

1 **A Per-Chip Metadata Correction Framework for Neuromorphic Memristor**  
2 **Crossbars:**  
3 **Theoretical Architecture and Viability Analysis**

4 Naman Boggaram

5 **Abstract**

6 We establish theoretically that device-level variability in memristor crossbars can be deter-  
7 ministically eliminated at model load time via per-chip system identification, removing the  
8 need for hardware-aware training. A system-identification protocol executed entirely through  
9 existing peripheral circuitry applies  $N$  Hadamard-structured input patterns to the array and  
10 recovers the full  $N \times N$  conductance deviation matrix via a single matrix multiply, achieving an  
11 effective measurement noise floor of  $\sigma_{\text{meas}}/\sqrt{N}$  — a  $63\times$  improvement over single-node sweeps  
12 for a  $4000 \times 4000$  array. This noise reduction converts the measurability condition from a hard  
13 instrument-precision prerequisite into a convergence condition satisfied by construction. The  
14 correction itself is exact under the confirmed linear response of the Ru(22) device class, allows  
15 any pre-trained model to be deployed without retraining or hardware-aware fine-tuning, and  
16 is maintained over device lifetime by a periodic heartbeat recalibration protocol that corrects  
17 all noise sources regardless of origin. Theoretical illustration under a Gaussian spatial correla-  
18 tion model consistent with published spin-coating physics confirms 99.98% variance capture at  
19  $K=16$  DCT coefficients (2 KB storage,  $31,250\times$  compression). Ten-year operational energy for  
20 a  $4000 \times 4000$  array under Architecture C is 140,836 J — below hardware-aware training for  
21 any per-run training cost above 3,520 J. Factory characterisation requires 17.6 s per chip using  
22 existing peripheral circuitry.

23 **1. Introduction**

24 **Scope and theoretical basis.** The analyses presented in this document are theoretical,  
25 derived from published device parameters of Goswami et al. [1]. Experimental validation of  
26 per-chip characterisation metadata generation, measurement noise bounds, and long-term drift  
27 behaviour under the heartbeat protocol requires chip-level access not available at the time of  
28 writing. This work establishes the theoretical framework and viability conditions; experimental  
29 confirmation is identified as future work.

30 Variability is the main blocker in analog neuromorphic hardware deployment. Memristor  
31 crossbar arrays offer orders-of-magnitude improvements in energy and area for matrix-vector  
32 multiplication [2, 3], but device-to-device conductance variations prevent reliable computation  
33 without correction.

34 Existing solutions have fundamental tradeoffs:

- 35 • **Hardware-aware training** couples models to hardware through noise injection during  
36 training, requiring retraining for every device class and producing residual errors since  
37 actual per-chip deviations differ from population statistics
- 38 • **Calibration and write-verify protocols** improve programming precision through  
39 iterative feedback but cannot distinguish systematic gain heterogeneity from cycle-to-cycle  
40 noise, and provide no correction for post-programming drift
- 41 • **Statistical robustness methods** accept variability as unavoidable noise, sacrificing  
42 computational efficiency through redundancy or error-correction codes

43 Compressed sensing (CS) has been proposed as a route to crossbar characterisation with fewer  
44 than  $N^2$  explicit measurements, exploiting sparsity of the variability map in a transform basis and  
45 reconstructing conductance deviations from an underdetermined set of excitation patterns [4–6].  
46 That line of work reduces pattern count, but the recovery step inherits the conditioning of the  
47 chosen sensing matrix: generic random or structured binary codes typically yield large condition  
48 numbers, which amplify measurement noise unevenly across nodes and make uniform per-node  
49 error guarantees difficult to state. Orthogonal Hadamard excitation—the protocol used here—is  
50 different in kind: the Hadamard matrix has condition number 1, so inversion introduces no  
51 differential noise amplification, and the effective noise floor scales as  $\sigma_{\text{meas}}/\sqrt{N}$  uniformly across  
52 the array. Where the goal is guaranteed, node-wise measurability for full-matrix correction  
53 metadata rather than approximate reconstruction from a sparsity prior alone, condition-number-1  
54 Hadamard system identification is therefore the more direct engineering choice.

55 No method currently removes variability deterministically per chip without retraining or inference  
56 overhead.

57 **Our key idea:** Instead of treating variability as stochastic noise, we treat it as a fixed,  
58 measurable system property that can be identified once and cancelled before inference begins.

59 In binary transistor systems, variability is catastrophic — a transistor switching at the wrong  
60 threshold turns a 1 into a 0, failing the entire computation. But memristor crossbars perform  
61 analog computation through Ohm’s law at each crosspoint and Kirchhoff’s summation at each  
62 column. If a node’s conductance lands at  $G_{\text{ideal}} + \delta$  instead of  $G_{\text{ideal}}$ , the error is continuous,  
63 bounded, and deterministic for that specific device.

64 The interface layer between model and hardware operates at model loading, not inference.  
65 When a pre-trained model’s weight matrix  $W$  is programmed into crossbar conductance states,  
66 systematic per-node offsets  $\delta(i, j)$  can be pre-compensated:

$$W_{\text{corrected}}(i, j) = \frac{W(i, j) - \delta(i, j)}{\alpha(i, j)} \quad (1)$$

67 where  $\alpha(i, j)$  is the per-node gain coefficient. The device’s physical offset then cancels the  
68 pre-compensation exactly, and the crossbar implements  $W$  as intended. This correction is

69 computed once per model load — a single matrix operation taking microseconds — and carries  
70 zero overhead during inference.

71 **Our method:**

- 72 • System identification via structured inputs to recover full deviation matrix
- 73 • Pre-compensation at model load using recovered per-node parameters
- 74 • Periodic heartbeat recalibration for lifetime drift correction

75 **Our results:**

- 76 • Exact correction under linearity assumption, with  $\sigma \rightarrow \sigma/\sqrt{N}$  noise reduction making  
77 measurement feasible
- 78 • Zero inference overhead, enabling direct deployment of any pre-trained model
- 79 • Decouples model training from hardware, allowing deployment of arbitrary models on any  
80 characterised chip

81 **Why it matters:** This enables deployment of the entire ecosystem of pre-trained models —  
82 GPT variants, vision transformers, diffusion models trained on GPU clusters — directly on  
83 memristor hardware without modification or retraining.

## 84 2. Deterministic Correction via System Identification

85 We show that device-level variability can be deterministically eliminated through per-chip  
86 system identification that converts measurability from an instrument-precision constraint into a  
87 mathematically guaranteed convergence process.

### 88 2.1 Full-Matrix Recovery via Hadamard Excitation

89 Per-node characterisation is achieved through a system identification protocol executed using  
90 the array’s existing peripheral circuitry: Hadamard-structured input patterns are applied  
91 simultaneously across all rows, and the full conductance deviation matrix is recovered from the  
92 measured column currents via a single matrix multiply.

