
Noise Aware Finetuning for Analog Non-Linear Dot Product Engine

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As interest in analog acceleration for deep neural networks (DNNs) grows, ReRAM-
2 based Dot-Product Engines (DPEs) offer an energy-efficient solution for performing
3 vector-matrix multiplications (VMMs) in the analog domain. However, DPEs re-
4 quire Analog-to-Digital Converters (ADCs), which contribute significantly to area
5 and power overhead, and rely on digital logic for operations such as non-linear
6 activations. This work presents an ADC-less DNN accelerator that leverages Ana-
7 log Content Addressable Memory (ACAM) to replace ADCs and digital activation
8 units. By training decision trees to approximate activation functions and program-
9 ming them to ACAMs, the novel Non-linear DPE (NL-DPE) enables arbitrary
10 activation to be implemented directly in the analog domain, supporting a broader
11 range of future DNN architectures. Additionally, we explore the inherent noise
12 present in real devices of both DPE and ACAM and propose noise-aware finetuning
13 techniques that mitigate accuracy loss, demonstrating notable improvements.

14 1 Introduction

15 DNNs have significantly grown in size and complexity, driving an increasing demand for memory
16 bandwidth and computation. As DNN models continue to expand, the energy consumption associated
17 with frequent data movement between memory and processing units has become a significant
18 bottleneck, known as the *memory wall*. Such bottlenecks have led to the development of In-Memory
19 Computing (IMC) accelerators, where computation occurs directly within the memory. Among the
20 various IMC approaches, ReRAM-based analog computing stands out as one of the most promising
21 solutions due to its potential for high energy efficiency and scalability. ReRAM cells can be organized
22 in a crossbar structure to design DPEs, which perform VMMs in the analog domain in a single step,
23 achieving low energy consumption and high parallelism [Shafiee et al. (2016); Ankit et al. (2019)].

24 However, DNNs consist of more than just VMMs; they also require non-linear activations, which are
25 typically executed using digital logic. Therefore, DPE’s analog outputs must be converted into digital
26 signals via ADCs. Unfortunately, ADCs are both energy and area-inefficient, significantly impacting
27 the overall efficiency of ReRAM-based accelerators. ADCs can consume more than 30% of the chip
28 area and account for over 50% of the total power consumption, creating a substantial bottleneck in
29 achieving truly energy-efficient DNN accelerators [Shafiee et al. (2016); Ankit et al. (2019)].

30 Another challenge for ReRAM-based DNN accelerators is the inherent noise in analog computing.
31 While many existing works tend to overlook or underestimate its effect, our experiments demonstrate
32 that noise plays a crucial role in determining the accuracy of computations on real hardware. In fact,
33 without careful attention to noise, accuracy can degrade to unacceptable levels, as noted also in other
34 studies [Mao et al. (2022)]. Although some efforts have been made to address this issue [Momeni
35 et al. (2024)], novel computation primitives may require to carefully handle new sources of noise.

36 In this work, we propose the NL-DPE, an analog DNN accelerator design that eliminates the need for
 37 ADCs and digital activation units. We convert activation functions into decision trees, which are then
 38 mapped onto the programmable ACAM for analog computation. Since ACAM accepts analog inputs
 39 and produces digital outputs, ADCs are no longer required after the DPE. To address the reliability
 40 issues of analog computation, we measure noise directly from real devices and develop detailed
 41 noise models that account for various sources, including programming inaccuracies and conductance
 42 fluctuations, etc. Noise models are integrated into a fine-tuning process to minimize the accuracy gap
 43 caused by noise.

44 2 Analog Computing with ReRAM

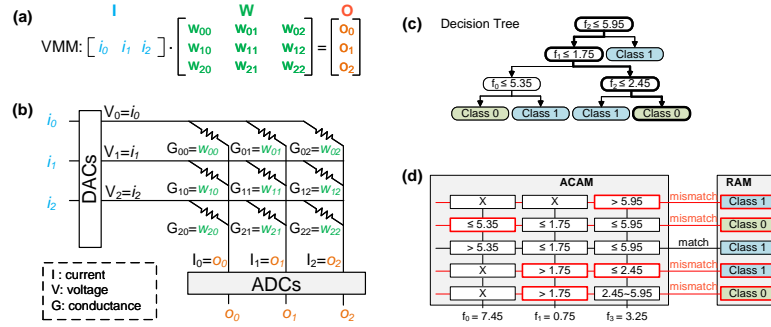


Figure 1: (a) An example of VMM computing. (b) Computing the VMM in DPE. (c) An example of a decision tree. (d) Mapping the decision tree onto ACAM and performing an inference.

