

This is the author-created version of the following work:

Lammie, Corey, Xiang, Wei, Linares-Barranco, Bernabé, and Rahimi Azghadi, Mostafa (2022) *MemTorch: an open-source simulation framework for memristive deep learning systems*. Neurocomputing, 485 pp. 124-133.

Access to this file is available from:

<https://researchonline.jcu.edu.au/72622/>

© 2022 Published by Elsevier B.V.

Please refer to the original source for the final version of this work:

<https://doi.org/10.1016/j.neucom.2022.02.043>

MemTorch: An Open-source Simulation Framework for Memristive Deep Learning Systems

Corey Lammie^a, Wei Xiang^b, Bernabé Linares-Barranco^c, Mostafa Rahimi
Azghadi^{a,*}

^a*College of Science and Engineering, James Cook University, Australia*

^b*School of Computing, Engineering and Mathematical Sciences, La Trobe University,
Australia*

^c*Institute of Microelectronics of Seville, IMSE-CNM, Parque Tecnológico de la Cartuja,
CSIC, University of Seville, Spain*

Abstract

Memristive devices have shown great promise to facilitate the acceleration and improve the power efficiency of Deep Learning (DL) systems. Crossbar architectures constructed using these Resistive Random-Access Memory (RRAM) devices can be used to efficiently implement various in-memory computing operations, such as Multiply Accumulate (MAC) and unrolled-convolutions, which are used extensively in Deep Neural Networks (DNNs) and Convolutional Neural Networks (CNNs). However, memristive devices face concerns of aging and non-idealities, which limit the accuracy, reliability, and robustness of Memristive Deep Learning Systems (MDLSSs), that should be considered prior to circuit-level realization. This Original Software Publication (OSP) presents *MemTorch*, an open-source¹ framework for customized large-scale memristive DL simulations, with a refined focus on the co-simulation of device non-idealities. MemTorch also facilitates co-modelling of key crossbar peripheral circuitry. MemTorch adopts a modernized software engineering methodology and integrates directly with the well-known PyTorch Machine Learning (ML) library.

Keywords: Memristors, RRAM, Non-Ideal Device Characteristics, Deep Learning, Simulation Framework

*Corresponding author

Email address: `mostafa.rahimiazghadi@jcu.edu.au` (Mostafa Rahimi Azghadi)

¹<https://github.com/coreylammie/MemTorch>

1. Introduction

MEMRISTIVE crossbar architectures [1] have been used to reduce the time complexity of Vector-Matrix Multiplications (VMMs) used in DNNs from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$, and in extreme cases to $\mathcal{O}(1)$ [2], facilitating the acceleration and improving the power efficiency of DL systems [3]. However, memristors are still considered an emerging technology, where their reliable manufacturing processes are yet to be achieved. As a result, DL architectures realized using memristor crossbars are putative to be prone to severe errors due to a number of device limitations including: finite discrete conductance states, device I/V non-linearity, failure, aging, cycle-to-cycle and device-to-device variability [4, 5]. Consequently, significant research efforts are being made to improve the reliability and robustness of memristive, or RRAM crossbars, used to perform *in-situ* learning [6, 7, 8, 9] and inference [2, 6, 10, 11, 12, 13] in DL systems. A general cross-platform, heterogeneous, high-level, customizable and open-source simulation framework with a refined focus on the co-simulation of device non-idealities could be used to conveniently build, rapidly prototype, and investigate device non-idealities in customized large-scale Memristive Deep Neural Networks (MDNNs) and MDLSs. In this OSP, we present such a framework, entitled *MemTorch*, for deep memristive learning using crossbar architectures. MemTorch is an open-source [14] simulation framework that integrates directly with the open-source PyTorch ML library that:

1. Facilitates the cross-platform development and distribution of large-scale [passive 0-Transistor 1-Resistor \(0T1R\)](#) and [active 1-Transistor 1-Resistor \(1T1R\)](#) memristive deep learning systems;
2. Places a large emphasis on modeling non-ideal, but inevitable, device characteristics in arbitrary and customizable device models;
3. Supports heterogeneous platforms such as Central Processing Units (CPUs) and Graphics Processing Units (GPUs);
4. Has a high-level Application Programming Interface (API), which is able to abstract performance-critical tasks described in various low-level languages.

2. Related Work

We compare MemTorch to other memristor-based DNN frameworks and inference accelerators, which are software-based and do not rely on SPICE

Table 1: Comparison of MemTorch to other memristor-based DNN simulation frameworks and inference accelerators.*Does not support GPU-accelerated inference and/or parameter mapping.†Models are shared using Google Drive without Application Programming Interfaces (APIs).

Simulation framework	Open-source	GPU	Pretrained DNN conversion	Programming language(s)
RAPIDNN [19]		✓*	✓	C++
MNSIM [20]			✓	Not Specified
PUMA [11]			✓	C++
DL-RSIM [21]		✓	✓	Python
PipeLayer [6]		✓*	✓	C++
Tiny but Accurate [22]	✓†		✓	MATLAB
Ultra-Efficient Memristor-Based DNN Framework [23]	✓†		✓	C++, MATLAB
Non-ideal Resistive Synaptic Device Characteristic Simulation Framework [24]		✓	✓	Python
Neurosim [25], NeuroSim+ [26], and DNN+NeuroSim [16, 17]	✓	✓	✓	C++, Python
IBM Analog Hardware Acceleration Kit [18]	✓	✓	✓	Python, C++, CUDA
MemTorch	✓	✓	✓	Python, C++, CUDA

36 modeling, in Table 1. [More exhaustive comparisons are performed in \[15\].](#)
37 Software-based frameworks and inference accelerators use a combination of
38 programming languages to simulate the behavior of memristive devices. Among
39 previous works, [DNN+NeuroSim \[16, 17\]](#) and [the IBM Analog Hardware Ac-](#)
40 [celeration Kit \[18\]](#) are the most similar offerings, which integrate with both
41 [PyTorch](#) and/or [Tensorflow](#), and can be used to account for non-ideal de-
42 vice characteristics. However, [they are](#) largely concerned with algorithm-
43 to-hardware mapping, and [are](#) designed to evaluate [training](#) and inference
44 accuracy with hardware constraints. [They are not](#) designed to model any
45 arbitrary device non-idealities for any behavioral device model. MemTorch,
46 on the other hand, [emphasizes](#) the co-simulation of non-ideal device charac-
47 teristics and generic behavioral device models with stochastic parameters for
48 higher flexibility to simply account for process variance.

49 3. Software Framework

50 The MemTorch simulation framework is programmed in C++, CUDA
51 and Python, with a Python interface. Performance critical tasks are per-
52 formed using either C++ or CUDA, for CPU or GPU execution, respectively;
53 otherwise Python is used. MemTorch relies heavily on the open source Py-
54 Torch [27] ML framework, and uses the C++ and Python PyTorch APIs
55 extensively to abstract low-level operations. Consequently, it supports na-
56 tive CPU and GPU operations.

57 3.1. Software Architecture

58 MemTorch is made up of seven distinct sub-modules. General utility
59 functions, such as data loaders or generic functions, are grouped within mem-

60 `torch.utils`. The `memtorch.bh` sub-module encapsulates all crossbar mod-
61 els, crossbar mapping and programming methods, crossbar tile mapping and
62 programming methods, memristor models, memristor model window func-
63 tions, models for all non-ideal device characteristics, quantization methods,
64 and methods to generate stochastic parameters. The `memtorch.mn` sub-
65 module mimics `torch.nn` and defines equivalent `memristivetorch.nn.Module`
66 layers. `memtorch.mn` currently extends `torch.nn.Linear`, `torch.nn.Conv1d`,
67 `torch.nn.Conv2d`, and `torch.nn.Conv3d`. `memtorch.mn.Module.patch_-`
68 `model` can be used to either instantiate new layers, or to patch existing in-
69 stances. `memtorch.mn.Module.patch_model()` iterates through and patches
70 all named modules within classes extending from `torch.nn.Module` and adds
71 a `self.tune_()` method, in addition to other helper methods, to the class
72 instance of the model that automatically patches each selected named mod-
73 ule.