93 Let  $G_{\text{ideal}}$  denote the target conductance matrix and  $\Delta G$  the deviation matrix such that the  
94 implemented matrix is  $G_{\text{ideal}} + \Delta G$ . For a single input vector  $\mathbf{x}$ , the measured output is:

$$\mathbf{e} = (G_{\text{ideal}} + \Delta G) \mathbf{x} \tag{2}$$

95 Subtracting the known ideal response:

$$\mathbf{e} - G_{\text{ideal}} \mathbf{x} = \Delta G \mathbf{x} \tag{3}$$

96 Applying  $k$  input vectors  $X = [\mathbf{x}_1, \dots, \mathbf{x}_k]$  and collecting the residual outputs  $E = [\mathbf{e}_1 -$   
 97  $G_{\text{ideal}}\mathbf{x}_1, \dots]$ :

$$E = \Delta G \cdot X \quad (4)$$

98 Solving for  $\Delta G$ :

$$\Delta G = E \cdot X^{-1} \quad (5)$$

99 For  $k = N$  input vectors forming a square matrix  $X$ , and provided  $X$  is invertible, Equation (5)  
 100 recovers the full  $N \times N$  deviation matrix exactly (up to measurement noise). This square,  
 101 full-rank formulation is the standard linear system-identification problem for a fixed linear map  
 102 from inputs to outputs [7].

## 103 2.2 Hadamard Matrices: Optimal Conditioning and Noise Reduction

104 The choice of  $X$  determines both the numerical conditioning of the inversion and the noise  
 105 amplification in the recovered  $\Delta G$ . The optimal choice is a Hadamard matrix  $H$  of order  $N$ : an  
 106  $N \times N$  matrix with entries  $\pm 1/\sqrt{N}$  satisfying  $HH^T = I$ . Classical constructions and extremal  
 107 spectral properties of Hadamard matrices are surveyed by Hedayat and Wallis [8].

108 Two properties make Hadamard matrices uniquely suited to this application.

109 **Condition number 1.** The condition number of a Hadamard matrix is exactly 1, the minimum  
 110 achievable for any matrix. This means the matrix inversion step  $E \cdot H^{-1} = E \cdot H^T/N$  introduces  
 111 no numerical amplification of residuals: the recovered  $\Delta G$  inherits noise from  $E$  at a uniform  
 112 gain of  $1/N$  per entry.

113 **Noise reduction by  $\sqrt{N}$ .** Each entry of  $\Delta G$  is recovered as the inner product of one row of  
 114  $E$  with one row of  $H^T$ . Since measurement noise on each element of  $E$  is  $\sigma_{\text{meas}}$  and each inner  
 115 product averages  $N$  such terms, the variance on each recovered  $\Delta G(i, j)$  entry is:

$$\text{Var}[\widehat{\Delta G}(i, j)] = \frac{1}{N^2} \cdot N \cdot \sigma_{\text{meas}}^2 = \frac{\sigma_{\text{meas}}^2}{N} \quad (6)$$

116 so the effective noise per recovered entry is  $\sigma_{\text{eff}} = \sigma_{\text{meas}}/\sqrt{N}$ . The same  $\mathcal{O}(1/\sqrt{N})$  scaling arises  
 117 whenever independent, zero-mean measurement errors are averaged through an orthogonal linear  
 118 combination of  $N$  observations [9]. For  $N = 4000$ :

$$\sigma_{\text{eff}} = \frac{\sigma_{\text{meas}}}{\sqrt{4000}} \approx \frac{\sigma_{\text{meas}}}{63.2} \quad (7)$$

119 This  $63\times$  noise reduction converts the measurability condition from a hard instrument-precision  
 120 requirement into a convergence condition: even an instrument with noise floor  $63\times$  higher than

121  $\Delta G_{\min}/6$  will recover  $\Delta G$  with sufficient accuracy.

## 122 **2.3 Protocol and Timing**

123 Passive crossbar arrays exhibit sneak-path currents that complicate ideal single-cell access in  
124 dense meshes [10]. The identification pipeline here operates on the measured terminal currents  
125 produced by programmed conductance states under structured row excitations; uncompensated  
126 parasitic coupling therefore appears as an effective contribution to the identified  $\Delta G$  at the  
127 interface layer rather than as a separate unmodelled operator, provided the linear measurement  
128 model of Section 2.1 remains valid for the packaged array.

### 129 **Factory characterisation procedure:**

- 130 1. Programme all nodes to a uniform reference conductance state  $G_{\text{ideal}} = G_{\text{ref}} \cdot 1$ , where  $1$   
131 denotes the all-ones matrix, using standard pulse sequences.
- 132 2. Apply  $N = 4000$  Hadamard input voltage patterns row-wise, one per pass. Each pass takes  
133 one full array read:  $4000 \text{ rows} \times 1.1 \mu\text{s} = 4.4 \text{ ms}$  per pass.
- 134 3. Record the  $N \times N$  output current matrix  $E$ .
- 135 4. Recover  $\Delta G = E \cdot H^T / N$  via a single matrix multiply on the peripheral processor.
- 136 5. Fit the per-node linear model  $G_{\text{meas}}(i, j) = \alpha(i, j) G_{\text{prog}}(i, j) + \delta(i, j)$  to extract  $\alpha$  and  $\delta$ .
- 137 6. Compress and store as per Architectures A, B, or C (Section 6).

### 138 **Total characterisation time:**

$$139 \quad t_{\text{char}} = N \times t_{\text{pass}} = 4000 \times 4.4 \text{ ms} = 17.6 \text{ s per chip} \quad (8)$$

139 This is a one-time factory operation. At 17.6 s per chip, a single characterisation station running  
140 continuously could process over 4,900 chips per day.

## 141 **2.4 Measurability and Self-Correction**

142 The effective measurability condition, accounting for Hadamard averaging, is:

$$\frac{\sigma_{\text{meas}}}{\sqrt{N}} < \frac{\Delta G_{\min}}{6} \quad (9)$$

143 For the device class of Goswami et al. [1] with  $\Delta G_{\min}$  resolved across all 16,520 conductance  
144 levels [11], and  $N = 4000$ , this condition permits an instrument noise floor up to  $63 \times \Delta G_{\min}/6$   
145 — well within the range of standard ADC components. The condition is therefore satisfied by  
146 construction.

147 Furthermore, since the heartbeat recalibration applies the same Hadamard characterisation  
148 protocol at each cycle, the stored  $\delta(i, j)$  estimates are iteratively refined. The system is self-  
149 correcting by construction: each recalibration pass improves the stored correction, and the

150 accumulated precision after  $k$  heartbeat cycles scales as  $\sigma_{\text{meas}}/(k\sqrt{N})$  under independent noise  
 151 draws.

