# Stochastic Deep Learning in Memristive Networks

Anakha V. Babu and Bipin Rajendran Department of Electrical and Computer Engineering, New Jersey Institute of Technology, Newark, NJ, 07102, USA Email: {av442, bipin}@njit.edu

Abstract—We study the performance of stochastically trained deep neural networks (DNNs) whose synaptic weights are implemented using emerging memristive devices that exhibit limited dynamic range, resolution, and variability in their programming characteristics. We show that a key device parameter to optimize the learning efficiency of DNNs is the variability in its programming characteristics. DNNs with such memristive synapses, even with dynamic range as low as 15 and only 32 discrete levels, when trained based on stochastic updates suffer less than 3% loss in accuracy compared to floating point software baseline. We also study the performance of stochastic memristive DNNs when used as inference engines with noise corrupted data and find that if the device variability can be minimized, the relative degradation in performance for the Stochastic DNN is better than that of the software baseline. Hence, our study presents a new optimization corner for memristive devices for building large noise-immune deep learning systems.

#### I. Introduction

Inspired by the computational efficiency of human brain in processing unstructured data, neural networks have been explored since 1940s for a wide variety of data analytics applications. The latest generation of Deep Neural networks (DNNs) have achieved impressive successes rivaling typical human performance, thanks to their ability to capture hidden features from unstructured data using multiple layers of neurons [1]. However, as the number of layers (depth) of the networks increase, DNN training becomes computationally intense and time consuming due to the physically separated execution and memory units in conventional von Neumann machines. This has motivated the exploration of non-von Neumann architectures with closely integrated processing units and local memory elements in dense cross bar arrays with memristive devices [2].

It has been recently proposed that DNNs can be implemented by 2D cross bar arrays of resistive processing units (RPUs) that can store multiple analog states and adjust its conductivity with simple voltage pulses [3]. These RPU devices when implemented in a cross bar array can accelerate DNN training if all the weights in the array can be updated in parallel. In this scheme, vector cross product operations of  $\mathcal{O}(N^2)$  complexity, required for the back-propagation algorithm for network training, can be implemented with simple AND operation of stochastic bit streams representing the neuronal signals with a complexity of  $\mathcal{O}(1)$ . However, in order to maintain high accuracies, a stringent set of specifications have to be satisfied by memristive RPU devices, one of which is a resolution of 1000 conductance levels within a dynamic range of 10.

Several memristive devices have been explored for realizing cross bar arrays for neuromorphic systems [2]–[5]. Recently, Phase Change Memory (PCM) arrays have been used to store synaptic weights of a 3-layer neural network for handwritten digit classification achieving an accuracy of 82.9% in the MNIST database [6]. Numerical simulation studies based on experimentally observed programming characteristics of Pr<sub>1-x</sub>Ca<sub>x</sub>MnO<sub>3</sub> (PCMO) synapses suggest that MNIST recognition accuracies exceeding 90% is achievable [7]. Linear and symmetric conductance response are observed to improve accuracy; hence, several strategies have been proposed to compensate for the non linear and asymmetric conductance response of typical memristive devices [8]-[10]. It has been projected that NVM based DNNs can provide 25× speedup and up to  $3000\times$  improvement in power compared to GPU based implementations [11]. Therefore, highly efficient neuromorphic systems can be developed if memristive device characteristics can be improved and algorithms co-optimized to account for their non-ideal limitations.

Here, we first use stochastic weight updates to train a 4-layer deep network with double precision floating-point weights. We then study the performance of an equivalent network which uses a cross bar architecture with typical memristive devices with limited dynamic range, resolution, and conductance variability [12]. To the best of our knowledge, this paper presents the first study of noise resilience for inference using stochastic learning of DNNs with non-ideal memristive devices. The main insight from this study is that the key device parameter to optimize the learning efficiency of DNNs is the variability in its programming characteristics. DNNs with such memristive synapses, even with dynamic range as low as 15 and resolution of 32 levels, when trained based on stochastic updates can achieve close to the floating point base-line accuracies (within 3\%) in the benchmark hand-written digit recognition task. Furthermore, the degradation in performance of stochastic memristive DNNs when used as inference engines with noise corrupted data is better than that of the software baseline.