74 The `memtorch.cpp` sub-module encapsulates all Python-wrapped C++
75 extensions, whereas the `memtorch.cu` sub-module encapsulates all Python-
76 wrapped CUDA extensions. Currently, MemTorch uses [C++ and CUDA](#)
77 [bindings to perform inference for both active and passive modular tiled ar-](#)
78 [chitectures, and to parallelize quantization operations.](#) The use of bind-
79 ings can be disabled, and legacy python methods (developed in previous
80 versions of MemTorch) can be used instead using the `use_bindings` argu-
81 ment, when patching `torch.nn.Module` instances. `memtorch.examples` sub-
82 module encapsulates general-usage examples and supporting scripts. The
83 `memtorch.map` sub-module encapsulates all mapping and tuning algorithms
84 used when programming and tuning memristive crossbar arrays. Finally, the
85 `memtorch.submodule` sub-module encapsulates all external git sub-modules
86 that MemTorch uses. We review and present the algorithms and models
87 that are currently built into MemTorch in Appendix A, and our approach to
88 modeling non-ideal device characteristics in Appendix B.

89 *3.2. Software Functionalities*

90 Complete examples demonstrating the functionality of MemTorch are
91 publicly accessible². ReadTheDocs³ is used to explain all functionalities.

²<https://github.com/coreylammie/MemTorch/tree/master/memtorch/examples>

³<https://memtorch.readthedocs.io/en/latest/>

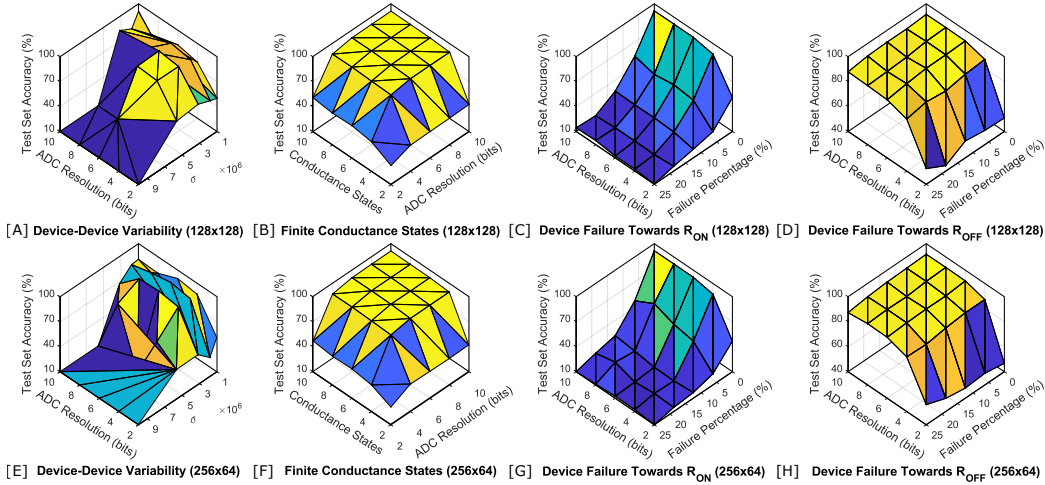


Figure 1: Simulation results when exemplar device non-idealities are considered for classifying CIFAR-10 dataset for two different modular crossbar tile sizes, 128x128 and 256x64. While MemTorch can be used to simulate both passive and active architectures, for demonstration purposes, in this figure, only active architectures are considered.

92 4. Implementation and Empirical Results

93 In Fig. 1, we present exemplar large-scale deep learning simulations to investigate the performance degradation due to device-device variability, finite
 94 number of conductance states, and device failure when two different modular
 95 crossbar sizes are considered. A separate case study is not presented, as
 96 we have previously used MemTorch in other works to perform hand gesture
 97 classification [28], epileptic seizure prediction [29], to develop an empirical
 98 metal-oxide device endurance and retention model [30], and to develop an
 99 extended Design Space Exploration (DSE) methodology for RRAM archi-
 100 tectures [31]. Prior to simulation, all convolutional and linear layers from a
 101 pre-trained MobileNetV2 for CIFAR-10 were converted to memristive equiv-
 102 alent layers, as explained in detail in ⁴, and documented in Appendix C.
 103

104 5. Illustrative Example

105 To demonstrate MemTorch’s intuitive design, we depict a typical use-
 106 case work flow in Fig. 2. Here, `torch.nn.Linear` layers are converted to

⁴https://github.com/coreylammie/MemTorch/tree/master/memtorch/examples/Exemplar_Simulations.ipynb

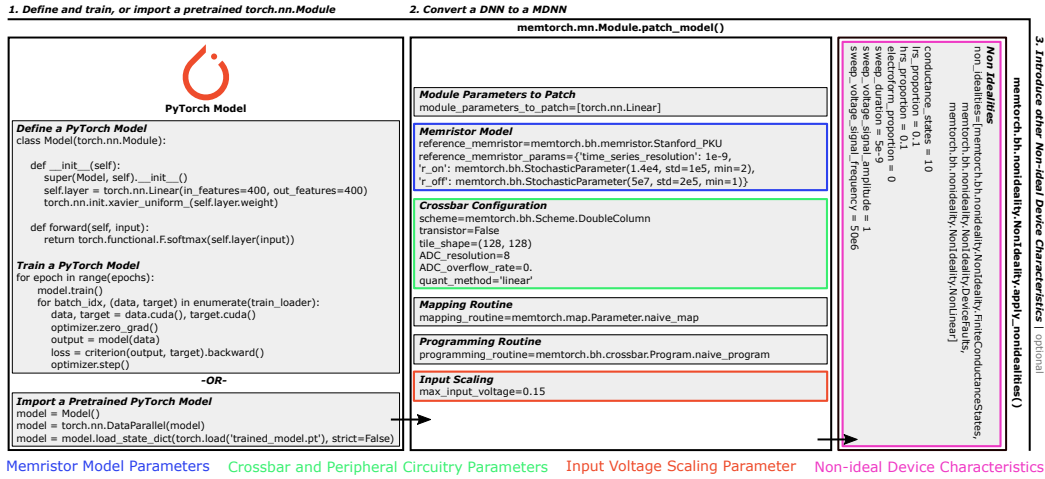


Figure 2: Illustration of a typical use-case workflow in MemTorch.

107 equivalent memristive layers constructed using modular crossbar tiles, that
 108 each contain (128×128) devices, which represent weights using a double-
 109 column parameter representation scheme. Inputs are scaled between $\pm 0.15V$,
 110 and 8-bit Analog to Digital Converters (ADCs) are used to read out col-
 111 umn currents. The Stanford PKU RRAM model [32] is used to model
 112 TiN/Hf(Al)O/Hf/TiN devices from [33]. Three other non ideal device char-
 113 acteristics were also accounted for including a finite number (10) of discrete
 114 conductance states, device faults, and non-linear I/V device behavior.

115 6. Conclusion

116 We presented an open-source simulation framework, entitled *MemTorch*,
 117 for large-scale deep memristive crossbar architectures. We showed that Mem-
 118 Torch is designed with a focus to integrate any desired behavioral or exper-
 119 imental device model, and introduce arbitrary device non-idealities, while
 120 co-simulating crossbar and peripheral circuitry. We compared MemTorch to
 121 similar works, detailed its package structure, and performed exemplar sim-
 122 ulations to demonstrate its functionality. We hope that MemTorch will be
 123 adopted and expanded by the community to advance memristive deep learn-
 124 ing research and development endeavours.

125 Acknowledgements

126 CL acknowledges the JCU DRTPS, [CAS Society Pre-Doctoral Grant](#),
127 [and the IBM PhD Fellowship](#). MRA acknowledges a JCU Rising Star ECR
128 Leadership Grant.

