# Hardware-aware training for large-scale and diverse deep learning inference workloads using in-memory computing-based accelerators

Malte J. Rasch[*1], Charles Mackin[2], Manuel Le Gallo[3], An Chen[2], Andrea Fasoli[2], Frédéric Odermatt[3], Ning Li[1], S.R. Nandakumar[3], Pritish Narayanan[2], Hsinyu Tsai[2], Geoffrey W. Burr[2], Abu Sebastian[3], and Vijay Narayanan[1]

[1]*IBM Research, TJ Watson Research Center, Yorktown Heights, NY USA*
[2]*IBM Research Almaden, 650 Harry Road, San Jose, CA USA*
[3]*IBM Research Europe, 8803 Rüschlikon, Switzerland*

February 17, 2023

## Abstract

Analog in-memory computing (AIMC) – a promising approach for energy-efficient acceleration of deep learning workloads – computes matrix-vector multiplications (MVMs) but only approximately, due to nonidealities that often are non-deterministic or nonlinear. This can adversely impact the achievable deep neural network (DNN) inference accuracy as compared to a conventional floating point (FP) implementation. While retraining has previously been suggested to improve robustness, prior work has explored only a few DNN topologies, using disparate and overly simplified AIMC hardware models. Here, we use hardware-aware (HWA) training to systematically examine the accuracy of AIMC for multiple common artificial intelligence (AI) workloads across multiple DNN topologies, and investigate sensitivity and robustness to a broad set of nonidealities. By introducing a new and highly realistic AIMC crossbar-model, we improve significantly on earlier retraining approaches. We show that many large-scale DNNs of various topologies, including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformers, can in fact be successfully re-trained to show iso-accuracy on AIMC. Our results further suggest that AIMC nonidealities that add noise to the inputs or outputs, not the weights, have the largest impact on DNN accuracy, and that RNNs are particularly robust to all nonidealities.

## 1 Introduction

The ever-increasing compute needed to train and use DNNs [72] have made hardware latency and energy-efficiency a growing concern. However, conventional processor architectures (e.g. CPUs, GPUs, etc.) incessantly transfer data between memory and processing through the 'von Neumann bottleneck,' inducing time and energy overheads that significantly degrade latency and energy-efficiency. Numerous hardware concepts have been introduced to accelerate DNN training and/or inference [77, 35, 68], by approximating MVMs and other arithmetic with custom floating-point representations such as bfloat16 [84] or DLFloat [2], or with reduced-precision fixed-point arithmetic to quantize synaptic weights and activations [76, 15, 16, 67]. Model

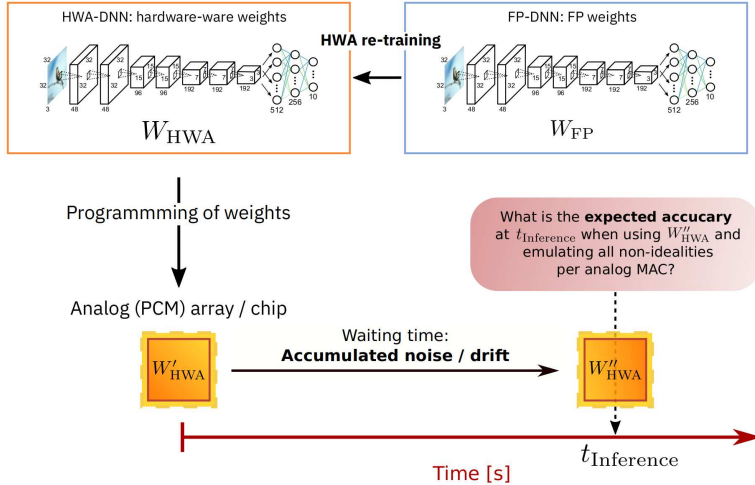---

*Correspondence: malte.rasch@ibm.com

Figure 1: In our hardware-aware (HWA) training setup, DNNs are first trained in 32bit floating-point (FP$_{32}$), then re-trained in a hardware-aware manner, by adding nonidealities and noise sources into the forward path and using SGD to improve the robustness to such generic nonidealities. HWA training is only performed once — no specific device or chip characteristics, such as failure maps, are taken into account during HWA training, so resulting models remain widely-deployable. This HWA-trained model is then programmed onto AIMC multiple times (here in simulation) and DNN accuracy is evaluated over time, taking into account conductance drift of PCM devices and read noise accumulation [56]

.

compression and sparsification techniques can further reduce compute requirements by pruning weights and/or activations [3, 27].