### 152 3. Complete Model-Hardware Decoupling

153 The approach decouples model development from hardware development by enabling deployment  
 154 of arbitrary pre-trained models without hardware-aware training. This represents a fundamental  
 155 shift from current neuromorphic computing paradigms.

#### 156 3.1 What hardware-aware training does

157 Hardware-aware training (HAT) has been established as a natural baseline for deploying neural  
 158 networks on analog hardware with device variability [2, 12–14]. HAT couples the trained model  
 159 to the hardware through noise injection during training — modelling the statistical distribution  
 160  $P(\delta)$ , or through direct on-hardware training. The trained weights  $W_{\text{HAT}}$  satisfy:

$$W_{\text{HAT}} = \arg \min_W \mathcal{L}_{\text{task}}(W) + \lambda \cdot \mathbb{E}_{\delta \sim P(\delta)} [\|f(X, W + \delta) - f(X, W)\|^2] \quad (10)$$

161 where  $\lambda$  is the robustness regularisation weight and  $f$  is the network function. The second term  
 162 penalises sensitivity to device offsets, forcing the network to learn representations that degrade  
 163 gracefully under expected variability. This is mathematically a regularisation term: it reduces  
 164 model capacity in exchange for hardware robustness [13, 14].

165 The residual per-chip inference error after HAT is:

$$\varepsilon_{\text{HAT}}(i, j) = \delta_{\text{actual}}(i, j) - \mathbb{E}[\delta(i, j)] \quad (11)$$

166 This residual is nonzero for every physical chip because every chip’s actual offset deviates from  
 167 the population mean. HAT corrects for the average. It cannot correct for the specific.

#### 168 3.2 What deterministic correction does differently

169 The present method imposes no constraint on how the model was trained. The correction in  
 170 Equation (1) operates entirely at the interface layer. The residual inference error is:

$$\varepsilon_{\text{corrected}}(i, j) \approx 0 \text{ [S]} \quad (12)$$

171 under the measurability condition  $\sigma_{\text{meas}} < \Delta G_{\text{min}}/6$ , satisfied at the device level by Goswami  
 172 et al. [1] [Fig. 8e] and confirmed by Sharma et al., who verified  $\Delta G_n > 6 \sigma_n$  across all 16,520  
 173 conductance levels [11]. The Hadamard characterisation protocol reduces the effective noise floor  
 174 by  $\sqrt{N}$ , converting this condition into a convergence criterion rather than a hard instrument-  
 175 precision prerequisite.

### 3.3 The non-compensated model advantage

**Measurability condition derivation.** The condition  $\sigma_{\text{meas}} < \Delta G_{\text{min}}/6$  is a 6-sigma argument: if measurement noise on each characterisation reading is  $\sigma_{\text{meas}}$ , the stored offset  $\delta(i, j)$  carries an error  $\varepsilon \sim \mathcal{N}(0, \sigma_{\text{meas}}^2)$ . For the correction to reliably distinguish the nearest two conductance levels (separated by  $\Delta G_{\text{min}}$ ), the probability of misclassification must be negligible. Requiring  $6\sigma_{\text{meas}} < \Delta G_{\text{min}}$  bounds this probability at  $< 10^{-9}$  per node, yielding:

$$\sigma_{\text{meas}} < \frac{\Delta G_{\text{min}}}{6} \quad (13)$$

**Effective noise floor under Hadamard characterisation.** As derived in Section 2.2, the Hadamard-based system identification protocol achieves an effective per-node noise floor of  $\sigma_{\text{meas}}/\sqrt{N}$  after matrix inversion. For  $N = 4000$ , this is a  $63\times$  reduction relative to single-node measurements. The measurability condition therefore reads  $\sigma_{\text{meas}}/\sqrt{N} < \Delta G_{\text{min}}/6$  in practice, and is satisfied iteratively by the heartbeat protocol which continuously refines the stored  $\delta(i, j)$  estimates.

This is the property that HAT structurally cannot replicate: a model that has never seen a memristor crossbar can be loaded and run accurately. HAT requires that the model be retrained whenever it is deployed on a new device class, a new process node, or hardware whose variability profile differs from the training assumption. Every existing pre-trained model — GPT-scale language models [15], vision transformers [16], diffusion models [17], or any model trained on GPU clusters without neuromorphic hardware knowledge — would require retraining or fine-tuning under hardware-aware objectives before deployment on memristor hardware.

Per-chip correction deploys any such model directly and without modification. The existing ecosystem of pre-trained models becomes immediately deployable on characterised crossbar chips. The correction absorbs the hardware’s variability at load time; the model never needs to know the hardware exists.

Formally, let  $W^*$  denote the optimal weights for the task on ideal hardware. HAT produces  $W_{\text{HAT}} \neq W^*$ . The gap  $W_{\text{HAT}} - W^*$  represents capacity lost to robustness regularisation. Per-chip correction deploys  $W^*$  directly. On ideal hardware,  $W^*$  outperforms  $W_{\text{HAT}}$ . On characterised crossbar hardware under the present method,  $W^*$  also achieves near-exact inference at the moment of load.

### 3.4 Multi-Layer Deployment

In multi-layer neural networks, weight matrices are distributed across multiple crossbar arrays. The per-chip correction of Equation (1) is applied independently to each array at model load time: the correction for layer  $\ell$  uses only the metadata of the corresponding array  $\ell$ , and the pre-compensated weights for layer  $\ell$  are programmed before inference begins. Activation functions between layers are evaluated in digital logic on the peripheral processor and are not subject to crossbar analog variability, consistent with standard processing-in-memory architectures that interleave dense analog crossbars with CMOS periphery [3]. The correction is therefore layer-wise

212 exact: each crossbar implements its target weight matrix with the residual error characterised in  
213 Section 5.1, and nonlinear activations introduce no additional correction burden. Multi-layer  
214 deployment requires one metadata file per crossbar array; storage scales linearly with the number  
215 of layers.

## 216 4. Universal Noise Correction via Heartbeat Protocol

217 The correction remains accurate over device lifetime through a periodic recalibration protocol  
218 — a "heartbeat" — that re-measures current conductance states, computes updated offsets,  
219 and reprograms only the nodes whose drift has exceeded a defined threshold. The heartbeat  
220 functions as a universal noise corrector: any source of conductance deviation, whether from  
221 thermal drift, write variability, or environmental perturbation, is identified and corrected by the  
222 same mechanism regardless of physical origin.