129 Appendix A. Algorithms and Models

130 This appendix reviews and presents the algorithms and models that are
131 currently built into MemTorch.

132 *Appendix A.1. Memristive Device Models*

133 Within MemTorch, we use [five](#) base memristive device models that ex-
134 tend the `memtorch.bh.Memristor.Memristor` base class for our simulations.
135 These include the linear ion drift model by [34]; the VTEAM model by [35],
136 which is a general model for voltage-controlled memristors that can be used
137 to fit a large range of experimental device data; the Stanford PKU RRAM
138 model [32], which describes switching performance for bipolar metal oxide
139 RRAM; and [two versions](#) of the data-driven Verilog-A RRAM model [36],
140 which expresses device current-voltage characteristics and resistive switching
141 rate as a function of the bias voltage and the initial resistive state of each
142 device. For each base model, finite differences is used to obtain a numerical
143 solution for each discretized time-step, dt . While only [five](#) base [memristive](#)
144 models are currently supported natively, others, [which can model the equiv-](#)
145 [alent conductance of a memristive device for an arbitrary applied voltage](#)
146 [signal](#), such as those modelling Phase Change Memory (PCM) or other de-
147 vice technology behavior, can easily be integrated modularly by extending
148 `memtorch.bh.Memristor.Memristor`.

149 *Appendix A.2. Window Functions*

150 Within memristive device models, window functions are widely employed
151 to restrict the changes of the internal state variables to specified intervals [37].
152 MemTorch currently natively supports the Biolek [38], Jogelkar [39], and
153 Prodromakis [40] window functions, and can easily be extended to support
154 others.

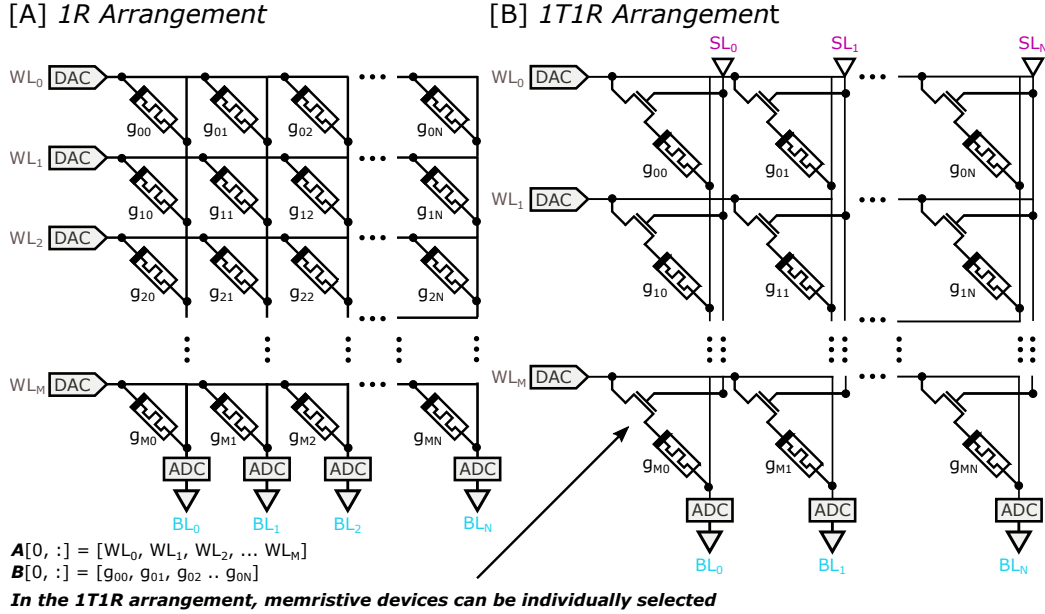


Figure A.3: Depiction of an $M \times N$ [A] 1R (0T1R) crossbar architecture and a [B] 1T1R crossbar architecture. Matrix-vector and matrix-matrix multiplication can be performed by encoding and presenting a scaled input vector or matrix \mathbf{A} as voltage signals to each row of the crossbar’s Word Lines (WLs). As shown in [A], assuming a linear I/V relationship, the total current in each column’s Bit Line (BL) is linearly proportional with the sum of the multiplication of the WL voltages and conductance values in that column, i.e., $\mathbf{BL}[0, :] \propto \mathbf{A}[0, :] \times \mathbf{B}$. In the 1T1R arrangement [B], individual memristive devices can be selected using Select Lines (SLs).

155 *Appendix A.3. Memristive Crossbar Architectures*

156 Memristive devices can be arranged within crossbar architectures to per-
 157 form VMMs, which are used extensively in forward and backward propa-
 158 gations within DNNs. There are two commonly used crossbar architecture
 159 configurations, namely 1-Transistor 1-Resistor (1T1R), and 1-Resistor (1R
 160 or /0T1R), which are both depicted in Fig. A.3. In 1T1R arrangements,
 161 one transistor is used to select and control each memristive device, whereas
 162 in 1R arrangements, rows and columns of memristive devices are positioned
 163 perpendicular to each other, with memristive devices sandwiched in-between.

164 The product of a vector and a matrix or, in a more general form, two
 165 matrices, \mathbf{A} of size $(M \times C)$ and \mathbf{B} of size $(C \times N)$, can be computed
 166 using a crossbar-architecture, as illustrated in Fig. A.3, where \mathbf{A} represents
 167 input voltage signals and \mathbf{B} is encoded within the crossbar as memristor

168 conductances. Separate ADCs can be used to read out the current of each
 169 column in parallel, as depicted, or sample and hold circuits can be used in
 170 conjunction with a single ADC per crossbar, that can be used to read out
 171 the current of each column sequentially using Time-division Multiplexing
 172 (TDM). As the output current of each column is linearly proportional to the
 173 elements of \mathbf{AB} , a linear constant, K , is used to correlate the ADC readout
 174 of each column accordingly. By separately presenting each row of \mathbf{A} to the
 175 crossbar through WLS, all rows of \mathbf{AB} can be computed.

Because memristors cannot be programmed to have negative conduc-
 tances, within MDNNs, weight matrices can either be represented using two
 devices per weight [41], as described by (A.1),

$$\mathbf{AB} = K \sum_{i=0}^C \mathbf{A}[i, :](g_{\text{pos}}[i, j] - g_{\text{neg}}[i, j]), \text{ for } j = 0 \text{ to } N, \quad (\text{A.1})$$

or using a single device per weight [42, 43] using complex weight mapping
 algorithms or current mirrors, as described by (A.2),

$$\mathbf{AB} = K \sum_{i=0}^C \mathbf{A}[i, :](g[i, j] - g_m), \text{ for } j = 0 \text{ to } N. \quad (\text{A.2})$$

176 For the single-column case, the current through g_m , used to mirror a current
 177 $-2V/(\bar{R}_{\text{ON}} + \bar{R}_{\text{OFF}})$ to each crossbar, is copied to each column and subtracted
 178 from all memristor columns. This current can be realized using a diode-
 179 connected NMOSFET by adjusting the NMOSFET channel width so that it
 180 has a passive resistance g_m . From this stage forward, we refer to the weight
 181 matrix representation methodology adopted, that is, whether two devices
 182 are used to encode each weight, i.e, differential weight mapping, one device
 183 is used to encode each weight, or another configuration is used to encode
 184 weight matrices, as the parameter representation scheme.

185 *Appendix A.3.1. Modular memristive crossbar tiles*

186 Mapping complete unrolled neural network layers into large memristive
 187 crossbar architectures often results in poor performance. This is due to
 188 non-ideal device characteristics that introduce substantial current variability
 189 when accumulated currents from columns with a large number of devices are
 190 read out. When one or two large crossbars are used, for single-column and
 191 double-column parameter representation schemes, respectively, they cannot

192 easily be modularized because customized crossbar shapes are required to
 193 represent each individual layer. Instead of using large crossbars, modular
 194 crossbar tiles [44] can be used that map layers into multiple uniformly sized
 195 crossbars, commonly referred to in literature as *crossbar tiles*.