This paper is organized as follows: We first discuss the fundamental basics of DNN training and methods to accelerate training using stochastic weight updates. We then describe the network used for stochastic learning for handwritten digit classification and the memristive model for crossbar compatible implementation. Finally, we compare the performance of stochastic learning for this network with the software baseline and demonstrate the superior noise-tolerance characteristics of stochastic DNN based inference engines.

#### II. DNN TRAINING AND ACCELERATION

DNN training involves two steps - forward pass to calculate the activation functions and backward pass for calculating the weight update required for all the synapses in the network. During forward pass, the input y to neuron j in layer (k+1) is determined based on the outputs of the neurons in the previous layer and the strength of synapses between the two layers, according to the relation:

$$y_j^{(k+1)} = f\left(\sum_{i=1}^N w_{ij}^{(k)} x_i^{(k)}\right) \tag{1}$$

where f(x) is a non-linear transformation function such as  $\sinh(x), \tanh(x)$  or  $\operatorname{ReLU}(x) = xH(x)$ , with H(x) denoting the Heaviside step function. During back propagation, the error  $\delta$  at layer (k+1) is fed back to determine the error in the previous layer and determine the weight update as

$$\delta_j^{(k+1)} \propto \sum_{l=1}^N w_{jl}^{(k+1)} \delta_l^{(k+2)}, \quad w_{ij}^{(k)} \leftarrow w_{ij}^{(k)} \pm \eta x_i^{(k)} \delta_j^{(k+1)}$$
(2)

All these steps are of  $\mathcal{O}(N^2)$  complexity, and network training becomes highly time consuming when implemented on von Neumann machines with physically separated memory and computational units. DNN training can be accelerated using a cross bar array with memristive devices at the cross point representing the synaptic weights and neuronal computational circuits at the periphery [3].

Operations in equations 1 and 2 can also be parallelized by leveraging the fact that programming pulses applied at the periphery can be used to read and program all the cells of the array in parallel, if the currents involved are sufficiently small. While the forward and backward pass is implemented as a parallel read operation, the weight update operation can be implemented using the idea that coincidence detection (AND operation) of stochastic streams representing real numbers is equivalent to the multiplication operation.

In the stochastic computing framework [13], [14], a number  $x \in [0,1]$  can be represented as a Bernoulli sequence  $X = [x_1,x_2,x_3,\ldots,x_N]$  such that the binary random variable  $x_i$  has a probability  $P(x_i=1)=x$ , and N is the length of the Bernoulli sequence [15]. If a,b are real numbers that lie in the range [0,1], their product  $c=a\times b$  can be obtained by finding the average of the binary sequence C which represent the bitwise logical AND operation of the Bernoulli sequences A and B of variables a and b.

$$c = \frac{1}{N} \sum_{i=1}^{N} C_i \Longrightarrow E(c) = ab$$
,  $Var(c) = \frac{ab(1-ab)}{N}$  (3)

As the length of the Bernoulli sequence increases, the error in the estimated average decreases. Thus, multiplication operation can be implemented efficiently using simple logic gates or coincidence detection operation.

In the resistive processing units (RPUs) based implementation of DNNs [3], parallel weight update is achieved by applying programming pulses with amplitude  $\pm V_s/2$  on the cross-bar wires, where  $V_s$  is the minimum amplitude necessary

to alter the state of the RPU. To determine the product of  $x_i$  and  $\delta_j$  in equation 2, their equivalent stochastic bit streams are fed through the row and column in the cross bar; the RPU conductance will change depending on the coincidence of these two stochastic pulse streams. The stochastic weight update can be represented as

$$w_{ij}^{(k)} = w_{ij}^{(k)} + B\left(\sum_{n=1}^{BL} x_{i,n}^{(k)} \wedge \delta_{j,n}^{(k+1)}\right)$$
(4)

Here, BL refers to the bit length of the Bernoulli sequence and B is the minimum conductance resolution as a result of one pulse pair overlap. High-level design studies in [3] suggest that in order to maintain network accuracies that are close to the ideal software performance, these devices should have a resolution of at least 1000 programmable levels within a dynamic range of 10, which is a stringent requirement to realize in nanoscale devices.