AIMC using non-volatile memory (NVM) elements is a promising mixed-signal approach for DNN acceleration [11, 71, 10], with weights stored using crossbar arrays of tuneable conductive elements. This enables approximate MVM computation directly in-memory, by applying activation vectors (as voltages or pulse durations) to the crossbar array, and then reading out analog physical quantities (instantaneous current or accumulated charge) [52, 12, 53]. As a 'non-von Neumann' architecture, AIMC performs MVM operations at the location of the stored weights, in a highly-parallel, fast, and energy-efficient manner [12] – but only approximately.

The success of reduced-precision digital accelerators proved that DNNs can tolerate surprisingly coarse approximations of the underlying MVMs. While naive direct weight-quantization invariably leads to DNN accuracy loss, original model accuracies can often be recovered when DNNs are re-trained in a quantization-aware manner, even for aggressive reductions in precision. Weight-quantization into as few as 2-4 bits can often be tolerated without significant accuracy reduction [42, 15, 54]. This observation led to the development of quantization-aware training (QAT) methods, now commonly used when deploying DNNs onto reduced-precision digital hardware (e.g. [1]).

In general, since reducing MVM precision decreases the representational power of each DNN layer as compared to a full FP implementation, performance naturally suffers once the function approximation becomes too coarse for the task at hand [54]. In practice, QAT is known to have limits, and MVM minimum-precision requirements vary across each DNN topology. For instance, the first and last layers of CNNs are invariably implemented with high precision FP, even in studies claiming CNN iso-accuracy at very low fixed-point precision [76, 15].

Up to now, it has been unclear how and to what degree DNNs can be re-trained to maintain accuracy on emerging AIMC technology. The successes of QAT cannot be directly translated onto AIMC, since the MVM approximations arise from fundamentally different concepts. In AIMC, weights are represented by conductances that are physical properties of NVM devices. In many materials, such as PCM [9, 45], resistive random-access memory

(ReRAM) [33, 34], conductive-bridge RAM (CBRAM) [49, 21], or electro-chemical random-access memory (ECRAM) [47, 59], these conductances are effectively continuous physical quantities, and stored weights are not quantized.

That said, effective AIMC weight-precision is impacted by various nonidealities, including thermal and 1/f noise, randomization during physical switching induced by electrical and thermal fluctuations, material inhomogenities [14], and device-to-device variability introduced during device fabrication or operation. These issues cause both MVM read-out [56] and the writing or programming of the conductances [60, 79, 51] to be erroneous and non-deterministic. Worse yet, conductances can evolve over time after programming [6, 4, 7]. Finally, any nonlinearities within the analog circuitry performing summation will further degrade MVM precision. Such imperfections include 'IR-drop' voltages on wires and transistors, restrictions on input (output) dynamic-range imposed by discretization and saturation of the digital-to-analog converter (DAC) (analog-to-digital converter (ADC)) components, and random noise or variability in the circuitry.

Whereas QAT gets challenging as precision is deterministically reduced, MVM approximation in AIMC is tied to non-deterministic signal-to-noise ratio. A number of previous studies have shown that noise-aware training – simple injection of noise onto weights or activations during DNN training – can make DNNs more robust for AIMC deployment [79, 36, 87, 24, 38, 73]. However, such studies have typically been limited to one or two exemplary DNNs of a particular type (e.g., CNNs) using only a limited subset of nonidealities such as NVM noise. Other critical AIMC system aspects such as output noise, saturation, and circuit nonlinearities have been neglected. Moreover, since each study makes different hardware and NVM-device choices, it is difficult to generalize, compare, or combine them. Thus more realistic and standardized AIMC crossbar-models – which can support comparison of AIMC accuracy for hardware-aware trained DNN models across studies – are needed.