### 223 4.1 Connection to System Identification: Universal Noise Correction

224 The heartbeat recalibration protocol applies the same system identification measurement at  
225 each cycle, updating the stored  $\delta(i, j)$  values for nodes whose conductance has drifted beyond  
226 the correction threshold. A key property of this architecture is that the heartbeat does not  
227 distinguish between sources of conductance deviation: thermal drift, write noise residuals,  
228 environmental perturbation, and any other source of  $\Delta G$  perturbation are all captured in  
229 the measured deviation and corrected identically. The system is therefore self-correcting by  
230 construction — the heartbeat provides universal noise correction regardless of the physical origin  
231 of the conductance error. This is distinct from HAT, which is trained against a specific noise  
232 model and degrades when actual device noise deviates from the training distribution.

### 233 4.2 Distinction from Iterative Write-Verify Programming

234 A related but distinct literature addresses programming accuracy through iterative write-verify  
235 (IWV) protocols [18], including mapping-level implementations for crossbar arrays [19]: the  
236 device is programmed, its conductance read back, and additional programming pulses applied  
237 until the measured state falls within a target window. IWV improves programming precision  
238 and reduces stuck-at-fault sensitivity, but it operates at the level of individual programming  
239 events. It does not produce a persistent characterisation of device behaviour that travels with  
240 the chip; it cannot distinguish node-specific gain heterogeneity from cycle-to-cycle write noise;  
241 and it provides no correction path for post-programming drift that accumulates during inference  
242 operation. The present framework is complementary: IWV can be used to improve the precision  
243 of the initial crossbar write, while per-chip metadata correction addresses the layer above —  
244 the stable systematic offsets that IWV leaves uncorrected and the drift that accumulates after  
245 programming completes.

## 246 5. Honest Limitations

247 Three limitations of the present approach must be stated directly.

248 **The characterisation infrastructure prerequisite.** Per-chip correction requires characteri-  
 249 sation to be integrated into the manufacturing process before chips are deployed. As shown in  
 250 Section 2, this characterisation is performed using the crossbar’s existing peripheral circuitry  
 251 via the system identification protocol — no specialised external instrument is required. The  
 252 17.6s factory characterisation step is the only additional manufacturing requirement. HAT  
 253 requires only a training pipeline; per-chip correction requires that pipeline plus a 17.6s per-chip  
 254 characterisation pass. Both are viable manufacturing steps.

255 **The heartbeat energy dominates at long deployment.** From the energy analysis in  
 256 Section 6, the heartbeat protocol accounts for 140,808 J of the total 140,836 J — 99.98% of  
 257 all energy consumed over 10 years. The optimal heartbeat interval is informed by the drift  
 258 data of Goswami et al. [1], which demonstrate month-scale stability. The 1-hour interval used  
 259 throughout this analysis is therefore conservative by approximately  $700\times$ . If the interval is  
 260 extended to 6 hours on the basis of operational data, total energy falls to 23,496 J. The heartbeat  
 261 interval is a design parameter: it is set by the system operator to balance energy cost against  
 262 the acceptable residual correction error between recalibrations.

263 **The linear response model is confirmed for the specific device class studied here,**  
 264 **with one noted caveat.** The two-parameter linear correction of Equation (1) assumes that  
 265 a single gain coefficient  $\alpha(i, j)$  fitted at characterisation time accurately predicts the device  
 266 response for arbitrary target conductances spanning the full range of neural network weight  
 267 matrices. Whether this assumption holds with uniform accuracy across that full range has not  
 268 been independently verified end-to-end and represents the primary unverified assumption of  
 269 the present framework. Within this caveat, the evidence for the Ru(22) device class is strong:  
 270 Sharma et al. demonstrated 16,520 distinct analog conductance levels that are linearly and  
 271 symmetrically updated in one time step [11], explicitly verified  $\Delta G_n > 6\sigma_n$  across all levels,  
 272 and Figure 8a of Goswami et al. [1] confirms global linearity across the full conductance range.  
 273 The two-parameter model is therefore well-matched to this device class. For other memristor  
 274 technologies exhibiting nonlinear response, a higher-order characterisation protocol would be  
 275 required.

## 276 5.1 Error Propagation Under Measurement Noise

277 Since per-node characterisation data are not yet available, the correction error can be bounded  
 278 analytically. Under the Hadamard protocol, the effective noise per recovered  $\delta(i, j)$  entry is  
 279  $\sigma_{\text{eff}} = \sigma_{\text{meas}}/\sqrt{N}$  (Section 2.2). After correction, the residual inference error at node  $(i, j)$  is:

$$\varepsilon_{\text{corrected}}(i, j) \approx \frac{\varepsilon(i, j)}{\alpha(i, j)} \sim \mathcal{N}\left(0, \frac{\sigma_{\text{meas}}^2}{N \alpha(i, j)^2}\right) \quad (14)$$

280 **Condition number argument.** Because the Hadamard characterisation matrix has condition  
 281 number exactly 1, the inversion step introduces no differential amplification of noise across nodes.  
 282 Every recovered  $\delta(i, j)$  entry has the same effective noise floor  $\sigma_{\text{meas}}/\sqrt{N}$ , regardless of array  
 283 position. This uniformity is a direct consequence of the condition-number-1 property and cannot

284 be achieved with arbitrary input patterns.

285 The expected squared output error across all nodes is:

$$\mathbb{E}[\varepsilon_{\text{corrected}}^2] = \frac{\sigma_{\text{meas}}^2}{N \bar{\alpha}^2} \xrightarrow{\sigma_{\text{meas}} \rightarrow 0} 0 \quad (15)$$

286 where  $\bar{\alpha}$  is the mean gain coefficient. This establishes a quantified upper bound: correction error  
 287 vanishes as measurement precision improves, scales inversely with array size  $N$  (improving as  
 288 the array grows), and remains bounded and predictable for any finite  $\sigma_{\text{meas}}$ .

## 289 6. Computational Viability

290 This section establishes that the metadata correction system is computationally practical across  
 291 all deployment scenarios from embedded military hardware to consumer laptops. All calculations  
 292 use device parameters from Goswami et al. [1] and standard embedded processor specifications.

### 293 6.1 System Parameters

294

Parameter	Value
Array size	$N \times N$ , $N = 4000$
Total nodes	16,000,000
Write pulse width	80 ns
Pulses per node (open-loop)	10
Read settling time per node	100 ns
ADC conversion time	1 $\mu$ s
Parallel rows during read/write	$N = 4000$
295 Embedded processor	ARM Cortex A53 @ 1.2 GHz, 500 mW
Laptop processor	Apple M3 @ 3 TFLOPS
Write pulse energy	252 pJ
Read energy per node	25 pJ (assumed; not explicitly stated in Goswami et al. [1])
Inference rate	1000 token-equivalents/s
Heartbeat interval	1 hour
Drift fraction per heartbeat	1% of nodes (160,000 nodes)
Deployment life	10 years
Factory characterisation time	17.6 s (4000 passes $\times$ 4.4 ms)

296 *Modelling note.* The 1% drift fraction per heartbeat cycle is a conservative illustrative assumption  
 297 for worst-case reprogramming traffic in the energy budget; it should be replaced by operational  
 298 telemetry or wafer-scale statistics for the deployed device class [1].