One large crossbar of size $(M \times N)$ can be mapped using $\text{ceil}(M/S_0) \times$
 $\text{ceil}(N/S_1)$ crossbar tiles, each with a size of $(S_0 \times S_1)$, where the total uti-
 lization, ρ , of all crossbar tiles can be determined using (A.3),

$$\rho = \frac{MN}{\text{ceil}(M/S_0)\text{ceil}(N/S_1)S_0S_1}. \quad (\text{A.3})$$

196 Duplication of crossbar tiles and TDM can be used to regulate mapping to
 197 improve the array utilization and computation time by balancing latency
 198 among layers [45].

199 *Appendix A.3.2. Memristor crossbar programming*

200 The conductance of memristive devices can be altered between a low resis-
 201 tance state R_{ON} and a high resistance state R_{OFF} , by applying programming
 202 voltage pulses with different intervals and amplitudes. While individual de-
 203 vices within crossbars can be selected and programmed within 1T1R cells, in
 204 1R arrangements, when a voltage is applied to a specific device, a non-zero
 205 voltage (usually half that of the nominal programming pulse amplitude) is
 206 applied to all other devices in the same row and column. Consequently, vari-
 207 ous multistage programming [46, 47, 48, 49] and corrective methods [50, 2, 1],
 208 which can use analog voltage wave-forms, are often used to ensure the dif-
 209 ference between the programmed conductance states and the conductance
 210 states-to-program are within an acceptable tolerance.

211 *Appendix A.3.3. Memristor crossbar tuning*

212 The total current of each column in an ideal memristive crossbar is linearly
 213 proportional to the output elements of the VMM resultant vector. Conse-
 214 quently, after each DNN layer’s weights are programmed into a crossbar or
 215 group of tiles, linear regression can be used to correlate the output current
 216 of each column with any desired output to determine K for the crossbar or
 217 group of tiles, given a randomly generated input matrix that is sufficiently
 218 large. On account of device-device variations and device failures, further tun-
 219 ing is often required to recover accuracy loss and mitigate variances between
 220 intended and actual device conductance values. Tuning methods can either

Algorithm 1 Memristor crossbar programming algorithm.

Input: Array containing all continuous weights in a given layer, \mathbf{w} , HRS/LRS ratio, p_L .

Output: Equivalent memristive crossbars conductance values, \mathbf{g} , indexed using i and j .

```

 $\mathbf{w} = \text{abs}(\mathbf{w})$ 
 $\mathbf{w} = \text{descending\_order}(\mathbf{w})$ 
 $s = \text{size}(\mathbf{w})$ 
 $\text{index} = \text{int}(p_L \cdot s)$ 
 $\mathbf{w}_{\max} = \mathbf{w}[\text{index}]$ 
 $\mathbf{w}_{\min} = \mathbf{w}_{\max} / (\mathbf{R}_{\text{OFF}} / \mathbf{R}_{\text{ON}})$ 
 $\mathbf{w} = \text{clip}(\mathbf{w}, \mathbf{w}_{\min}, \mathbf{w}_{\max})$ 
 $\mathbf{g}[i, j] = \frac{(R_{\text{ON}} - R_{\text{OFF}}) \cdot (\sigma(\mathbf{w})[i, j] - \mathbf{w}_{\min})}{|\mathbf{w}|_{\max} - \mathbf{w}_{\min}} + R_{\text{OFF}}$ 

```

221 be used pre-programming [51], to improve robustness and reduce susceptibil-
 222 ity to error, or post-programming by retraining device-specific conductance
 223 values [22].

224 *Appendix A.3.4. Memristor crossbar weight mapping*

Weights, denoted using \mathbf{w} , within unrolled convolutional layers [52] and linear layers can be mapped to equivalent conductance values, \mathbf{g} , using (A.4).

$$\mathbf{g}[i, j] = \frac{(g_{\text{ON}} - g_{\text{OFF}})(\sigma(\mathbf{w})[i, j] - \mathbf{w}_{\min})}{|\mathbf{w}|_{\max} - \mathbf{w}_{\min}} + g_{\text{OFF}}, \quad (\text{A.4})$$

225 where \mathbf{w}_{\min} represents the minimum weight value to encode, and \mathbf{w}_{\max} rep-
 226 represents the maximum weight value to encode. $g_{\text{ON}} = 1/R_{\text{ON}}$ and $g_{\text{OFF}} =$
 227 $1/R_{\text{OFF}}$. When two crossbars are used to represent weight, crossbars con-
 228 taining positive components will have $\sigma(\mathbf{w}) = \mathbf{w}[\mathbf{w} \geq 0]$, while crossbars
 229 containing negative components will have $\sigma(\mathbf{w}) = \mathbf{w}[\mathbf{w} \leq 0]$. When a single
 230 crossbar is used to represent weights, $\sigma(\mathbf{w}) = \mathbf{w} - g_m$.

231 To reduce the inner weight gap in a given device, Algorithm 1 in [10]
 232 can be used to exclude a small proportion, p_L , of weights with the absolute
 233 largest values to reduce the variability effect of non-ideal memristive devices.

234 *Appendix A.3.5. Passive memristor crossbar architectures*

When modeling passive memristor crossbar architectures, source and line resistances should be accounted for. MemTorch utilizes a comprehensive crossbar array model with solutions for source and line resistances, as described in [53], to solve for node voltages \mathbf{V} , within passive crossbar architectures in each simulation time-step. Specifically, linear matrix algebra is used to solve (A.5)

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \mathbf{V} = \mathbf{E}, \quad (\text{A.5})$$

235 where for a crossbar of size $(M \times N)$, \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} are all $(MN \times MN)$
236 matrices, and \mathbf{E} is a $(2MN \times 1)$ vector. All matrices and vectors in (A.5)
237 are derived and defined in [53]. As the concatenated \mathbf{ABCD} matrix is
238 sparse, traditional linear matrix algebra methods cannot be easily used to
239 solve (A.5), because they require a prohibitive amount of memory. Instead,
240 sparse supernodal LU factorization with partial pivoting for general matrices
241 [54] is used, as it is parallelizable, and has demonstrated the best empirical
242 numerical performance, compared to related techniques, e.g., QR decomposition
243 [55].

244 *Appendix A.4. C++ and CUDA Bindings*

245 Numerous performance critical operations, including tiled VMMs, linear
246 matrix algebra, and quantization, are accelerated using C++ and CUDA
247 bindings. PYBIND11_MODULE() is a method within the pybind11⁵ python library
248 [56] that exposes C++ types in Python to enable seamless operability
249 between C++11, CUDA, and Python. This library is used within MemTorch
250 to overload method pointers and to expose C++ and CUDA functions to the
251 developed Python API. The Eigen [57] C++ template library is used extensively
252 to perform various linear algebra operations, as many Eigen functions
253 can be compiled for use within CUDA kernels using `__device__ __host__`
254 `_` function type qualifiers.

255 **Appendix B. Modeling Non-Ideal Device Characteristics**

256 Non-ideal device characteristics can either be encapsulated within device-
257 specific memristive models, or introduced to base (generic) models using

⁵<https://github.com/pybind/pybind11>

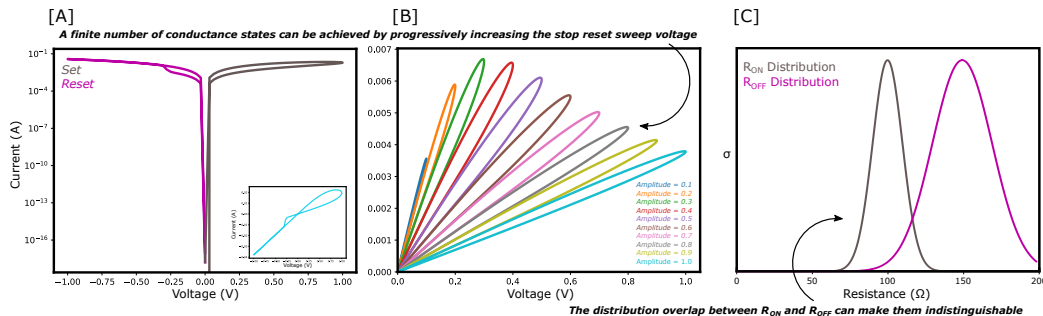


Figure B.4: Depiction of [A] device I/V characteristics, and [B] reset voltage double-sweeps demonstrating gradual switching from R_{ON} to R_{OFF} , which can be used to achieve 10 finite stable conductance states for the VTEAM model using the TEAM [58] model's parameters, with a linear dependence on w , achieved using sinusoidal signals with a fixed frequency of 50 MHz. [C] shows distributions of R_{ON} and R_{OFF} , which are caused by device-device variability, for a memristive device with $\bar{R}_{ON} = 100\Omega$ and $\bar{R}_{OFF} = 150\Omega$. In [C], overlapped regions are indistinguishable from each other.