Although some promising, small-sized DNN prototype demonstrations exist [83, 39, 86, 20, 57, 5, 88], it remains unclear how robust the AIMC deployment of realistically-sized AI workloads will be. How will the various nonidealities of AIMC hardware impact the DNN accuracy, across all the various topographies and thus application domains? And how much of the lost accuracy could be recovered by hardware-aware training? Which crossbar-array design choices will be most effective in maintaining accuracy? And if necessary, what degree of improved device-to-device uniformity might be required – through better NVM device-fabrication – in order for AIMC to succeed on all DNN models? A systematic study comparing the various DNN topographies in terms of robustness to AIMC nonidealities is needed.

In this paper, we establish a robust HWA framework by extending and improving existing training methods for AIMC to include previously neglected nonidealities (see Fig. 1 for an illustration). We define a standard inference model for PCM-based AIMC that can readily be extended to other types of NVM devices. We explore the functional suitability of AIMC across application domains by assessing the robustness of a wide set of DNN topographies. Finally, we estimate the individual impact of various AIMC nonidealities and gauge their relative importance for consideration in future hardware designs. Functions for our standard evaluation process are provided in an open-source IBM Analog Hardware Acceleration Toolkit (AIHWKit) [66], enabling future studies on noise robustness for AIMC to build seamlessly upon our work.

We find that various DNNs and AI workloads – ranging across image classification using CNNs, text-prediction and speech-to-text conversion using RNNs, and natural language processing using transformer networks – can actually be robustly deployed on AIMC given proper HWA training. We show – for the first time – iso-accuracy inference results (within 1% of the FP reference) using hardware-calibrated PCM models, for five out of the eleven AI workloads tested, even after 1 hour (or more) of conductance drift.

However, precision requirements are heterogeneous, and not all architectures reach this iso-accuracy target easily even after extensive HWA training, pinpointing the need for continued device improvement. We find that CNNs are typically much less robust to various nonidealities and design choices of AIMC hardware. Interestingly, RNNs – already well-suited for AIMC
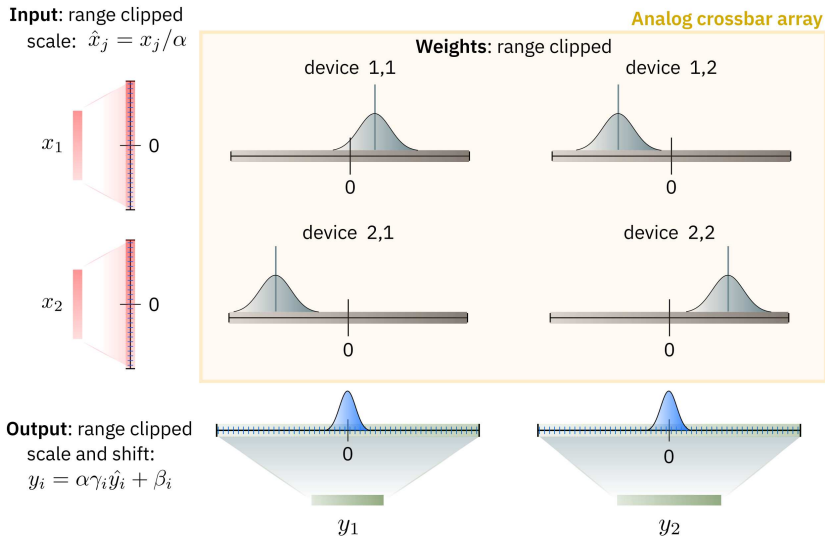
Figure 2: Our AIMC crossbar-model assumes that each array or "tile" approximates matrix-vector multiplication (MVM) $\widetilde{\mathbf{y}} \approx W\mathbf{x}$, using a scalar $\alpha$ (and vectors $\boldsymbol{\gamma}$ and $\boldsymbol{\beta}$) to scale into (from) quantized inputs (outputs) so that input, weight, and output ranges can remain fixed. Negative weights are programmed onto a different conductance for current subtraction, and output noise is fully represented (see Eqs. 3 and 6).

given their large, dense MVMs [32] – also seem to be the most robust to the finite signal-to-noise ratio (SNR) of AIMC hardware. We further show that among various nonidealities tested, the sensitivity to additive system noise at the output of each crossbar-array is the most critical for achieving good accuracy.

