, L. Xu, F. Ru, and S. Quan, "Research on perceptron neural network based on memristor," *IEEE Transactions on Industrial Electronics*, vol. 71, no. 8, pp. 9649–9657, 2024.
- [23] Z. Yan, J. Chen, R. Hu, T. Huang, Y. Chen, and S. Wen, "Training memristor-based multilayer neuromorphic networks with sgd, momentum and adaptive learning rates," *Neural Networks*, vol. 128, pp. 142–149, 2020.
- [24] R. Feng, J. Li, S. Xie, and X. Mao, "Efficient training method for memristor-based array using 1t1m synapse," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 70, no. 7, pp. 2410–2414, 2023.
- [25] S. Chartier and M. Boukadoum, "A bidirectional heteroassociative memory for binary and grey-level patterns," *IEEE Transactions on Neural Networks*, no. 2, pp. 385–396, 2006.
- [26] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [27] L. Chua and S. M. Kang, "Memristive devices and systems," *Proceedings of the IEEE*, vol. 64, no. 2, pp. 209–223, 1976.
- [28] C. Yakopcic, T. M. Taha, G. Subramanyam, and R. E. Pino, "Generalized memristive device spice model and its application in circuit design," *IEEE TCAD*, vol. 32, no. 8, 2013.
- [29] A. S. Oblea, A. Timilsina, D. Moore, and K. A. Campbell, "Silver chalcogenide based memristor devices," in *The 2010 International Joint Conference on Neural Networks (IJCNN)*, 2010, pp. 1–3.
- [30] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010.
- [31] K. Miller, K. S. Nalwa, A. Bergerud, N. M. Neihart, and S. Chaudhary, "Memristive behavior in thin anodic titania," *IEEE Electron Device Letters*, no. 7, pp. 737–739, 2010.
- [32] J. E. David Greenfield, Arezu Sarvestani and P. Baunach, "Nasa: Asteroids classification," 2018.
- [33] D. Dua and C. Graff, "UCI machine learning repository," 2017.
- [34] E. Rosenthal, S. Greshnikov, D. Soudry, and S. Kvatinisky, "A fully analog memristor-based neural network with online gradient training," in *ISCAS*, 2016, pp. 1394–1397.
- [35] M. Shevgoor, N. Muralimanohar, R. Balasubramonian, and Y. Jeon, "Improving memristor memory with sneak current sharing," in *2015 33rd IEEE International Conference on Computer Design (ICCD)*, 2015, pp. 549–556.
- [36] M. N. Mondal, S. Sur-Kolay, and B. B. Bhattacharya, "Test optimization in memristor crossbars based on path selection," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 42, no. 1, pp. 294–307, 2023.
- [37] S. He, J. Liu, H. Wang, and K. Sun, "A discrete memristive neural network and its application for character recognition," *Neurocomputing*, vol. 523, pp. 1–8, 2023.
- [38] Y. Yilmaz and F. Gul, "Design and optimisations of metal-oxide artificial synaptic device based machine learning model," *IEEE Transactions on Emerging Topics in Computational Intelligence*, pp. 1–11, 2024.
- [39] Y. Yilmaz, "Accuracy improvement in ag:a-si memristive synaptic device-based neural network through adadelat learning method on handwritten-digit recognition," *Neural Computing and Applications*, vol. 35, pp. 1–16, 09 2023.

APPENDIX

A. Glossary of parameters for memristor model

The parameters of the memristor model [28] for LTSpice used in Section V are described in Table XI.

Table XI: Parameters in Spice Model for Memristor [28]

Parameter	Relation to Physical Behaviors
a_1 and a_2	Closely related to the thickness of the dielectric layer in a memristor device, as more electrons can tunnel through a thinner barrier leading to an increase in conductivity.
b	Determines how much curvature is seen in the I-V curve relative to the applied voltage. This relates to how much of the conduction in the device is Ohmic and how much is due to the tunnel barrier.
A_p and A_n	These control the speed of ion (or filament) motion. This could be related to the dielectric material used since oxygen vacancies have different mobility depending which metal-oxide they are contained in.
V_p and V_n	These represent the threshold voltages. There may be related to the number of oxygen vacancies in a device in its initial state. A device with more oxygen vacancies should have a larger current draw that may lead to a lower switching threshold if switching is assumed to be based on the total charge applied
α_p, α_n, x_p and x_n	These determine where the state variable motion is no longer linear, and determine the degree to which the state variable motion is dampened. This could be related to the electrode metal used on either side of the dielectric film since the metals chosen may react to the oxygen vacancies differently.