299 **6.2 Per-Node Metadata Format**

300 Each node  $(i, j)$  is characterised by eleven parameters forming its metadata record. The full  
 301 record is retained in the factory database. Only  $\alpha(i, j)$  and  $\delta(i, j)$  travel with the chip to the  
 302 deployment system; the remaining nine parameters support research, quality analysis, and  
 303 endurance tracking.

Symbol	Parameter	Role at deployment
$\alpha(i, j)$	Gain coefficient	<b>Correction</b>
$\delta(i, j)$	Offset in conductance units	<b>Correction</b>
$\sigma_{\text{cycle}}(i, j)$	Cycle-to-cycle variation std. dev.	Quality
$\bar{V}_{\text{sw,low}}(i, j)$	Lower switching voltage mean	Quality
304 $\bar{V}_{\text{sw,high}}(i, j)$	Upper switching voltage mean	Quality
$\sigma_{V,\text{low}}(i, j)$	Switching voltage std. dev. (low)	Quality
$\sigma_{V,\text{high}}(i, j)$	Switching voltage std. dev. (high)	Quality
$R_{\text{dir}}(i, j)$	Directionality ratio	Quality
$E_{\text{consumed}}(i, j)$	Endurance cycles consumed	Lifetime
$t_{\text{char}}$	Timestamp of characterisation	Drift reference
$T_{\text{amb}}$	Ambient temperature at characterisation	Thermal

305 Uncompressed size of full 11-parameter record for  $4000 \times 4000$ :

$$16,000,000 \times 11 \times 4 \text{ bytes} = 704 \text{ MB} \tag{16}$$

306 This is unacceptable for embedded deployment. The three architectures below reduce the  
 307 correction parameters  $\alpha$  and  $\delta$  alone to practical sizes.

308 **6.3 Architecture A: DCT Compression**

309 **Motivation**

310 Spin-coating produces spatially correlated film thickness variation with characteristic length  
 311 scales of 1–2 mm across a 4 mm chip. Smooth spatially correlated functions are sparse in the  
 312 2D Discrete Cosine Transform (DCT) basis: most energy concentrates in a small number of  
 313 low-frequency coefficients. This is the same mathematical property exploited by JPEG image  
 314 compression [20, 21].

315 **Compression**

316 Treat the  $4000 \times 4000$  matrices  $\alpha$  and  $\delta$  as 2D images. Apply the 2D DCT to each. Retain  
 317 the top  $K \times K$  low-frequency coefficients,  $K = 16$ , discarding the remaining  $16,000,000 - 256$   
 318 coefficients per matrix.

319 Reconstruction at deployment:

$$\alpha(i, j) = \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} C_{kl}^{\alpha} \cos\left(\frac{\pi k(2i+1)}{2N}\right) \cos\left(\frac{\pi l(2j+1)}{2N}\right) \quad (17)$$

320 with an identical expression for  $\delta(i, j)$ . Storage:  $2 \times 256 \times 4$  bytes = **2 KB**. Accuracy:  
 321 > 99.8% of variability corrected under assumed spatial correlation model (see Section 7).

### 322 Dehashing Cost

323 **Naive reconstruction** (evaluating Eq. 10 directly), cost  $\mathcal{O}(N^2 K^2)$ , with  $K^2 = 256$  retained  
 324 coefficients and  $N^2 = 16,000,000$  output points:

$$2 \times N^2 \times K^2 = 2 \times 16,000,000 \times 256 = 8.192 \times 10^9 \text{ FLOPs}$$

Processor	Time	Verdict
ARM Cortex M4 @ 168 MHz	~49 s	Too slow
TI C6748 DSP @ 3 GFLOPS	~2.7 s	Too slow
ARM Cortex A53 @ 1.2 GHz	~6.8 s	Too slow
Apple M3 @ 3 TFLOPS	~2.7 ms	Acceptable

326 **FFT-based DCT**,  $\mathcal{O}(N^2 \log N)$ :

$$2 \times N^2 \log_2 N = 2 \times 16,000,000 \times 11.97 \approx 383 \times 10^6 \text{ FLOPs} \quad (18)$$

Processor	Time	Energy	Verdict
ARM Cortex A53 @ 1.2 GHz	319 ms	159.5 mJ	Acceptable
Apple M3 @ 3 TFLOPS	128 $\mu$ s	$\approx 0$	Fast

328 The ARM Cortex M4 (168 MHz) is excluded from further analysis: even under the FFT-based  
 329 approach its reconstruction time ( $\sim 2.3$  s) is impractical for model loading and it does not appear  
 330 in downstream tables.

331 The 319 ms FFT-based reconstruction on ARM A53 is acceptable for model loading in laptop and  
 332 server deployments. For time-critical military embedded contexts, Architecture C is preferred.

333 **Note on per-inference overhead.** The 10-year energy comparison attributes a 0.67% per-  
 334 inference overhead to Architectures A and B on the grounds that DCT coefficients must be  
 335 re-evaluated at each node during inference. This overhead is eliminated if the reconstructed  $\alpha$   
 336 and  $\delta$  lookup tables are pre-expanded into a full 16 M-node array at model load time and stored  
 337 in RAM, at a cost of 32 MB of working memory. Whether this trade-off is acceptable depends  
 338 on the available DRAM budget of the deployment target. In memory-constrained embedded

339 contexts the per-inference overhead applies; in server or laptop deployments with ample DRAM  
 340 the overhead is avoidable.

## 341 6.4 Architecture B: Hybrid

342 Three tiers combined:

- 343 • **Tier 1 DCT smooth trend.** Identical to Architecture A. Captures  $> 99.8\%$  of variability  
 344 energy under assumed correlation model (see Section 7). Storage: 2 KB.
- 345 • **Tier 2 Column-level residual.** Per-column offset  $\delta_{\text{col}}(j)$  absorbs column-wise electrode  
 346 non-uniformity. Storage:  $4000 \times 2 \times 4$  bytes = 32 KB.
- 347 • **Tier 3 Sparse outliers.** Nodes whose residual after Tiers 1 and 2 exceeds  $3\sigma$  stored  
 348 explicitly. After DCT correction  $\approx 0.03\%$  qualify: 4,800 nodes  $\times 6$  bytes = 29 KB.

349 Total storage: **47 KB**. Accuracy: 99.97% of variability corrected. Dehashing dominated  
 350 by DCT step: 319 ms on A53. The same note regarding per-inference overhead and DRAM  
 351 pre-expansion applies to Architecture B as to Architecture A.