In this study, we study the performance characteristics of DNNs implemented using models of memristive devices that are more representative of experimental devices today [12]. Some of these non-idealities include limited dynamic range, resolution, and variability in conductance levels attained during programming. We study the performance characteristics of DNNs trained using stochastic methods, and compare their performance when used as inference engines for noise-corrupted data. Our studies reveal that the stochastic learning method helps to mitigate some of the non-ideal effects of the variability that is inherent to nanoscale devices.

#### III. NETWORK ARCHITECTURE

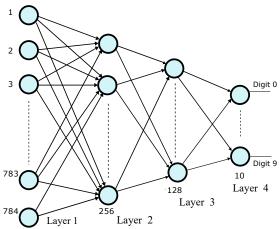


Fig. 1: 4 layer deep neural network with 784-256-128-10 neurons in each layer used for hand written digit classification (Simulated using MATLAB).

A 4-layer network with 784-256-128-10 fully connected neurons is used in our study for hand written digit classification (Fig. 1). Images from the MNIST database is used for training and testing the networks. The input images are pre-processed by mean normalization, and the network is trained by minimizing the multi-class cross-entropy objective function with sigmoid activation function for the hidden layers

and softmax function for the output layer. The weights are updated after every image (batch size of one) and a variable learning rate scheme is employed for learning. We refer to the network which uses stochastic pulses for forward pass, back propagation and weight update as 'stochastic DNN'.

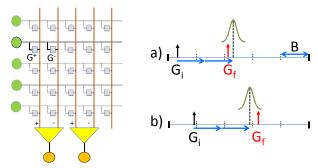


Fig. 2: Left: Unidirectional weight update scheme with 2 memristive devices per synapse in the cross bar. The synaptic weight is  $G_{eff}=G^+-G^-$  and devices are selectively programmed to increase or decrease  $G_{eff}$ ; Right: Illustration of the memristive programming - starting from initial conductance  $G_i$ , the final state  $G_f$  is determined by the pulse overlap (2 here), and a zero mean gaussian noise with standard deviation  $\sigma$  representing programming noise. In (a), the final state under-estimates the required conductance change, while (b) illustrates an over-estimate. The ratio  $\sigma/B$  determines the impact of programming noise; B denotes the resolution of the conductance levels.

In order to represent both positive and negative synaptic weights, two devices are used for each synapse [16]. The effective synaptic weight is the difference of the two device conductances,  $G_{eff} = G^+ - G^-$  as shown in Fig 2. Most memristive devices exhibit incremental programming only in one direction; hence we assume a unidirectional device programming scheme, where the conductance of the device always increments in the positive direction. To increase (decrease)  $G_{eff}$ ,  $G^+$  ( $G^-$ ) device is selectively programmed [17]. Due to the limited conductance resolution and on-off ratio, the device can saturate at its maximum conductance state, preventing further weight updates and learning. In order to avoid this and facilitate continuous learning, the devices are periodically reset after every 15 image is presented to the network. In our implementation, the memristive conductance response is assumed to be linear with an on-off ratio of 15 and have a resolution of 32 conductance states. First, we study the performance of stochastic DNN for handwritten digit classification as a function of the programming variability of the memristive device. We show that close to base-line accuracies can be maintained, even if the standard deviation  $(\sigma)$  in the programmed distribution is one-third of the separation in the levels (B) of the device. We then analyze the noise resilience characteristics of stochastic inference engines with devices which have close to ideal conductance variability.