## 2 Results

### 2.1 Analog IMC standard MVM model

Our standard AIMC crossbar-model (see Figs. 2 and 3, and Eqs. 3 and 6 in Methods) encapsulates the critical nonidealities incurred during MVM operations, including the fixed dynamic-ranges of physical inputs (limited by maximum pulse-duration), weights (limited by maximum conductance), and outputs (limited by maximum output current). Floating point inputs $x$ are mapped to the fixed input range by a global scalar, $\alpha$, which is optimized for each crossbar during HWA training, then held fixed for inference. Such optimization avoids issues created if $\alpha$ is chosen poorly (Fig. 3E). Similarly optimized scales ($\gamma_i$) and offsets ($\beta_i$) map ADC-counts of each output column to MVM-outputs $\widetilde{y}_i$ (see Eq. 3) that can be passed to subsequent digital compute for auxiliary operations (activation functions, etc.) [32]. We further assume a number of nonidealities, such as PCM programming errors and drift(Fig. 3A), IR-drops within the array (Fig. 3C), weight read noise, and system noise (Fig. 3B and Eq. 6). Each of these parameters has been carefully calibrated to existing PCM hardware [56]. All parameter settings of the AIMC crossbar-model are summarized in Table A.1.

We quantify MVM errors in computing $\widetilde{\mathbf{y}}$ with respect to the ideal outcome $\mathbf{y}$ through $\epsilon_M$, the ratio of the $l_2$-norm of the deviation $(\mathbf{y} - \widetilde{\mathbf{y}})$ relative to the $l_2$-norm of the ideal outcome $\mathbf{y}$ (see Eq. 20). Figure 3B shows that, even after including PCM drift, the effective MVM error of the standard AIMC crossbar-model we will use throughout the paper roughly corresponds to 4bit fixed-point quantization of weights or inputs.
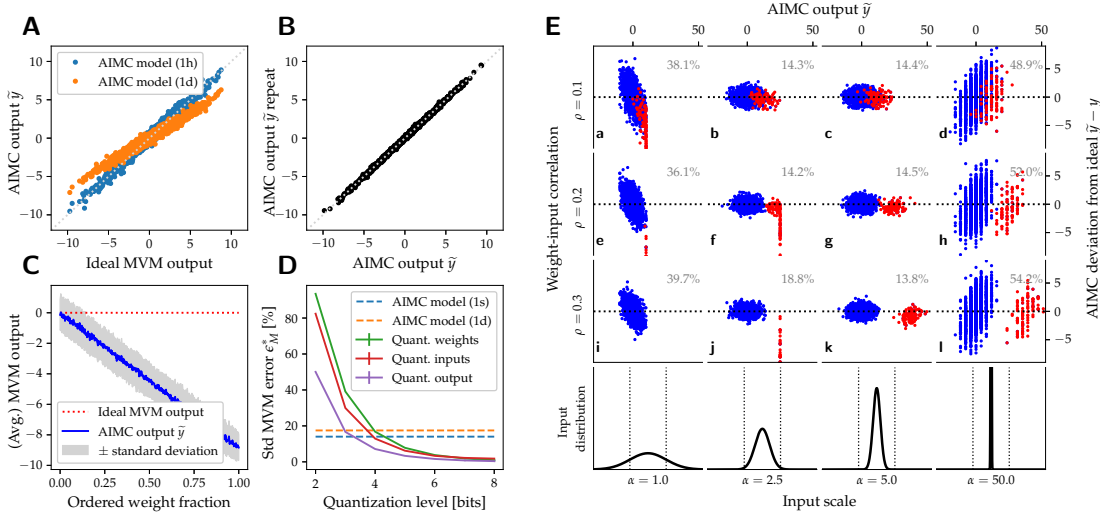
Figure 3: (A) Correlations between AIMC outputs – for MVMs performed between Gaussian random weight matrices and uniform random inputs – and the ideal ($FP_{32}$) expected results reveal significant deviations, due primarily to weight-programming errors and PCM conductance drift (shown here without any mean-drift compensation [4]). (B) Short-term noise sources induce cycle-to-cycle noise for repeated MVM calculations even with the same programmed weight matrix. (C) "IR-drops" due to finite wire resistance result from input-position dependency of the accumulated AIMC column-currents, and can cause outputs to further deviate from the ideal MVM. In an extreme case, an expected 0 output – the correct result when a linearly-graded weight matrix (ranging from -1 to 1 in order) is read with a constant input on all rows – can actually deviate drastically due