45 Fig. 1(a) shows an example VMM operation ($\mathbf{O} = \mathbf{I} \cdot \mathbf{W}$). $\mathbf{O} \in \mathbb{R}^3$, $\mathbf{I} \in \mathbb{R}^3$ and $\mathbf{W} \in \mathbb{R}^{3 \times 3}$.
 46 Fig. 1(b) illustrates computing the VMM operation in the DPE with a 3×3 crossbar. A ReRAM cell
 47 is placed at every intersection of the horizontal wires and vertical wires. Weight elements (W_{ij}) are
 48 programmed as ReRAM conductance (G_{ij}), and input elements (i_i) are represented by voltage (V_i)
 49 on horizontal wires. Following Kirchhoff's law and Ohm's law, the current (I) from each vertical
 50 wire conveys the dot-product result of the input vector and the weight matrix. The analog VMM
 51 operation is completed in a single step within the DPE, achieving high computing parallelism.

52 Recently, ReRAM has been used to build Analog Content Addressable Memories (ACAMs), ac-
 53 celerating tree-based machine learning algorithms in the analog domain [Li et al. (2020); Pedretti
 54 et al. (2021b, 2023)]. Unlike digital CAM, which is limited to comparing a single input bit with a
 55 stored bit, an ACAM cell can compare an analog input value against a stored analog range. To map a
 56 trained decision tree of Fig. 1(c) onto an ACAM array, we traverse each leaf node back to the root,
 57 storing the feature thresholds along the path in a row of ACAM cells (Fig. 1(d)). Wildcard cells ('X')
 58 indicates that the particular feature is irrelevant along that path, thus the full range is programmed.
 59 Given an input feature vector during inference, if all feature values fall within the ranges stored in a
 60 row, the corresponding match line will be activated, retrieving the predicted class from an adjacent
 61 RAM. ACAMs have thus analog input, but digital output representing a match of mismatch.

62 3 Non-Linear Dot Product Engine

63 Fig. 2 shows a two-layer snippet from
 64 a larger neural network mapped to
 65 a conventional DPE (a) and the pro-
 66 posed NL-DPE (b). The first layer
 67 is followed by a Sigmoid activation
 68 function, while the second layer has
 69 no activation. Digital-to-Analog Con-
 70 verters (DACs) are required to convert
 71 the digital input into analog signals for
 72 computing VMM with the weights in
 73 the DPE. In conventional DPE-based
 74 accelerators (Fig. 2(a)), activations are
 75 typically performed in the digital do-

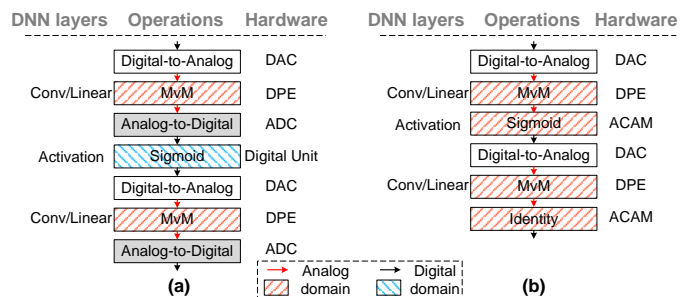


Figure 2: Mapping DNN on (a) conventional DPE and (b) NL-DPE.

76 main, thus ADCs are then needed to convert the analog output from the DPE back into digital form.
 77 In our proposed NL-DPE (Fig. 2(b)), we use ACAM to compute activations. As ACAM accepts
 78 analog input, ADCs are no longer needed after the DPE. The output of our ACAM is a digital signal,
 79 DACs are still needed before the next DPE. To replace the ADC, ACAM is programmed as an identity
 80 function, a special case of activation.