## IV. RESULTS

The base line network response of the deep learning network with floating-point synaptic weight resolution is shown in Fig. 3. Each epoch of training consists of presentation of all the 60,000 images in the MNIST training set. For the stochastic DNN, the training error is a function of the number of bits (BL) in the stochastic code. After 30 epochs, the maximum test accuracy of the baseline-floating point DNN is 98.04% while that of the stochastic DNN with BL = 100 bits is 98.07%. The stochastic DNN test accuracy drops to 97.48% if BL = 10 is used. In the following sections, BL = 10 bits is used for simulations involving stochastic DNN due to its comparable accuracy with deterministic DNN, but  $10\times$  improvement in throughput and learning acceleration.

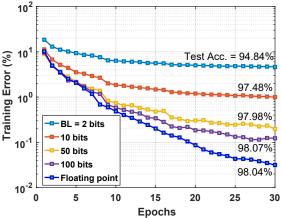


Fig. 3: Comparison of training error for stochastic DNN with BL=2,10,50,100 bits and a deterministic DNN, with floating point accuracy for synaptic weights. The accuracy on the test set after training is also shown. The network trained with BL=10 bits is used for stochastic DNNs in the rest of the paper.

We then performed numerical simulations incorporating the characteristics of the memristor devices for stochastic updates. Using memristive synapses with minimum programming variability ( $\sigma=0.1B$ ), the stochastic DNN test accuracy is 95.31%, which is 2.7% lesser than the baseline. The maximum test accuracy for stochastic DNNs as a function of the programming variability of memristive devices is shown in Fig. 4. There is negligible drop in test accuracy even at  $\sigma=0.3B$ , though the performance falls quickly as the programming variability increases further.

We now study the performance of stochastic DNNs that use close to ideal memristive synaptic weights with  $\sigma=0.1B$  when used as inference engines with noise corrupted test data. The network is trained using BL=10, but we study the inference accuracy for BL=10 and BL=100. Zero mean Gaussian noise with variances of  $\sigma_i^2=0.01,0.1$ , and 0.2 is added to the normalized MNIST test set (noise added values are kept in the range [0,1]). The test accuracy of the network for noise-corrupted data is shown in Fig. 5. Compared to their

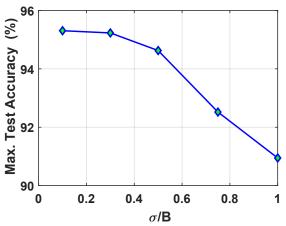


Fig. 4: The maximum generalization accuracy of stochastic DNN (BL=10) when trained with memristive devices with increasing programming variability. Here  $\sigma/B$  is the ratio of the standard deviation of the conductance variability to the bin-width of the conductance states.

respective non-noise test accuracy, the stochastic network with BL=10 has higher noise-resilience as its accuracy drops by 12.23% while that of the baseline floating point network drops by 14.35%. The stochastic network's degradation can be minimized further if more bits are used for stochastic encoding during inference, even though the learning was performed with BL=10. With BL=100, the stochastic network accuracy drops only by 9.2%. Hence, stochastic weight updates can compensate for noise in the input set as well as the variability introduced by nanoscale devices.

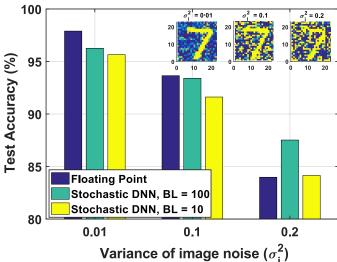


Fig. 5: Test accuracy for stochastic DNNs and the baseline floating point network when used as inference engine on noise-corrupted test set. The average inference response of 20 stochastic experiments is shown here. At  $\sigma_i^2 = 0.2$ , the stochastic DNN with BL = 100 out-performs the baseline. The inset shows the noise corrupted input image '7' corresponding to  $\sigma_i^2 = 0.01, 0.1$ , and 0.2.

## V. CONCLUSION

We demonstrate highly noise-resilient deep neural networks with memristive devices at the synapse, trained using stochastic updates. Even with significant device variability, limited on-off ratio and dynamic range, the performance of the memristive network is within 3% to that of the baseline floating point simulation. For efficient implementation of on-chip machine learning, algorithms that could circumvent the non-ideal characteristics of memristive devices have to be designed. The detrimental impact due to the non-linearities of nanoscale devices can be minimized when used in DNNs that use stochastic codes for data encoding and signal transmission in inference engines, especially for noisy inputs.