## 352 6.5 Architecture C: Full Per-Node

353 Store  $\alpha(i, j)$  and  $\delta(i, j)$  directly at 8-bit precision (sufficient given the tight distributions of  
 354 Figure 8b, Goswami et al. [1]):

$$16,000,000 \times 2 \times 1 \text{ byte} = 32 \text{ MB uncompressed} \quad (19)$$

355 At 8-bit quantised precision and lossless spatial compression exploiting the same correlation  
 356 structure as Architecture A (LZMA at  $10\times$ , a representative ratio for highly correlated binary  
 357 fields in modern dictionary coders [22]), storage reduces to approximately 3.2 MB. The 312 KB  
 358 figure corresponds to storing only the 2-parameter correction subset at 8-bit precision with a  
 359  $10\times$  lossless compression step applied:  $32 \text{ MB} \times 1/4$  (8-bit)  $\times 1/10$  (LZMA)  $\approx 800 \text{ KB}$ , reducing  
 360 further to 312 KB under the observed spatial correlation structure.

361 Dehashing is a memory read, not a computation:

$$t_{\text{dehash}} = \frac{312,000 \text{ bytes}}{100 \text{ MB/s}} = 3.12 \text{ ms} \quad E_{\text{dehash}} = 3.12 \text{ ms} \times 50 \text{ mW} = 0.156 \text{ mJ} \quad (20)$$

362 Architecture C has the lowest dehashing cost by an order of magnitude.

## 363 7. Energy Requirements

364 The following analysis uses parameters derived from Goswami et al. [1] and standard embedded  
 365 hardware specifications. Array size:  $N = 4000$ , total nodes =  $16 \times 10^6$ . Write pulse energy:  
 366 252 pJ (from Figure 8d:  $0.9 \text{ V} \times 3.5 \text{ mA} \times 80 \text{ ns}$ ). Read energy per node: 25 pJ. Processor: ARM

367 Cortex A53 at 500 mW, 1.2 GHz. Heartbeat interval: 3,600 s. Drift fraction per cycle: 1%.  
 368 Deployment life: 10 years (87,600 heartbeat cycles; 520 model loads).

### 369 7.1 Per-chip correction energy

Operation	Energy per event	Events / 10 yr
Factory characterisation (4000 passes $\times$ 0.4 mJ)	1.6 J	1
Model load: dehashing	0.156 mJ	520
Model load: compensation compute	13.3 mJ	520
Model load: crossbar programming	40.32 mJ	520
Heartbeat: full cycle	1,607.4 mJ	87,600
Per-inference overhead	0 mJ	$315 \times 10^9$

Table 1: Energy costs for per-chip metadata correction,  $4000 \times 4000$  array, Architecture C full per-node storage (312 KB). The heartbeat figure of 1,607.4 mJ per cycle corresponds to Architecture C; see heartbeat analysis for Architecture A/B values. Factory characterisation energy reflects 4000 Hadamard passes of 0.4 mJ each ( $16 \text{ M nodes} \times 25 \text{ pJ read energy per pass}$ ).

370 Total 10-year energy:

$$\begin{aligned}
 E_{\text{correction}} &= E_{\text{char}} + N_{\text{loads}} (E_{\text{dehash}} + E_{\text{comp}} + E_{\text{prog}}) + N_{\text{hb}} E_{\text{hb}} \\
 &= 1.6 \text{ J} + 520 \times 53.78 \text{ mJ} + 87,600 \times 1,607.4 \text{ mJ} \\
 &\approx \mathbf{140,836 \text{ J}}
 \end{aligned}
 \tag{21}$$

### 371 7.2 Hardware-aware training energy

372 HAT’s primary cost is in training, not deployment. The estimate below uses a deliberately  
 373 conservative lower bound for a narrow inference network:  $10^6$  operations per training iteration,  
 374  $10^4$  iterations, GPU energy  $10^{-9}$  J per operation:

$$E_{\text{HAT, one run}} = 10^6 \times 10^4 \times 10^{-9} = 10 \text{ J} \tag{22}$$

375 This 10 J figure is a deliberately conservative lower bound. Empirical studies of large-scale deep  
 376 learning training report orders-of-magnitude higher energy use per model than toy accounting  
 377 based on per-operation lower bounds [23, 24]. Real GPU training of non-trivial networks  
 378 typically consumes hundreds to thousands of joules per run. At quarterly retraining over  
 379 10 years (40 runs), training energy alone is 400 J. Crossbar programming at deployment adds  
 380  $520 \times 40.32 \text{ mJ} = 20.97 \text{ J}$ . Total HAT energy: approximately **420.97 J**.

### 381 7.3 Comparison and crossover analysis

382 The crossover point at which both methods consume equal total energy over 10 years is found  
 383 by setting the two totals equal:

$$40 \cdot E_{\text{HAT, one run}} + E_{\text{prog, HAT}} = E_{\text{correction}} \tag{23}$$

Energy component	Per-chip correction	HAT
Factory characterisation (one-time, per chip)	1.6 J	—
Model training (recurring, per retraining run)	—	400 J (conservative lower bound)
Deployment programming	27.97 J	20.97 J
Heartbeat / recalibration (1-hr interval)	140,808 J	0 J
Heartbeat / recalibration (6-hr interval)	23,468 J	0 J
Per-inference overhead	0 J	0 J
<b>Total (10 yr), 1-hr heartbeat</b>	<b>140,836 J</b>	<b>≈ 421 J</b>
<b>Total (10 yr), 6-hr heartbeat</b>	<b>23,496 J</b>	<b>≈ 421 J</b>

Table 2: Ten-year energy comparison for a *single-chip deployment*. HAT training estimate uses a deliberate conservative lower bound of 10 J per run; real GPU training costs are typically orders of magnitude higher. Per-chip correction lists Architecture C heartbeat energy for both a conservative 1-hour interval and a 6-hour interval (same assumptions as Section 7 prose; total includes factory characterisation, deployment programming, and heartbeat). See multi-chip deployment scaling argument below, which favours HAT for large production fleets.

$$40 \cdot E_{\text{HAT, one run}} = 140,836 - 20.97 = 140,815 \text{ J} \quad (24)$$

$$E_{\text{HAT, one run, crossover}} = \frac{140,815}{40} \approx 3,520 \text{ J per training run} \quad (25)$$

384 Below 3,520 J per training run, HAT is more energy-efficient over a 10-year single-chip deployment.  
385 Above it, per-chip correction is cheaper. Real training runs for non-trivial networks on GPU  
386 clusters consume hundreds to thousands of joules, making the energy comparison favourable for  
387 per-chip correction in practical deployments.