258 the `memtorch.bh.nonideality` sub-module. This sub-module can currently
 259 be used to introduce four non-ideal device characteristics to memristive de-
 260 vice models: device-device variability, finite number of discrete conductance
 261 states, device failure, and non-linear I/V device characteristics. We leave
 262 native support for modeling other non-ideal device characteristics to future
 263 releases. Three of the non-ideal device characteristics that are currently sup-
 264 ported by MemTorch are shown in Fig. B.4. Fig. B.4[A] depicts typical
 265 non-linear I/V device characteristics using a set-reset curve and an inset
 266 hysteresis loop. Fig. B.4[B] demonstrates gradual switching, which is used
 267 to achieve a finite number of stable conductance states, and Fig. B.4[C] shows
 268 overlapping distributions of R_{ON} and R_{OFF} , which is caused by device-to-
 269 device variability.

270 Appendix B.1. Device-to-device Variability

271 Device-to-device variability is modeled stochastically using
 272 `memtorch.bh.StochasticParameter`. Stochastic parameters are generated
 273 using the `memtorch.bh.StochasticParameter.StochasticParameter()`
 274 method, which accepts an arbitrary number of keyword arguments, that are
 275 used to sample from a `torch.distributions` each time a device model is in-
 276 stantiated. To model device-device variability, we use stochastic parameters
 277 to sample R_{ON} and R_{OFF} from a normal distribution with $\sigma R_{ON} = \sigma$
 278 $\sigma R_{OFF} = 2\sigma$. $\sigma R_{OFF} > \sigma R_{ON}$, as the variability of R_{OFF} has been demon-

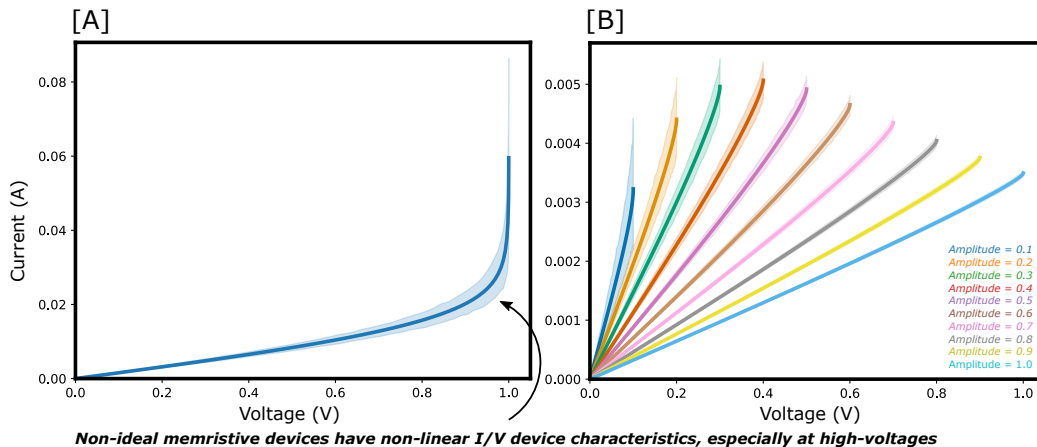


Figure B.5: Non-linear I/V characteristics for 100 devices (instances) of the VTEAM model using the TEAM [58] model’s parameters, with a linear dependence on w , achieved using sinusoidal signals with a fixed frequency of 50 MHz. R_{ON} and R_{OFF} were stochastically sampled from a normal distribution with $\bar{x} = 50, \sigma = 25$, and $\bar{x} = 1000, \sigma = 50$, respectively. [A] depicts I/V characteristics for devices with an infinite number of discrete conductance states. [B] depicts I/V characteristics for devices with a finite number of discrete conductance states.

279 stated to be larger than R_{ON} [59]. As depicted in Fig. B.4[C], device-device
 280 variability can cause the distribution of R_{ON} and R_{OFF} to overlap, resulting
 281 in R_{ON} and R_{OFF} occupying the same conductance regions.

282 *Appendix B.2. Cycle-to-cycle Variability*

283 Cycle-to-cycle (C2C) variability [60] is modeled stochastically, similarly
 284 to device-to-device variability, using stochastic parameters for R_{ON} and R_{OFF} .
 285 `memtorch.bh.nonideality.DeviceFaults.apply_cycle_variability()` is
 286 used to sample R_{ON} and R_{OFF} from a normal distribution with $\sigma R_{\text{ON}} = \sigma$
 287 and $\sigma R_{\text{OFF}} = 2\sigma$ after each SET RESET cycle.

288 *Appendix B.3. Finite Number of Discrete Conductance States*

289 Realistic memristive devices are non-ideal and have a finite number of
 290 stable discrete electrically switchable conductance states, bounded by a low-
 291 conductance semiconducting state R_{OFF} , and a high-conductance metallic
 292 state, R_{ON} [61]. Previous works have investigated evenly spaced conductance
 293 or resistance states, and have demonstrated that, assuming they are relatively
 294 uniformly distributed, the spacing between states is not critical [10].

295 Therefore, deterministic discretization [62] can be used to represent a
296 finite number of electrically switchable conductance states, as depicted in Fig.
297 B.4[B]. In order to efficiently quantize a tensor to a defined finite number of
298 quantization states, in which each element can have a different range, CUDA
299 kernels are used to perform a binary search on sorted tensors (generated using
300 the `linspace` algorithm in C++) containing defined quantization states in
301 $\mathcal{O}(n \log(n))$, where n is the number of quantized states.

302 *Appendix B.4. Device Failure*

303 Memristive devices are susceptible to failure, by either failing to electro-
304 form at a pristine state, or becoming stuck at high or low resistance states [10].
305 MemTorch incorporates a specific function for accounting for device failure
306 in simulating DL systems. Given a `nn.Module`, `memtorch.bh.nonideality.`
307 `DeviceFaults.apply_device_faults()` sets the conductance of a propor-
308 tion of devices within each crossbar to R_{ON} or R_{OFF} . It is assumed that the
309 total proportion of devices set to R_{OFF} is equal to the proportion of devices
310 that fail to electroform at pristine states plus the proportion of devices stuck
311 at a high resistance state. However, these proportions and the ratio of device
312 failures can be manipulated as desired. Devices are chosen at random using
313 `np.random.choice()`.

314 *Appendix B.5. Non-linear I/V Characteristics*

315 Non-ideal memristive devices have non-linear I/V device characteristics,
316 especially at high voltages, which are difficult to accurately and efficiently
317 model [10]. We demonstrate these characteristics using Fig. B.5[A], by
318 depicting the I/V curve of the VTEAM model between 0–1V using the
319 TEAM [58] model’s parameters. The `memtorch.bh.nonideality.NonLinear`
320 `.apply_non_linear()` method can be used to efficiently model non-linear
321 device I/V characteristics during inference for devices with an infinite num-
322 ber of discrete conductance states, and for devices with a finite number of
323 conductance states. For cases where devices are not simulated using their
324 internal dynamics, it is assumed that the change in conductance during read
325 cycles is negligible.

326 *Appendix B.5.1. Devices with an infinite number of discrete conductance* 327 *states*

328 The `memtorch.bh.nonideality.NonLinear.apply_non_linear()`
329 method uses two methods to efficiently model non-linear device I/V charac-

330 teristics for devices with an infinite number of discrete conductance states
331 during inference:

- 332 1. During inference, each device is simulated for a single timestep,
333 `device.time_series_resolution`, using `device.simulate()`.
- 334 2. Post weight mapping and programming, the I/V characteristics of each
335 device are determined using a single reset voltage sweep. The I/V
336 characteristics of each device are stored, and used as Lookup Tables
337 (LUTs) to compute device output currents during inference.