to these position-dependent IR-drops. The deviation induced by these IR-drops increases as the fraction of ordered weights increases from 0 (completely unordered, typical case) to 1 (fully-ordered weights from -1 to 1, extreme case). (D) The MVM-error $\epsilon_M^*$ (Eq. 20) of our standard PCM-based AIMC inference model (using the nonideal MVM of Eqs. 3 and 6 together with hardware-calibrated PCM conductance noise and drift Eqs.7-14; dotted lines) ≈15% roughly corresponds to fixed-point quantized digital (solid lines) at ∼ 4 bits. (E) Correlations of MVM deviation, $\tilde{y}_i - y_i$, vs. desired MVM output $y_i = w_{ij}x_j$ illustrate the importance of proper input scaling $\alpha$, for $w_{ij} \sim \mathcal{N}(0, 0.246)$ and $x_j \sim \mathcal{N}(0, 1)$. Red dots mimic the weight-to-activation correlations that SGD learning will produce, using $\tilde{x}_j = \rho w_{kj} + (1 - \rho)x_j$. Low $\alpha = 1$ leads to input clipping (panels a, e, i) and $\epsilon_M^*$ exceeding 35% (gray text). Intermediate $\alpha$ values can still lead to saturated outputs for correlated inputs, even without input clipping (panels f,j); excessive $\alpha$ values reduce clipping but increase $\epsilon_M^*$ dramatically (panels d, h, l). We optimize $\alpha$ during HWA training, then keep it fixed during AIMC inference, minimizing $\epsilon_M^*$ regardless of input correlation (panels c, g, k).

## 2.2 DNN accuracy impact when directly using AIMC

To test the effect of AIMC nonidealities on a variety of AI workloads, we consider 11 medium-to large-scale DNNs of various topologies as a benchmark set (see Table 1). These cover a wide spectrum of target-applications (image classification, natural language processing, speech-to-text), network topologies (convolutional, recurrent, transformer with attention), model sizes (from 0.3M to 108M parameters), crossbar utilization (from 4% to 86%), total number of MVMs per data input (from 2.2K to 240K), MVM sizes (from 0.1G to 96.5G flops), average weight-matrix reuse-factor per data-input (from 17 to 1285), and network depth (up to 121 layers). Our benchmark set thus covers a wide variety of network topologies and challenges for AIMC.

For comparison, we first directly map weights produced by standard SGD-training in $FP_{32}$ onto our standard AIMC crossbar-model and evaluate the resulting test error, to measure the accuracy drop (with respect to the $FP_{32}$ model) due to all the various AIMC nonidealities. For each crossbar, we first clip raw weight-values to 2.5× the standard deviation of the weight distribution, to disregard outliers and compactify the weight distribution. Output scales $\gamma_i$ are