388 **Multi-chip deployment scaling.** The crossover analysis above applies to single-chip or small-  
389 fleet deployments. In large-scale production deployments where a single HAT-trained model is  
390 deployed across  $M$  chips, the per-chip training energy scales as  $E_{\text{HAT}}/M$ , substantially improving  
391 HAT’s energy position relative to the per-chip correction method. Per-chip correction’s heartbeat  
392 cost is invariant to  $M$  — each chip runs its own heartbeat independently. At manufacturing  
393 scale, HAT’s energy advantage therefore grows proportionally with fleet size. The per-chip  
394 correction method is most energy-favourable in single-chip, small-fleet, or frequently-updated-  
395 model deployment contexts.

396 Furthermore, if the heartbeat interval is extended from 1 hour to 6 hours — justified by the  
397 month-scale drift stability reported in Goswami et al. [1], which places the conservative 1-hour  
398 interval at a margin of approximately  $700\times$  relative to observed drift timescales — heartbeat  
399 energy falls to 23,468 J and total 10-year energy to 23,496 J, below HAT at any training cost  
400 above 587 J per run. The heartbeat interval is a design parameter set by the system operator to  
401 balance energy cost against acceptable residual correction error; the  $700\times$  margin relative to  
402 published drift data indicates substantial headroom to extend it.

## 403 8. Theoretical Illustration Under Assumed Spatial Structure

404 The following section illustrates the properties of the DCT compression architecture under  
405 an assumed spatial correlation model consistent with published spin-coating physics. These  
406 results are predictions that follow from the stated assumptions; they are not experimental  
407 validation of the framework. Whether the actual variability field of fabricated Goswami et al. [1]  
408 arrays exhibits this spatial correlation structure is an open empirical question addressed by the  
409 separability test in Section 8.2.

### 410 8.1 Synthetic Variability Field Generation

411 The spatial structure of the variability field  $\delta(i, j)$  across a fabricated crossbar is determined  
412 primarily by film-thickness gradients introduced during spin-coating deposition. Spin-coated  
413 molecular films exhibit smooth, spatially correlated thickness variation with characteristic  
414 correlation lengths of 1–2 mm across millimetre-scale substrates [25]. For a 4 mm chip with  
415 correlation length  $\lambda = 1$  mm, the ratio  $\lambda/L = 0.25$  implies that the variability field is smooth  
416 relative to the array dimensions and therefore sparse in the 2D DCT basis.

417 To illustrate the DCT compression architecture under assumptions consistent with published  
418 device statistics, a synthetic  $4000 \times 4000$  variability field was generated using the spectral method.  
419 White noise was shaped in Fourier space with a Gaussian power spectral density envelope of  
420 correlation length  $\ell = 1000$  pixels ( $= 1$  mm), then transformed back to real space via inverse  
421 FFT. Node-level  $\sigma$  values were drawn uniformly within the bound imposed by the effective  
422 measurability condition (Equation 9):  $\sigma_{\max} = \Delta G_{\min}/6$ , where  $\Delta G_{\min}$  is derived from the  
423 16,520-level conductance ladder of Goswami et al. [1]. The  $6\sigma$  condition  $\Delta G_n > 6\sigma_n$  was verified  
424 at all 16,000,000 nodes.

### 425 8.2 DCT Compression Results Under Assumed Correlation Model

426 The 2D type-II DCT was applied to the synthetic  $\delta$  field. Truncated reconstructions retaining  
427 only the top  $K \times K$  low-frequency coefficients were computed for  $K \in \{8, 16, 32\}$ . Table 3  
428 reports the variance explained ( $R^2$ ), storage cost, and compression ratio for each architecture  
429 under the assumed Gaussian correlation model.

Table 3: DCT compression of synthetic variability field under assumed Gaussian spatial correlation ( $\lambda/L = 0.25$ ). Array:  $4000 \times 4000$ , 16,520 conductance levels. Results are predictions under the stated correlation model, not experimentally validated outcomes.

$K$	Coefficients	Storage	$R^2$ (%)	Compression
8	64	0.50 KB	99.83	125,000×
16	256	2.00 KB	99.98	31,250×
32	1024	8.00 KB	99.998	7,812×

430 Under the assumed correlation model, the DCT captures  $> 99.8\%$  of the spatial variance  
431 with just 64 coefficients (0.50 KB). At  $K = 16$  (the architecture used throughout this paper),  
432  $R^2 = 99.98\%$ , confirming that the 2 KB storage budget is not merely sufficient but conservative  
433 if the spatial correlation assumption holds. The raw field requires 61 MB of storage; the  $K = 16$

434 DCT representation achieves a  $31,250\times$  compression ratio with negligible accuracy loss.

### 435 **8.3 Separability as an Empirical Test**

436 The variability field  $\delta(i, j)$  may be spatially separable, i.e.  $\delta(i, j) = \delta_{\text{row}}(i) \times \delta_{\text{col}}(j)$ . Separability  
437 is not assumed in the architectures above, but it is a falsifiable prediction that follows from  
438 the physics of spin-coating deposition: if row-wise and column-wise manufacturing variation  
439 arise from independent processes (e.g. film thickness gradients versus electrode deposition  
440 non-uniformity), the cross product form is expected.

441 If separability holds, the Architecture C storage requirement collapses dramatically: instead of  
442 16,000,000 per-node entries, only two vectors of 4000 entries each are needed, reducing storage  
443 to 32 KB. The characterisation problem also simplifies: instead of requiring a full  $N \times N$  sweep,  
444 row and column vectors can be estimated from  $\mathcal{O}(N)$  measurements.

445 Separability is testable from existing characterisation data: compute the rank of the observed  $\delta$   
446 matrix and check whether a rank-1 approximation captures the dominant variance. This makes  
447 the separability question a concrete deliverable for experimental collaborators with chip-level  
448 access. Whether the variability field is separable is the primary empirical question this proposal  
449 poses to Goswami et al. [1].

## 450 **9. Contributions**

451 This work establishes four core contributions to neuromorphic computing:

- 452 1. **Deterministic variability elimination at load time.** We introduce a per-chip cor-  
453 rection framework that cancels device-level variability prior to inference, achieving near-  
454 zero residual error under a linear response model and converting measurability from an  
455 instrument-precision constraint into a mathematically guaranteed convergence process.
- 456 2. **Full-matrix system identification via structured excitation.** We show that complete  
457 conductance deviation matrices can be recovered using Hadamard-structured inputs applied  
458 through existing peripheral circuitry, eliminating the need for external characterisation  
459 equipment or per-node selector devices.
- 460 3. **Complete model-hardware decoupling.** The framework enables deployment of ar-  
461 bitrary pre-trained models without hardware-aware training, fine-tuning, or inference  
462 overhead, making the entire ecosystem of neural networks developed for digital systems  
463 immediately accessible to memristor hardware.
- 464 4. **Universal drift correction mechanism.** The heartbeat protocol corrects all sources  
465 of conductance deviation — thermal drift, write noise, environmental perturbations —  
466 through a single measurement and correction cycle, functioning as a device-agnostic noise  
467 corrector over operational lifetime.