338 *Appendix B.5.2. Devices with a finite number of discrete conductance states*

339 The `memtorch.bh.nonideality.NonLinear.apply_non_linear()`
340 method effectively models non-linear I/V characteristics for devices with a
341 finite number of discrete conductance states by determining the I/V charac-
342 teristics of each device post weight mapping and programming during several
343 single reset voltage sweeps. Fig. B.5[B] depicts sweeps for 100 stochastic de-
344 vices with 10 finite discrete conductance states. These are stored and used as
345 LUTs to compute device output currents during inference, where each I/V
346 curve corresponds to each finite discrete conductance state. In Fig. B.5[B],
347 the smallest voltage amplitude corresponds to the finite conductance state
348 closest to R_{ON} , whereas the largest voltage amplitude corresponds to the
349 finite conductance state closest to R_{OFF} .

350 Appendix C. Exemplar Simulation Details

351 For all simulations performed to obtain the results presented in Fig. 1,
352 we followed the following training and test procedure. We first augmented
353 a pretrained MobileNetV2 CNN trained using the CIFAR-10 training set.
354 All convolutional and linear layers within the network were sequenced with
355 batch-normalization layers with fixed affine parameters to normalize out-
356 puts. The network was trained until improvement on the validation set was
357 negligible (for 100 epochs) with a batch size of $\mathfrak{S} = 256$. The initial learning
358 rate was $\eta = 1e - 1$, which was decayed by an order of magnitude every 40
359 training epochs. Stochastic Gradient Descent (SGD) was used to optimize
360 network parameters and Cross Entropy (CE) [63] was used to determine net-
361 work losses. The network achieved $> 90\%$ accuracy on the CIFAR-10 test
362 set.

363 When implementing the MDNNs, each memristive layer’s weights were
364 mapped to a double column line crossbar architecture adopting a 1T1R ar-
365 rangement. Linear regression was used to correlate the output current of
366 each column and its corresponding output to determine K for each crossbar,
367 given a randomly generated input matrix sampled from a uniform distribu-
368 tion between ± 1.0 . For linear layers, the random inputs had a size of $(8 \times$
369 $\text{in_features})$, while for convolutional layers the random inputs had a size of
370 $(8 \times \text{in_channels} \times 32 \times 32)$. Unless otherwise stated, inputs to memristive
371 layers were scaled from -0.3 to 0.3, to emulate voltage signals between $\pm 0.3\text{V}$,
372 which were applied to the word-lines of each memristive crossbar. All device
373 models originated from the VTEAM model, with $\bar{R}_{\text{ON}} = 1.4\text{e}4\Omega$ and \bar{R}_{OFF}
374 $= 5\text{e}7\Omega$, to model TiN/Hf(Al)O/Hf/TiN devices from [33].

375 Implementations are investigated using modular crossbar tiles of size
376 128×128 and 256×64 , as these have previously been demonstrated to be
377 effective in terms of utilization and power efficiency [45]. While power and
378 latency balancing is beyond the scope of MemTorch 1.1.5, 256×64 tile size
379 enables higher operation throughput and more analog operations per ADC
380 compared to 128×128 tile size [45]. However, the area utilization may be
381 lower for arrays with more than 64 columns considering the number of out-
382 put channels.

383 References

- 384 [1] M. Hu, H. Li, Y. Chen, Q. Wu, G. S. Rose, R. W. Linderman, Memristor
385 Crossbar-Based Neuromorphic Computing System: A Case Study, *IEEE*
386 *Transactions on Neural Networks and Learning Systems* 25 (2014) 1864–
387 1878.
- 388 [2] C. Lammie, O. Krestinskaya, A. James, M. R. Azghadi, Variation-aware
389 Binarized Memristive Networks, in: *Proc. 26th IEEE International Con-*
390 *ference on Electronics, Circuits and Systems (ICECS)*, Genoa, Italy,
391 2019, pp. 490–493.
- 392 [3] M. Rahimi Azghadi, Y. Chen, J. Eshraghian, J. Chen, C. Lin, A. Amir-
393 soleimani, A. Mehonic, A. Kenyon, B. Fowler, J. Lee, Y. Chang, Comple-
394 mentary Metal-Oxide Semiconductor and Memristive Hardware for Neu-
395 romorphic Computing, *Advanced Intelligent Systems* 2 (2020) 1900189.

- 396 [4] S. Mittal, A Survey of ReRAM-Based Architectures for Processing-
397 In-Memory and Neural Networks, *Machine Learning and Knowledge*
398 *Extraction* 1 (2018) 75–114.
- 399 [5] G. C. Adam, A. Khiat, T. Prodromakis, Challenges Hindering Memris-
400 tive Neuromorphic Hardware from Going Mainstream, *Nature commu-*
401 *nications* 9 (2018) 5267.
- 402 [6] L. Song, X. Qian, H. Li, Y. Chen, PipeLayer: A Pipelined ReRAM-
403 Based Accelerator for Deep Learning, in: *Proc. IEEE International Sym-*
404 *posium on High Performance Computer Architecture (HPCA)*, Austin,
405 TX, 2017, pp. 541–552.
- 406 [7] R. Hasan, T. M. Taha, C. Yakopcic, On-chip Training of Memristor
407 Based Deep Neural Networks, in: *Proc. International Joint Conference*
408 *on Neural Networks (IJCNN)*, Anchorage, AK, 2017, pp. 3527–3534.
- 409 [8] C. Lammie, J. K. Eshraghian, W. D. Lu, M. R. Azghadi, Memris-
410 tive Stochastic Computing for Deep Learning Parameter Optimization,
411 *IEEE Transactions on Circuits and Systems II: Express Briefs* (2021).
- 412 [9] O. Krestinskaya, K. N. Salama, A. P. James, Learning in Memristive
413 Neural Network Architectures Using Analog Backpropagation Circuits,
414 *IEEE Transactions on Circuits and Systems I: Regular Papers* 66 (2019)
415 719–732.
- 416 [10] A. Mehonic, D. Joksas, W. H. Ng, M. Buckwell, A. J. Kenyon, Simula-
417 tion of Inference Accuracy Using Realistic RRAM Devices, *Frontiers in*
418 *Neuroscience* 13 (2019) 593.
- 419 [11] A. Ankit, I. E. Hajj, S. R. Chalamalasetti, G. Ndu, M. Foltin, R. S.
420 Williams, P. Faraboschi, W. Hwu, J. P. Strachan, K. Roy, D. S. Milojevic,
421 PUMA: A Programmable Ultra-efficient Memristor-based Accelerator
422 for Machine Learning Inference, *CoRR* abs/1901.10351 (2019).
- 423 [12] H. Jeong, L. Shi, Memristor devices for neural networks, *Journal of*
424 *Physics D: Applied Physics* 52 (2018) 023003.
- 425 [13] H. Tsai, S. Ambrogio, P. Narayanan, R. M. Shelby, G. W. Burr, Re-
426 cent Progress in Analog Memory-based Accelerators for Deep Learning,
427 *Journal of Physics D: Applied Physics* 51 (2018) 283001.