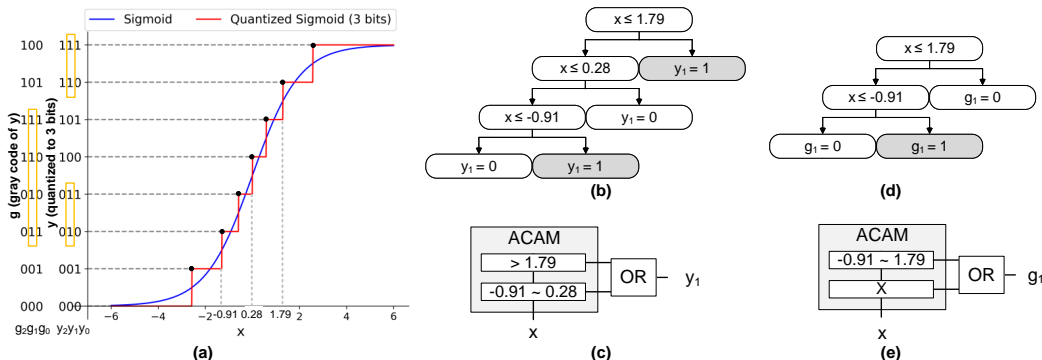


Figure 3: (a) Analytical and 3-bit quantized sigmoid function using conventional binary (y) and Grey code (g) format. (b) Trained decision tree to predict the second MSB of the activation with conventional binary representation (y_1) and (c) its mapping on ACAM. (d) Trained decision tree to predict the second MSB of the activation with Grey coding (g_1) and (e) its mapping on ACAM.

81 Fig. 3(a) demonstrates how ACAM is used to compute the Sigmoid activation with a 3-bit output. In
 82 this example, we assume 0 and 1 are quantized to 000_b and 111_b respectively, and the remaining 6
 83 values are evenly distributed between 0 and 1, as shown by the y axis in Fig. 3(a). Each output bit
 84 (y_2 , y_1 , and y_0) is computed using a separate decision tree, treating each as a binary classification
 85 task. The training dataset contains only one feature (input x), with the output bit as the target. For
 86 example, to train the decision tree for the second most significant bit (MSB) (y_1), the target is set
 87 to 1 when the input x falls within the range of -0.91 to 0.28 , or when x exceeds 1.79 , as shown in
 88 Fig. 3(a). Fig. 3(b) depicts the trained decision tree and Fig. 3(c) its mapping onto an ACAM array.
 89 Unlike traditional decision trees, which generalize to predict unseen inputs, these trees memorize the
 90 exact patterns from the training data, forcing overfitting. Note that the output of the ACAM already
 91 represents the bit value, there is no need to access to an attached memory as opposed to a Look Up
 92 Table (LUT) approach [Zhu et al. (2022)].

93 Since activations are typically monotonic functions, the less significant bits in the output tend to
 94 toggle more frequently between 0 and 1, leading to deeper decision trees. To mitigate this, we
 95 propose an encoding scheme based on Gray code to reduce bit toggling. Fig. 3(a) also shows the Gray
 96 code for the output y on the left (denoted as g), where only one bit changes between consecutive
 97 values. Fig. 3(d) shows the decision tree trained using Gray code as the output format and Fig. 3(e)
 98 its mapping to ACAM.

99 To enable subsequent computations, the Gray code output from the ACAM must be converted back
 100 to its original binary form, as illustrated in the Appendix Section A.

101 4 Mapping DNN layers into NL-DPE

102 Although ReRAM is inherently analog devices, most previous work maps DNN weights onto DPEs
 103 in a digital manner [Shafiee et al. (2016)]. Typically, this involves quantizing the full-precision
 104 floating-point DNN weights into fixed-point format (e.g., 8 bits) and using ReRAM cells to store
 105 portions of the quantized weights. Fig. 4(a) shows an example with four ReRAM cells used to store
 106 an 8-bit weight, with each cell representing 2 bits, an approach called bit slicing [Shafiee et al. (2016)]
 107 However, this approach is suboptimal, as it doesn't fully use the analog potential of ReRAM.

108 In our design, we take advantage of the analog properties of ReRAM to map unquantized DNN
 109 weights (Fig. 4 (b)), with a technique often referred to as *analog slicing* [Pedretti et al. (2021c); Song
 110 et al. (2024)]. Programming ReRAM cells involves an iterative program-and-verify (PV) process to
 111 reach a target conductance. However, achieving an exact value can be time-consuming, so a tolerance

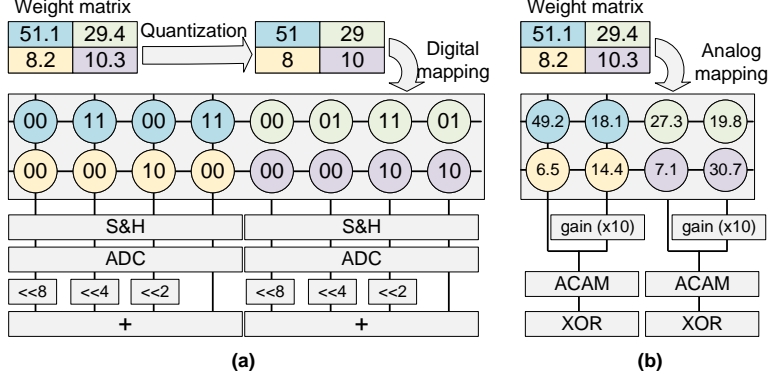


Figure 4: Mapping a weight matrix onto DPE in (a) digital form and (b) analog form.