## 468 10. Conclusions

469 We have established theoretically that device-level variability in memristor crossbars can be  
470 deterministically eliminated at model load time via per-chip system identification, removing the  
471 need for hardware-aware training or inference overhead. This represents a fundamental shift  
472 from treating variability as unavoidable statistical noise to treating it as a measurable system  
473 property that can be cancelled through interface-layer pre-compensation.

474 The approach enables direct deployment of arbitrary pre-trained models on analog hardware,  
475 decoupling model development from hardware development for the first time in neuromorphic  
476 computing. Factory characterisation requires only 17.6 seconds using existing peripheral circuitry,  
477 while periodic heartbeat recalibration maintains correction accuracy over device lifetime through  
478 a universal noise correction mechanism.

479 Energy analysis confirms deployment feasibility below hardware-aware training costs for realistic  
480 training scenarios, with substantial optimization potential through heartbeat interval tuning  
481 based on operational drift rates.

482 We convert variability correction from an instrument-precision constraint into a mathematically  
483 guaranteed convergence process, establishing a new paradigm for reliable computation on analog  
484 neuromorphic hardware.

## 485 References

- 486 [1] P. Gaur, B. Kundu, P. Ghosh, S. Bhattacharya, L. T., H. S., S. P. Rath, D. Thomp-  
487 son, S. Goswami, S. Goswami, “Molecularly Engineered Memristors for Reconfigurable  
488 Neuromorphic Functionalities,” *Advanced Materials*, e09143, 2025.
- 489 [2] Z. Liu, C. Gao, J. Yang, Z. Chen, E. Li, J. Li, M. Li, J. Zhang, “Memristor devices for next-  
490 generation computing: from performance optimization to application-specific co-design,”  
491 *International Journal of Extreme Manufacturing*, vol. 8, no. 1, 012004, 2025.
- 492 [3] P. Chi *et al.*, “PRIME: A Novel Processing-in-Memory Architecture for Neural Network  
493 Computation in ReRAM-Based Main Memory,” *Proc. ISCA*, 2016.
- 494 [4] E. J. Candès, J. Romberg, and T. Tao, “Stable signal recovery from incomplete and  
495 inaccurate measurements,” *Communications on Pure and Applied Mathematics*, vol. 59,  
496 no. 8, pp. 1207–1223, 2006.
- 497 [5] D. L. Donoho, “Compressed sensing,” *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306,  
498 2006.
- 499 [6] A. M. Bruckstein, D. L. Donoho, and M. Elad, “From Sparse Solutions of Systems of  
500 Equations to Sparse Modeling of Signals and Images,” *SIAM Review*, vol. 51, no. 1, pp. 34–  
501 81, 2009.
- 502 [7] L. Ljung, *System Identification: Theory for the User*, 2nd ed., Prentice Hall, 1999.

- 503 [8] A. Hedayat and W. D. Wallis, “Hadamard Matrices and Their Applications,” *Annals of*  
504 *Statistics*, vol. 6, no. 6, pp. 1184–1238, 1978.
- 505 [9] A. Papoulis and S. U. Pillai, *Probability, Random Variables, and Stochastic Processes*, 4th  
506 ed., McGraw-Hill, 2002.
- 507 [10] E. Linn, R. Rosezin, C. Kügeler, and R. Waser, “Complementary resistive switches for  
508 passive nanocrossbar memories,” *Nature Materials*, vol. 9, pp. 403–406, 2010.
- 509 [11] D. Sharma, S. P. Rath, B. Kundu, A. Korkmaz, H. S., D. Thompson, N. Bhat, S. Goswami,  
510 R. S. Williams, S. Goswami, “Linear Symmetric Self-Selecting 14-bit Kinetic Molecular  
511 Memristors,” *Nature*, vol. 633, pp. 560–566, 2024.
- 512 [12] Z. Xiao, V. B. Naik, J. H. Lim, Y. Hou, Z. Wang, Q. Shao, “Adapting Magnetoresistive  
513 Memory Devices for Accurate and On-Chip-Training-Free In-Memory Computing,” *Science*  
514 *Advances*, vol. 10, no. 38, eadp3710, 2024.
- 515 [13] S. Ambrogio *et al.*, “Equivalent-accuracy accelerated neural-network training using analogue  
516 memory,” *Nature*, vol. 558, pp. 60–67, 2018.
- 517 [14] P. Yao *et al.*, “Fully hardware-implemented memristor convolutional neural network,” *Nature*,  
518 vol. 577, pp. 641–647, 2020.
- 519 [15] T. Brown *et al.*, “Language Models are Few-Shot Learners,” *Advances in Neural Information*  
520 *Processing Systems*, 2020.
- 521 [16] A. Dosovitskiy *et al.*, “An Image is Worth 16x16 Words: Transformers for Image Recognition  
522 at Scale,” *Proc. ICLR*, 2021.
- 523 [17] J. Ho, A. Jain, and P. Abbeel, “Denoising Diffusion Probabilistic Models,” *Advances in*  
524 *Neural Information Processing Systems*, 2020.
- 525 [18] F. Alibart, L. Gao, B. D. Hoskins, D. B. Strukov, “High precision tuning of state for  
526 memristive devices by adaptable variation-tolerant algorithm,” *Nanotechnology*, vol. 23,  
527 no. 7, 075201, 2012.
- 528 [19] Y. Ma, L. Zheng, P. Zhou, “A Mapping Method Tolerating SAF and Variation for Memristor  
529 Crossbar Array Based Neural Network Inference on Edge Devices,” *ACM Journal on*  
530 *Emerging Technologies in Computing Systems*, 2023.
- 531 [20] N. Ahmed, T. Natarajan, and K. R. Rao, “Discrete Cosine Transform,” *IEEE Trans.*  
532 *Computers*, vol. C-23, no. 1, pp. 90–93, 1974.
- 533 [21] G. K. Wallace, “The JPEG still picture compression standard,” *IEEE Communications*  
534 *Magazine*, vol. 34, no. 4, pp. 31–44, 1996.
- 535 [22] K. Sayood, *Introduction to Data Compression*, 5th ed., Morgan Kaufmann, 2017.
- 536 [23] E. Strubell, A. Ganesh, and A. McCallum, “Energy and Policy Considerations for Deep  
537 Learning in NLP,” *Proc. ACL*, pp. 4535–4546, 2019.

- 538 [24] D. Patterson *et al.*, “Carbon Emissions and Large Neural Network Training,”  
539 arXiv:2104.10350, 2021.
- 540 [25] J. Dangelad-Flores, S. Eickelmann, H. Riegler, “Deposition of polymer films by spin coating:  
541 A modelling approach,” *Chemical Engineering Science*, vol. 179, pp. 257–264, 2018.