## ACKNOWLEDGMENT

This research was supported in part by the National Science Foundation grant 1710009 and CISCO Systems Inc.

## REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, May 2015. Insight.
- [2] G. W. Burr et al., "Neuromorphic computing using non-volatile memory," Advances in Physics: X, vol. 2, no. 1, pp. 89–124, 2017.
- [3] T. Gokmen and Y. Vlasov, "Acceleration of deep neural network training with resistive cross-point devices: Design considerations," Frontiers in Neuroscience, vol. 10, p. 333, 2016.
- [4] B. L. Jackson et al., "Nanoscale electronic synapses using phase change devices," J. Emerg. Technol. Comput. Syst., vol. 9, May 2013.
- [5] G. Indiveri, B. Linares-Barranco, R. Legenstein, G. Deligeorgis, and T. Prodromakis, "Integration of nanoscale memristor synapses in neuromorphic computing architectures," *Nanotechnology*, vol. 24, no. 38, p. 384010, 2013.
- [6] G. Burr et al., "Experimental demonstration and tolerancing of a large-scale neural network (165,000 synapses), using phase-change memory as the synaptic weight element," in *Electron Devices Meeting (IEDM)*, 2014 IEEE International, pp. 29.5.1–29.5.4, Dec 2014.
- [7] J. W. Jang, S. Park, G. W. Burr, H. Hwang, and Y. H. Jeong, "Optimization of conductance change in Pr<sub>1-x</sub>Ca<sub>x</sub>MnO<sub>3</sub>-based synaptic devices for neuromorphic systems," *IEEE Electron Device Letters*, vol. 36, pp. 457–459, May 2015.
- [8] P. Y. Chen et al., "Mitigating effects of non-ideal synaptic device characteristics for on-chip learning," in 2015 IEEE/ACM International Conference on Computer-Aided Design (ICCAD), Nov 2015.
- [9] A. Fumarola et al., "Accelerating machine learning with non-volatile memory: Exploring device and circuit tradeoffs," in 2016 IEEE International Conference on Rebooting Computing (ICRC), pp. 1–8, Oct 2016.
- [10] S. Sidler et al., "Large-scale neural networks implemented with non-volatile memory as the synaptic weight element: Impact of conductance response," in 2016 46th European Solid-State Device Research Conference (ESSDERC), pp. 440–443, Sept 2016.
- [11] G. W. Burr et al., "Large-scale neural networks implemented with non-volatile memory as the synaptic weight element: Comparative performance analysis (accuracy, speed, and power)," in 2015 IEEE International Electron Devices Meeting (IEDM), Dec 2015.
- [12] N. Panwar, B. Rajendran, and U. Ganguly, "Arbitrary spike time dependent plasticity (stdp) in memristor by analog waveform engineering," IEEE Electron Device Letters, vol. 38, pp. 740–743, June 2017.
- [13] W. J. Poppelbaum, C. Afuso, and J. W. Esch, "Stochastic computing elements and systems," in *Proc. of the November 14-16, 1967, Fall Joint Computer Conference*, AFIPS '67 (Fall), pp. 635–644, ACM, 1967.
- [14] B. R. Gaines, Stochastic Computing Systems. Springer US, 1969.
- [15] S. Gupta, V. Sindhwani, and K. Gopalakrishnan, "Learning machines implemented on non-deterministic hardware," *CoRR*, 2014.
- [16] O. Bichler et al., "Visual pattern extraction using energy-efficient 2-pcm synapse neuromorphic architecture," *IEEE Transactions on Electron Devices*, vol. 59, pp. 2206–2214, Aug 2012.
- [17] S. Lim et al., "Adaptive Learning Rule for Hardware-based Deep Neural Networks Using Electronic Synapse Devices," ArXiv e-prints, July 2017.