112 threshold δ is used to stop once the error between the programmed conductance (G) and its target
 113 (W) is close enough (i.e., $|W - G| < \delta$). Typically, an error ϵ remains ($\epsilon = W - G$). In analog
 114 slicing a first Most Significant Cell (MSC) is iteratively programmed. In the example of Fig. 4(b),
 115 considering the target weight $W_{00} = 51.1\mu\text{S}$ the programmed MSC is $G_{00}^M = 49.2\mu\text{S}$, resulting in
 116 an error $\epsilon_{00}^M = 51.1 - 49.2 = 2.7\mu\text{S}$.

117 After programming all the MSCs, the resulting error is mapped into a second ReRAM cell, namely the
 118 least significant cell (LSC). First, the maximum error ϵ_{max}^M is used to compute a *gain* for programming
 119 the error exploiting the conductance full-scale range. For example, if the conductance range is 0 to
 120 $100\mu\text{S}$ and the maximum error is $10\mu\text{S}$, a gain $g = 10$ can be considered, as in our example. Thus,
 121 the LSC corrects the ϵ_{ij}^M error by being programmed to a conductance that represents the scaled-up
 122 error. In our example, the LSC is programmed targeting $27\mu\text{S}$. At the end of program-and-verify
 123 algorithm, the actual conductance of the LSC might be $G_{00}^L = 18.1\mu\text{S}$, with a programming error
 124 $\epsilon_{00}^L = 27 - 18.1 = 8.9\mu\text{S}$. During inference, both cells receive the input at the same time, and
 125 their output currents are combined, with the LSC's current scaled down to compensate for the error,
 126 attenuating it by g^{-1} . Thus, the final programming error for each weight is

$$\begin{aligned}
 \epsilon_{ij}^{TOT} &= W_{ij} - (G_{ij}^M - g^{-1}G_{ij}^L) \\
 &= W_{ij} - (G_{ij}^M + g^{-1}[g(W_{ij} - G_{ij}^M) + \epsilon_{ij}^L]) \\
 &= W_{ij} - G_{ij}^M - W_{ij} + G_{ij}^M - g^{-1}\epsilon_{ij}^L \\
 &= g^{-1}\epsilon_{ij}^L
 \end{aligned} \tag{1}$$

127 reducing the conductance error by a factor $g \gg 1$. In our example, the final noise is thus $51.1 -$
 128 $(49.2 + 18.1 \div 10) = 0.09\mu\text{S}$. We use the programming strategy described in [Mao et al. (2022)] to
 129 make the actual conductance always smaller than the target value.

130 For our experiments, we considered TaOx ReRAM devices integrated in 28 nm CMOS technology
 131 [Sheng et al. (2019)] and extracted a detailed noise model-based measurement on a fabricated crossbar
 132 test chip. The detailed noise model are described in Appendix Section B.

133 5 Noise Aware Finetuning for NL-DPE

134 With noise present in both DPE and ACAM of the NL-DPE, accuracy can be significantly affected.
 135 We propose a progressive approach that replaces ADCs and digital activation units with ACAM
 136 while accommodating noise. Fig. 5 outlines the process, which starts from a trained DNN model and
 137 produces a noise-tolerant model ready for deployment on hardware. In Step 1, activation quantization
 138 is applied, as DPE inputs must be processed bit by bit (i.e., input bit slicing [Shafiee et al. (2016)]).
 139 Step 2 introduces noise by sampling from the model in Appendix Section B and injecting it into
 140 the DPE, followed by a few iterations of noise aware fine-tuning (NAF) to compensate for DPE
 141 noise. This finetuning also helps to recover the accuracy lost due to quantization in Step 1. In
 142 Step 3, the quantization parameters are exported to generate the training data for the decision trees.
 143 Specifically, we need the activation function type to generate the graph similar to Fig. 3(a), and the