- 428 [14] C. Lammie, W. Xiang, B. Linares-Barranco, M. R. Azghadi, corey-
429 lammie/MemTorch: Initial Release, 2020. URL: [https://doi.org/10.](https://doi.org/10.5281/zenodo.3760696)
430 [5281/zenodo.3760696](https://doi.org/10.5281/zenodo.3760696). doi:10.5281/zenodo.3760695.
- 431 [15] C. Lammie, W. Xiang, M. Rahimi Azghadi, Modeling and simulating
432 in-memory memristive deep learning systems: An overview of current
433 efforts, *Array* 13 (2022) 100116.
- 434 [16] X. Peng, S. Huang, Y. Luo, X. Sun, S. Yu, DNN+NeuroSim: An End-
435 to-End Benchmarking Framework for Compute-in-Memory Accelerators
436 with Versatile Device Technologies, in: *IEEE International Electron*
437 *Devices Meeting*, 2019.
- 438 [17] X. Peng, S. Huang, H. Jiang, A. Lu, S. Yu, DNN+NeuroSim V2.0: An
439 End-to-End Benchmarking Framework for Compute-in-Memory Accel-
440 erators for On-chip Training, 2020. [arXiv:2003.06471](https://arxiv.org/abs/2003.06471).
- 441 [18] M. J. Rasch, D. Moreda, T. Gokmen, M. Le Gallo, F. Carta, C. Gold-
442 berg, K. El Maghraoui, A. Sebastian, V. Narayanan, A Flexible and
443 Fast PyTorch Toolkit for Simulating Training and Inference on Ana-
444 log Crossbar Arrays, in: *Proceedings of the IEEE International Con-*
445 *ference on Artificial Intelligence Circuits and Systems (AICAS)*, 2021.
446 doi:10.1109/AICAS51828.2021.9458494.
- 447 [19] M. Imani, M. Samragh, Y. Kim, S. Gupta, F. Koushanfar, T. Rosing,
448 RAPIDNN: In-Memory Deep Neural Network Acceleration Framework,
449 *CoRR* abs/1806.05794 (2018).
- 450 [20] L. Xia, B. Li, T. Tang, P. Gu, P. Chen, S. Yu, Y. Cao, Y. Wang,
451 Y. Xie, H. Yang, MNSIM: Simulation Platform for Memristor-Based
452 Neuromorphic Computing System, *IEEE Transactions on Computer-*
453 *Aided Design of Integrated Circuits and Systems* 37 (2018) 1009–1022.
- 454 [21] M. Lin, H. Cheng, W. Lin, T. Yang, I. Tseng, C. Yang, H. Hu, H. Chang,
455 H. Li, M. Chang, DL-RSIM: A Simulation Framework to Enable Reliable
456 ReRAM-based Accelerators for Deep Learning, in: *Proc. IEEE/ACM*
457 *International Conference on Computer-Aided Design (ICCAD)*, San
458 Diego, CA, 2018, pp. 1–8. doi:10.1145/3240765.3240800.
- 459 [22] X. Ma, G. Yuan, S. Lin, C. Ding, F. Yu, T. Liu, W. Wen, X. Chen,
460 Y. Wang, Tiny but Accurate: A Pruned, Quantized and Optimized

- 461 Memristor Crossbar Framework for Ultra Efficient DNN Implementa-
462 tion, arXiv e-prints (2019) arXiv:1908.10017.
- 463 [23] G. Yuan, X. Ma, C. Ding, S. Lin, T. Zhang, Z. S. Jalali, Y. Zhao,
464 L. Jiang, S. Soundarajan, Y. Wang, An Ultra-Efficient Memristor-Based
465 DNN Framework with Structured Weight Pruning and Quantization Us-
466 ing ADMM, arXiv e-prints (2019) arXiv:1908.11691.
- 467 [24] X. Sun, S. Yu, Impact of Non-Ideal Characteristics of Resistive Synaptic
468 Devices on Implementing Convolutional Neural Networks, *IEEE Journal*
469 *on Emerging and Selected Topics in Circuits and Systems* 9 (2019) 570–
470 579.
- 471 [25] P. Chen, X. Peng, S. Yu, NeuroSim: A Circuit-Level Macro Model for
472 Benchmarking Neuro-Inspired Architectures in Online Learning, *IEEE*
473 *Transactions on Computer-Aided Design of Integrated Circuits and Sys-*
474 *tems* 37 (2018) 3067–3080.
- 475 [26] P. Chen, X. Peng, S. Yu, NeuroSim+: An integrated device-to-algorithm
476 framework for benchmarking synaptic devices and array architectures,
477 in: *IEEE International Electron Devices Meeting, 2017*, pp. 6.1.1–6.1.4.
478 doi:10.1109/IEDM.2017.8268337.
- 479 [27] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin,
480 A. Desmaison, L. Antiga, A. Lerer, Automatic differentiation in Py-
481 Torch, in: *NIPS Autodiff Workshop, 2017*.
- 482 [28] M. R. Azghadi, C. Lammie, J. K. Eshraghian, M. Payvand, E. Donati,
483 B. Linares-Barranco, G. Indiveri, Hardware Implementation of Deep
484 Network Accelerators Towards Healthcare and Biomedical Applications,
485 *IEEE Transactions on Biomedical Circuits and Systems* 14 (2020) 1138–
486 1159.
- 487 [29] C. Lammie, W. Xiang, and M. R. Azghadi, Towards Memristive Deep
488 Learning Systems for Real-time Mobile Epileptic Seizure Prediction, in:
489 *Proc. IEEE International Symposium on Circuits and Systems (ISCAS)*,
490 Daegu, South Koera., 2021.
- 491 [30] C. Lammie, M. R. Azghadi, D. Ielmini, Empirical Metal-Oxide RRAM
492 Device Endurance and Retention Model For Deep Learning Simulations,
493 *Semiconductor Science and Technology* (2021).

- 494 [31] C. Lammie, J. K. Eshraghian, C. Li, A. Amirsoleimani, R. Genov, W. D.
495 Lu, M. R. Azghadi, Design Space Exploration of Dense and Sparse
496 Mapping Schemes for RRAM Architectures (2022).
- 497 [32] Z. Jiang, S. Yu, Y. Wu, J. H. Engel, X. Guan, H. . P. Wong, Verilog-a
498 compact model for oxide-based resistive random access memory (rram),
499 in: International Conference on Simulation of Semiconductor Processes
500 and Devices (SISPAD), Yokohama, Japan., 2014, pp. 41–44.
- 501 [33] A. Fantini, L. Goux, A. Redolfi, R. Degraeve, G. Kar, Y. Y. Chen,
502 M. Jurczak, Lateral and Vertical Scaling Impact on Statistical Performances and Reliability of 10nm TiN/Hf(Al)O/Hf/TiN RRAM Devices,
503 in: Symposium on VLSI Technology, 2014.
- 504
- 505 [34] D. B. Strukov, G. S. Snider, D. R. Stewart, R. S. Williams, The Missing
506 Memristor Found, *Nature* 453 (2008) 80–83.
- 507 [35] S. Kvatinsky, M. Ramadan, E. G. Friedman, A. Kolodny, VTEAM: A
508 General Model for Voltage-Controlled Memristors, *IEEE Transactions*
509 *on Circuits and Systems II: Express Briefs* 62 (2015) 786–790.
- 510 [36] I. Messaris, A. Serb, S. Stathopoulos, A. Khiat, S. Nikolaidis, T. Prodromakis, A Data-Driven Verilog-A ReRAM Model, *IEEE Transactions on*
511 *Computer-Aided Design of Integrated Circuits and Systems* 37 (2018)
512 3151–3162.
513
- 514 [37] V. A. Slipko, Y. V. Pershin, Importance of the Window Function Choice
515 for the Predictive Modelling of Memristors, *CoRR* abs/1811.06649
516 (2018).
- 517 [38] Z. Biolek, D. Biolek, V. Biolková, Spice Model of Memristor With Non-
518 linear Dopant Drift, *Radioengineering* (2009) 210–214.
- 519 [39] Y. N. Joglekar, S. J. Wolf, The Elusive Memristor: Properties of Basic
520 Electrical Circuits, *European Journal of Physics* 30 (2009) 661–675.
- 521 [40] T. Prodromakis, B. P. Peh, C. Papavassiliou, C. Toumazou, A Versatile
522 Memristor Model With Nonlinear Dopant Kinetics, *IEEE Transactions*
523 *on Electron Devices* 58 (2011) 3099–3105.