| DNN | Type | # par. | # mapped | # tiles | util. | # MVM | ⟨reuse⟩ | flops |
|---|---|---|---|---|---|---|---|---|
| ResNet-32 CF10 | **C** | 0.36M | 0.36M | 34 | 4.0% | 14K | 435. | 0.1G |
| WideRN-16 CF100 | **C** | 11M | 11M | 68 | 61.7% | 19K | 286. | 3.1G |
| ResNet-18[†]ImNet | **C** | 11.7M | 11.2M | 75 | 56.8% | 40K | 533. | 3.4G |
| ResNet-50[†]ImNet | **C** | 25.6M | 23.4M | 149 | 60.0% | 74K | 497. | 7.9G |
| DenseNet-121[†]ImNet | **C** | 8M | 6.9M | 266 | 9.8% | 143K | 537. | 5.4G |
| WideRN-50[†]ImNet | **C** | 69M | 66.8M | 296 | 86.0% | 101K | 342. | 22.6G |
| BERT-base MRPC | **T** | 108M | 85M | 486 | 67.1% | 61K | 126. | 21.8G |
| Albert-base MRPC | **T** | 12M | 7.8M | 48 | 61.7% | 61K | 1285. | 21.8G |
| Speech[□]SWB300 | **L** | 30M | 30M | 153 | 74.8% | 2.5K | 17. | 0.9G |
| LSTM PTB | **L** | 19.8M | 13.3M | 88 | 57.5% | 2.2K | 26. | 0.7G |
| RNN-T | **L** | 57M | 57M | 304 | 71.6% | 240K | 790.[*] | 96.5G |

Table 1: A wide range of DNNs topologies (Type) and sizes are studied, including CNNs (**C**), LSTMs (**L**), and Transformers (**T**). For each DNN-model and dataset, model size is quantified by number of parameters (# par.); number of parameters mapped to analog crossbars (# mapped); number of 512 × 512 crossbars (# tiles) needed for naive mapping (each weight matrix gets at least 1 tile); overall utilization of the devices within the tiles (util.); total number of MVMs per input data (# MVM); average tile-reuse for one input data (⟨reuse⟩; [*] for maximal input length in dataset); and the number of $FP_{32}$ operations in the mapped MVM for one input data (FLOPS). [†]first conv-layer and last FC layer in $FP_{32}$, [□]additional hidden-Markov-model used as decoder.

initially estimated according to the absolute maximum value for each column (see Eq. 21). To adjust the digital parameters of our standard AIMC crossbar-model for these directly-mapped-from-software weights, we use our HWA-training flow – but without any weight noise injection, with weight learning rates set to zero, and for only 1000 batches. As expected, such "direct" mapping of DNNs onto AIMC, without any additional retraining of the weights, generally results in significant increases in test error (accuracy drop) in comparison to the floating point reference (Table 2).

Direct comparison of accuracy values between DNNs is complicated by the fact that these various AI tasks exhibit different worst-case (random guessing) and best-case (well-trained DNN model) accuracies. To quantify and compare accuracy drop across different topologies, we therefore define a normalized relative accuracy $\mathcal{A}_*^{1h}$, which re-scales the AIMC test error $\epsilon_{test}^{1h}$ (at 1 hour PCM drift) by the distance between the original $FP_{32}$ test error and the "chance" test error from random guessing, as follows:

$$\mathcal{A}_*^{1h} = 1 - \frac{\epsilon_{test}^{1h} - \epsilon_{test}^{FP}}{\epsilon_{chance} - \epsilon_{test}^{FP}}. \tag{1}$$

Thus a value of $\mathcal{A}_*^{1h} = 100\%$ means that the AIMC DNN achieves the same accuracy as the $FP_{32}$ reference model (no accuracy drop at all), while a value of $\mathcal{A}_*^{1h} = 0\%$ implies that the AIMC crossbar-model is so inaccurate that it is indistinguishable from random guessing.

Ideally, deploying a given DNN in an AIMC system should have no impact on model accuracy. We define our iso-accuracy target as $\mathcal{A}_*^{1h} > 99\%$, allowing less than a 1% drop in accuracy, as judged relative to the distance between the $FP_{32}$ reference accuracy and the chance (random guessing) accuracy-floor. Table 2 shows that direct AIMC mapping fails to achieve this iso-accuracy target for almost all of the DNNs tested, establishing both the challenge posed by the nonidealities existing in AIMC (as compactly encapsulated by our standard crossbar-model, Figs. 2 and 3), as well as the need for HWA training methods that can greatly improve the robustness and reduce these accuracy drops.