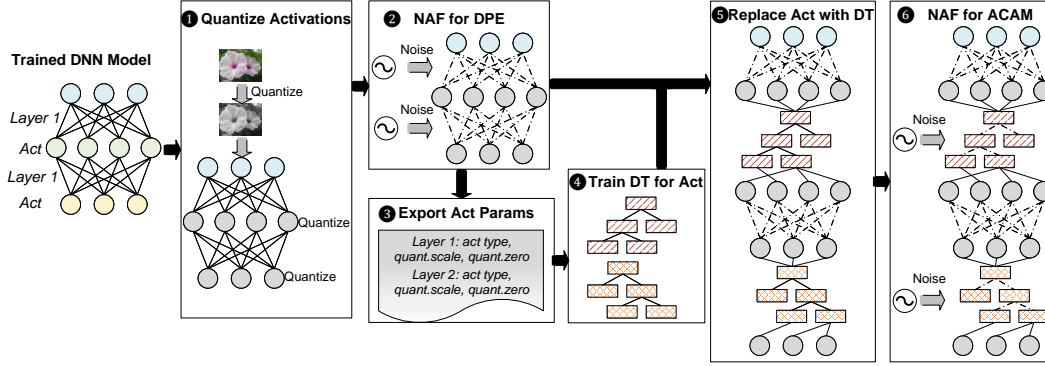


Figure 5: Steps of replacing activations with decision trees and using NAF to tolerate noise in NL-DPE.

144 quantization parameters (i.e., scaling factor and zero point) that are used to generate the output bit
 145 format (equivalent to y in Fig. 3(a)), which are then converted into Gray code format. Step 4 involves
 146 training decision trees for each activation using tools like scikit-learn. In Step 5, these decision trees
 147 replace the original activations in the DNN. Finally, in Step 6, we sample again from the noise model
 148 in Appendix Section B and inject noise into the decision tree thresholds, performing additional NAF
 149 iterations to account for ACAM noise.

150 6 Results

151 Table 1 shows inference accuracy for three neural network models on MNIST and CIFAR-10 datasets,
 152 all utilizing DPE for VMM and ACAM for ADCs and activation. We first consider an ideal DPE
 153 and ideal activation functions, in which only inputs and outputs are quantized. Then we introduce
 154 noise in the ReRAM conductance according to the model previously presented and programming
 155 with analog slicing, which results in an accuracy drop. We then implement Noise Aware Finetuning
 156 (NAF) as described in Section 5. We then train decision trees to approximate the activation functions
 157 and ADCs, we introduce noise for the ACAMs and perform NAF for the ACAMs as well. NAF
 158 helps recover the accuracy both for DPE and ACAMs. Note that we did not exhaustively explore all
 159 training hyperparameters, the accuracy was largely recovered, suggesting that further optimization of
 160 hyperparameters could potentially close this gap. We envision more sophisticated training algorithms
 for the NL-DPE employing differentiable ACAMs [Pedretti et al. (2022)].

Table 1: Inference accuracy of the proposed NL-DPE, also considering the noise in DPE and ACAM.

	LeNet	ResNet18	ResNet50
	MNIST	Cifar10	Cifar10
Baseline (FP32)	99.15	92.59	92.78
Quantized Activations	99.15	92.43	92.86
+ Noisy DPE	99.117	89.99	89.03
+ NAF for Noisy DPE	99.133	90.24	90.23
+ ACAM	99.067	90.09	89.69
+ Noisy ACAM	98.483	89.13	87.53
+ NAF for noisy ACAM	98.917	84.88	86.34

161

162 NL-DPE completely removes the ADC and digital activation unit in conventional DPE accelerators.
 163 We compare the power and area between our proposed NL-DPE computing unit and conventional DPE-
 164 based computing unit (including ADC and activation units needed by DPE) in Table 2. We take the
 165 ADC data from [Shafiee et al. (2016)] and scaled down to 28nm to match the technology node in NL-
 166 DPE. For digital activation unit, we take the data from FlexSFU [Reggiani et al. (2023)], a piecewise
 167 linear (PWL) based implementation to approximate arbitrary activation in digital domain. While