- 524 [41] F. Alibart, E. Zamanidoost, D. B. Strukov, Pattern Classification by
525 Memristive Crossbar Circuits using ex situ and in situ Training, *Nature*
526 *Communications* 4 (2013) 2072.
- 527 [42] S. N. Truong, K.-S. Min, New Memristor-based Crossbar Array Archi-
528 tecture with 50-% Area Reduction and 48-% Power Saving for Matrix-
529 vector Multiplication of Analog Neuromorphic Computing, *Journal of*
530 *semiconductor technology and science* 14 (2014) 356–363.
- 531 [43] J. Lee, J. K. Eshraghian, K. Cho, K. Eshraghian, Adaptive Precision
532 CNN Accelerator Using Radix-X Parallel Connected Memristor Cross-
533 bars, *arXiv e-prints* (2019) arXiv:1906.09395.
- 534 [44] D. J. Mountain, M. R. McLean, C. D. Krieger, Memristor Crossbar
535 Tiles in a Flexible, General Purpose Neural Processor, *IEEE Journal on*
536 *Emerging and Selected Topics in Circuits and Systems* 8 (2018) 137–145.
- 537 [45] Q. Wang, X. Wang, S. H. Lee, F. Meng, W. D. Lu, A Deep Neural
538 Network Accelerator Based on Tiled RRAM Architecture, in: *IEEE*
539 *International Electron Devices Meeting (IEDM)*, San Francisco, CA.,
540 2019, pp. 14.4.1–14.4.4.
- 541 [46] I. E. Ebong, P. Mazumder, Self-Controlled Writing and Erasing in a
542 Memristor Crossbar Memory, *IEEE Transactions on Nanotechnology*
543 10 (2011) 1454–1463.
- 544 [47] M. S. Tarkov, Mapping Neural Network Computations onto Memristor
545 Crossbar, in: *Proc. International Siberian Conference on Control and*
546 *Communications (SIBCON)*, Omsk, Russia, 2015, pp. 1–4. doi:10.1109/
547 *SIBCON.2015.7147235*.
- 548 [48] R. Hasan, C. Yakopcic, T. M. Taha, Ex-situ Training of Dense Mem-
549 ristor Crossbar for Neuromorphic Applications, in: *Proc. IEEE/ACM*
550 *International Symposium on Nanoscale Architectures (NANOARCH)*,
551 Beijing, China, 2015, pp. 75–81. doi:10.1109/NANOARCH.2015.7180590.
- 552 [49] K. Jo, C. Jung, K. Min, S. Kang, Self-Adaptive Write Circuit for Low-
553 Power and Variation-Tolerant Memristors, *IEEE Transactions on Nan-*
554 *otechnology* 9 (2010) 675–678.

- 555 [50] B. Feinberg, S. Wang, E. Ipek, Making Memristive Neural Network
556 Accelerators Reliable, in: Proc. IEEE International Symposium on High
557 Performance Computer Architecture (HPCA), Vienna, Austria, 2018,
558 pp. 52–65. doi:10.1109/HPCA.2018.00015.
- 559 [51] B. Zhang, N. Uysal, D. Fan, R. Ewetz, Handling Stuck-at-faults in
560 Memristor Crossbar Arrays Using Matrix Transformations, in: Proc.
561 24th IEEE Asia and South Pacific Design Automation Conference (ASP-
562 DAC), ASPDAC '19, ACM, New York, USA, 2019, pp. 438–443. doi:10.
563 1145/3287624.3287707.
- 564 [52] K. Chellapilla, S. Puri, P. Y. Simard, High Performance Convolutional
565 Neural Networks for Document Processing, 2006.
- 566 [53] A. Chen, A comprehensive crossbar array model with solutions for line
567 resistance and nonlinear device characteristics, IEEE Transactions on
568 Electron Devices 60 (2013) 1318–1326.
- 569 [54] X. S. Li, M. Shao, A Supernodal Approach to Incomplete LU Factor-
570 ization with Partial Pivoting, ACM Trans. Math. Softw. 37 (2011).
- 571 [55] P. Matstoms, Sparse QR Factorization in MATLAB, ACM Trans. Math.
572 Softw. 20 (1994) 136–159.
- 573 [56] W. Jakob, J. Rhineland, D. Moldovan, pybind11 —
574 Seamless Operability Between C++11 and Python, 2016.
575 <https://github.com/pybind/pybind11>.
- 576 [57] G. Guennebaud, B. Jacob, et al., Eigen v3, <http://eigen.tuxfamily.org>,
577 2010.
- 578 [58] S. Kvatinsky, E. G. Friedman, A. Kolodny, U. C. Weiser, TEAM:
579 ThrEshold Adaptive Memristor Model, IEEE Transactions on Circuits
580 and Systems I: Regular Papers 60 (2013) 211–221.
- 581 [59] O. Krestinskaya, A. Irmanova, A. P. James, Memristive Non-Idealities:
582 Is there any Practical Implications for Designing Neural Network Chips?,
583 in: Proc. IEEE International Symposium on Circuits and Systems
584 (ISCAS), Sapporo, Japan, 2019, pp. 1–5. doi:10.1109/ISCAS.2019.
585 8702245.

- 586 [60] E. Miranda, A. Mehonic, W. H. Ng, A. J. Kenyon, Simulation of Cycle-
587 to-Cycle Instabilities in SiO_x -Based ReRAM Devices Using a Self-
588 Correlated Process With Long-Term Variation, *IEEE Electron Device*
589 *Letters* 40 (2019) 28–31.
- 590 [61] W. Yi, S. E. Savel'ev, G. Medeiros-Ribeiro, F. Miao, M.-X. Zhang, J. J.
591 Yang, A. M. Bratkovsky, R. S. Williams, Quantized Conductance Coin-
592 cides with State Instability and Excess Noise in Tantalum Oxide Mem-
593 ritors, *Nature Communications* 7 (2016) 11142.
- 594 [62] S. Yu, Neuro-inspired Computing with Emerging Nonvolatile Memorys,
595 *Proceedings of the IEEE* 106 (2018) 260–285.
- 596 [63] Z. Zhang, M. R. Sabuncu, Generalized Cross Entropy Loss for Training
597 Deep Neural Networks with Noisy Labels, *CoRR* abs/1805.07836 (2018).

598 **Required Metadata**

599 **Current executable software version**

Nr.	(executable) Software metadata description	Please fill in this column
S1	Current software version	1.1.5
S2	Permanent link to executables of this version	https://github.com/coreylammie/MemTorch/releases/tag/v1.1.5
S3	Legal Software License	GPLv3
S4	Computing platform/Operating System	Linux, OS X, Microsoft Windows
S5	Installation requirements & dependencies	A working Python interpreter (≥ 3.6). If CUDA is True in <code>setup.py</code> , CUDA Toolkit (≥ 10.1) and Microsoft Visual C++ Build Tools are required. If CUDA is False in <code>setup.py</code> , Microsoft Visual C++ Build Tools are required. All Python requirements are listed in https://github.com/coreylammie/MemTorch/blob/master/requirements.txt
S6	If available, link to user manual - if formally published include a reference to the publication in the reference list	https://memtorch.readthedocs.io/en/latest/
S7	Support email for questions	coreylammie@gmail.com , corey.lammie@jcu.edu.au

Table C.2: Software metadata (optional)

600 **Current code version**

Nr.	Code metadata description	Please fill in this column
C1	Current code version	1.1.5
C2	Permanent link to code/repository used of this code version	https://github.com/coreylammie/MemTorch
C3	Legal Code License	GPLv3
C4	Code versioning system used	git
C5	Software code languages, tools, and services used	Python, C++, CUDA.
C6	Compilation requirements, operating environments & dependencies	A working Python interpreter (≥ 3.6) . If CUDA is True in <code>setup.py</code> , CUDA Toolkit (≥ 10.1) and Microsoft Visual C++ Build Tools are required. If CUDA is False in <code>setup.py</code> , Microsoft Visual C++ Build Tools are required. All Python requirements are listed in https://github.com/coreylammie/MemTorch/blob/master/requirements.txt
C7	If available Link to developer documentation/manual	https://memtorch.readthedocs.io/en/latest/
C8	Support email for questions	coreylammie@gmail.com , corey.lammie@jcu.edu.au

Table C.3: Code metadata (mandatory)