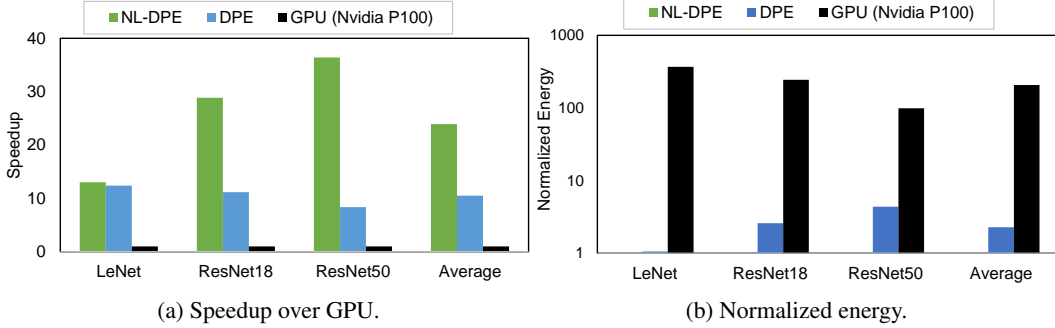


Figure 6: Speedup and normalized energy of running the benchmarks on different accelerators.

168 ACAM’s area is larger than that of ADC, it is significantly smaller than the FlexSFU. Considering the
 169 whole computing unit, NL-DPE’s area and power are both lower than that of DPE based computing.

Table 2: Power and area of ACAM, ADC and digital activation unit.

	Params	Spec	Power (mW)	Area (mm ²)
NL-DPE computing unit				
DPE	array size number	256 × 256 1	0.82	1.73
ACAM	array size number	16 × 1 8	0.18	0.009
Total			1	1.74
DPE computing unit				
DPE	array size number	256 × 256 1	0.82	1.73
ADC	resolution frequency	8bits 1.2GSps	1.71	0.001
FlexSFU	bit width	8bits	3.7	0.015
Total			6.23	1.75

170 The speedup and energy of running the whole model on three different accelerators are shown in
 171 Fig. 6. DPE is a conventional in-memory computing accelerator with ADC and digital activation units
 172 [Shafiee et al. (2016)], while NL-DPE replaces ADC and digital activation units with the proposed
 173 ACAM approach. We also test on NVidia P100 as the GPU baseline. Larger models have higher
 174 speedup and energy saving due to more activation can be paralleled using ACAM. On average,
 175 NL-DPE achieves 23× speedup and over 200× energy saving over GPU. Compared with DPE,
 176 NL-DPE achieves 2× speedup and 2.5× energy saving on average.

177 7 Conclusion

178 We presented the NL-DPE, a novel ADC-less analog in-memory computing primitive. NL-DPE
 179 builds on previous works on crossbar arrays and ACAM, by connecting them in the analog domain,
 180 with the former accelerating dot products and the latter activation functions, effectively implementing
 181 $f(xW)$ in the analog domain. Activations are approximated with decision trees which the ACAM
 182 can efficiently accelerate. We propose a Noise Aware Finetuning (NAF) routing to increase accuracy
 183 in the presence of ReRAM analog noise and we benchmark the proposed accelerator against the
 184 conventional DPE and GPU, reaching significant improvements. We envision NL-DPE as a new
 185 building block for in-memory computing, opening up new possibilities for model design fully
 186 exploiting the programmable analog non-linearity operation.

187 References

188 Ankit Aayush, Hajj Izzat El, Chalamalasetti Sai Rahul, Ndu Geoffrey, Foltin Martin, Williams

- 189 *R Stanley, Faraboschi Paolo, Hwu Wen-mei W, Strachan John Paul, Roy Kaushik, others* . PUMA:
190 A programmable ultra-efficient memristor-based accelerator for machine learning inference //
191 Proceedings of the twenty-fourth international conference on architectural support for programming
192 languages and operating systems. 2019. 715–731.
- 193 *Ielmini Daniele*. Resistive switching memories based on metal oxides: mechanisms, reliability and
194 scaling // *Semiconductor Science and Technology*. 2016. 31, 6. 063002.
- 195 *Li Can, Graves Catherine E, Sheng Xia, Miller Darrin, Foltin Martin, Pedretti Giacomo, Strachan*
196 *John Paul*. Analog content-addressable memories with memristors // *Nature communications*.
197 2020. 11, 1. 1638.
- 198 *Mao Ruibin, Wen Bo, Jiang Mingrui, Chen Jiezh, Li Can*. Experimentally-validated crossbar model
199 for defect-aware training of neural networks // *IEEE Transactions on Circuits and Systems II:*
200 *Express Briefs*. 2022. 69, 5. 2468–2472.
- 201 *Momeni Ali, Rahmani Babak, Scellier Benjamin, Wright Logan G, McMahan Peter L, Wanjura*
202 *Clara C, Li Yuhang, Skalli Anas, Berloff Natalia G, Onodera Tatsuhiko, others* . Training of
203 Physical Neural Networks // *arXiv preprint arXiv:2406.03372*. 2024.
- 204 *Pedretti Giacomo, Ambrosi Elia, Ielmini Daniele*. Conductance variations and their impact on the
205 precision of in-memory computing with resistive switching memory (RRAM) // *2021 IEEE*
206 *International Reliability Physics Symposium (IRPS)*. 2021a. 1–8.
- 207 *Pedretti Giacomo, Graves Catherine E, Serebryakov Sergey, Mao Ruibin, Sheng Xia, Foltin Martin,*
208 *Li Can, Strachan John Paul*. Tree-based machine learning performed in-memory with memristive
209 analog CAM // *Nature communications*. 2021b. 12, 1. 5806.
- 210 *Pedretti Giacomo, Graves Catherine E, Van Vaerenbergh Thomas, Serebryakov Sergey, Foltin Martin,*
211 *Sheng Xia, Mao Ruibin, Li Can, Strachan John Paul*. Differentiable content addressable memory
212 with memristors // *Advanced electronic materials*. 2022. 8, 8. 2101198.
- 213 *Pedretti Giacomo, Mannocci Piergiulio, Li Can, Sun Zhong, Strachan John Paul, Ielmini Daniele*.
214 Redundancy and analog slicing for precise in-memory machine learning—Part I: Programming
215 techniques // *IEEE Transactions on Electron Devices*. 2021c. 68, 9. 4373–4378.
- 216 *Pedretti Giacomo, Moon John, Bruel Pedro, Serebryakov Sergey, Roth Ron M, Buonanno Luca,*
217 *Gajjar Archit, Ziegler Tobias, Xu Cong, Foltin Martin, others* . X-TIME: An in-memory engine
218 for accelerating machine learning on tabular data with CAMs // *arXiv preprint arXiv:2304.01285*.
219 2023.
- 220 *Reggiani Enrico, Andri Renzo, Cavigelli Lukas*. Flex-sfu: Accelerating dnn activation functions by
221 non-uniform piecewise approximation // *2023 60th ACM/IEEE Design Automation Conference*
222 *(DAC)*. 2023. 1–6.
- 223 *Shafiee Ali, Nag Anirban, Muralimanohar Naveen, Balasubramonian Rajeev, Strachan John Paul,*
224 *Hu Miao, Williams R Stanley, Srikumar Vivek*. ISAAC: A convolutional neural network accelerator
225 with in-situ analog arithmetic in crossbars // *ACM SIGARCH Computer Architecture News*. 2016.
226 44, 3. 14–26.
- 227 *Sheng Xia, Graves Catherine E, Kumar Suhas, Li Xuema, Buchanan Brent, Zheng Le, Lam Sity,*
228 *Li Can, Strachan John Paul*. Low-conductance and multilevel CMOS-integrated nanoscale oxide
229 memristors // *Advanced electronic materials*. 2019. 5, 9. 1800876.
- 230 *Song Wenhao, Rao Mingyi, Li Yunning, Li Can, Zhuo Ye, Cai Fuxi, Wu Mingche, Yin Wenbo, Li Zongze,*
231 *Wei Qiang, others* . Programming memristor arrays with arbitrarily high precision for analog
232 computing // *Science*. 2024. 383, 6685. 903–910.
- 233 *Zhu Hanqing, Zhu Keren, Gu Jiaqi, Jin Harrison, Chen Ray T, Incorvia Jean Anne, Pan David Z.*
234 *Fuse and mix: MACAM-enabled analog activation for energy-efficient neural acceleration* //
235 *Proceedings of the 41st IEEE/ACM International Conference on Computer-Aided Design*. 2022.
236 1–9.

237 A Reconstruction of Gray coded output

238 Conversion from Gray code to binary can be achieved using simple logic gates: each binary bit is ob-
 239 tained by performing an XOR operation with all the higher-order bits in the Gray code representation,
 240 except the MSB which is the same between these two formats, i.e.,

$$b_i = \begin{cases} g_i & i = n - 1 \\ XOR(g_{n-1}, g_{n-2}, \dots, g_{i+1}) & i < n - 1 \end{cases}$$

241 where b_i and g_i are the i th bit in the binary and Gray code format. n is the output bit width. The
 242 overhead of these XOR gates are included in the power and area estimation of ACAM.

243 B ReRAM Noise Model

244 Since weight values are represented by ReRAM conductance, the conductance should ideally remain
 245 stable and exhibit linear properties. However, real ReRAM devices exhibit various stochastic behav-
 246 iors, such as inaccurate programming, conductance relaxation, and conductance fluctuation [Pedretti
 247 et al. (2021a)]. Inaccurate programming refers to the error between the actual conductance and the
 248 target value, determined by the programming tolerance. After programming, the conductance may
 249 drift before stabilizing—a phenomenon known as conductance relaxation. Conductance fluctuation
 250 is temporary noise that occurs during each read operation. Both relaxation and fluctuation are
 251 conductance-dependent. For example, evidence shows a constant standard deviation of fluctuations
 252 for high conductance states, with a linear decrease once a strong filament is formed and a low
 253 resistance state is reached [Ielmini (2016)].

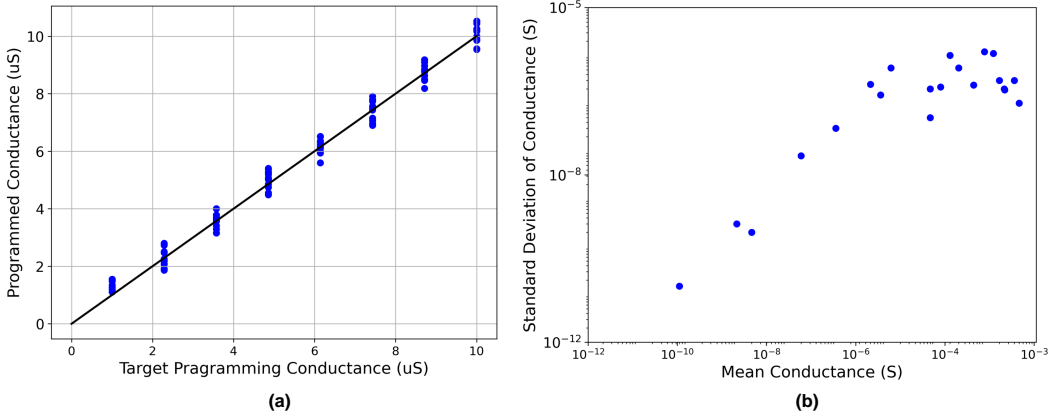


Figure 7: **(a)** Programming 8 levels at $10\mu\text{S}$ and below, with equal conductance spacing of $1.29\mu\text{S}$ and a fixed tolerance of $0.55\mu\text{S}$. **(b)** Standard deviation of read conductance as a function of mean conductance.

254 We use the same mathematical model from [Mao et al. (2022)] that accounts for all three ef-
 255 fects—inaccurate programming, conductance relaxation, and conductance fluctuation—and fit the
 256 model parameters using data from TaOx based ReRAM [Sheng et al. (2019)]. Fig. 7(a) shows the
 257 actual conductance values programmed at $10\mu\text{S}$ or below, with 8 levels spaced evenly at $1.29\mu\text{S}$ and
 258 a fixed tolerance of $0.55\mu\text{S}$. Fig. 7(b) illustrates the standard deviation (σ) of the read conductance
 259 as a function of the mean programmed conductance, demonstrating that while σ increases with the
 260 mean, it remains within $1\mu\text{S}$ across the entire range.

261 The conductance read out from the device is modeled as:

$$\mathbf{G}_{\text{read}} = \mathbf{G}_{\text{target}} + \mathbf{G}_{\text{relax}} + \mathbf{G}_{\text{fluc}} \quad (2)$$

262 $\mathbf{G}_{\text{target}}$ is target conductance converted from weight value w , i.e., $\mathbf{G}_{\text{target}} = w \cdot \mathbf{g}_{\text{ration}}$, we use
 263 $\mathbf{g}_{\text{ration}} = 75$. $\mathbf{G}_{\text{relax}}$ and \mathbf{G}_{fluc} are the noise caused by relaxation and fluctuation, respectively.
 264 After fitting the data, $\mathbf{G}_{\text{relax}}$ and \mathbf{G}_{fluc} are computed as:

$$\mathbf{G}_{\text{relax}} = N(0.064, 0.31^2) \quad (3)$$

$$\mathbf{G}_{\text{fluc}} = \exp(\log(\mathbf{G}_{\text{relax}} + \mathbf{G}_{\text{target}}) \cdot 0.75 - 2.93 + N(0, 0.98^2)) \cdot N(0, 1) \quad (4)$$

265 $N(\mu, \sigma^2)$ denotes the gaussian with mean μ and standard deviation σ .