| Direct mapping | Test Error in % | | | Normalized accuracy | |
|---|---|---|---|---|---|
| DNN | FP$_{32}$ | 1 hour | 1 year | $\mathcal{A}_*^{1h}$ [%] | $\mathcal{A}_*^{1y}$ [%] |
| ResNet-32 $_{\text{CF10}}$ | 5.80 | 11.68 $_{0.16}$ | 15.21 $_{0.40}$ | 93.0 | 88.8 |
| WideRN-16 $_{\text{CF100}}$ | 20.00 | 26.02 $_{0.10}$ | 29.42 $_{0.28}$ | 92.4 | 88.1 |
| ResNet-18$^\dagger_{\text{ImNet}}$ | 30.50 | 47.91 $_{0.21}$ | 55.00 $_{0.41}$ | 74.9 | 64.7 |
| ResNet-50$^\dagger_{\text{ImNet}}$ | 23.87 | 32.51 $_{0.09}$ | 38.18 $_{0.31}$ | 88.6 | 81.2 |
| DenseNet-121$^\dagger_{\text{ImNet}}$ | 25.57 | 47.38 $_{0.24}$ | 59.81 $_{0.71}$ | 70.7 | 53.9 |
| WideRN-50$^\dagger_{\text{ImNet}}$ | 21.53 | 40.07 $_{0.13}$ | 45.45 $_{0.30}$ | 76.3 | 69.5 |
| BERT-base $_{\text{MRPC}}$ | 14.60 | 15.99 $_{0.23}$ | 18.36 $_{0.30}$ | 97.3 | 92.7 |
| Albert-base $_{\text{MRPC}}$ | 15.08 | 31.50 $_{0.10}$ | 31.62 $_{0.03}$ | 67.8 | 67.5 |
| Speech$^\square_{\text{SWB300}}$ | 14.05 | 16.11 $_{0.02}$ | 16.32 $_{0.02}$ | 97.6 | 97.4 |
| LSTM $_{\text{PTB}}$ | 72.90 | 73.16 $_{0.01}$ | 73.28 $_{0.01}$ | **99.0** | 98.6 |
| RNN-T $_{\text{SWB300}}$ | 11.80 | 28.16 $_{0.10}$ | 29.08 $_{0.16}$ | 81.4 | 80.4 |

Table 2: Inference results using AIMC for the 11 benchmark DNNs when deployed directly without any weight retraining. Test errors in % ± standard error of mean (across 24 inference repeats) are shown after 1 hour and 1 year of PCM drift (center columns) and compared to the original FP$_{32}$ test error (leftmost column). Digital parameters needed for the AIMC crossbar-model are estimated by training briefly with the AIMC MVM in the forward pass (1000 batches), but without touching the directly-mapped analog weights. This helps adjust statistics of each batch norm to the new output distributions caused by the AIMC MVMs. During these 1000 batches, we estimate $\alpha$ by averaging the maximal absolute inputs for each batch during the first 500 batches, and then allow SGD to tune it further during the second half of the brief digital-parameter-only HWA-training. Righthand columns show normalized accuracy values after 1 hour ($\mathcal{A}_*^{1h}$) and 1 year ($\mathcal{A}_*^{1y}$), as scaled to the range between the FP$_{32}$ reference test-error and the test-error obtained by random guessing. Nearly all models fail to achieve iso-accuracy, as defined by > 99% in this normalized accuracy and indicated in bold font. Note that for BERT and Albert only one GLUE task (MRPC) is used here.

## 2.3 HWA training improves AIMC accuracy for all DNNs

Building on previous approaches (e.g. [36, 38, 24]), we set out to retrain these 11 DNNs in a hardware-aware (HWA) manner. In our methodology for HWA training followed by delayed inference (Fig. 1), each DNN is retrained with noise injection using SGD. But in contrast to earlier approaches, we incorporate a much more comprehensive and realistic set of software-simulated AIMC nonidealities, including dynamic-range limitations, weight-programming errors, PCM drift and system noise. Once a given DNN is trained and mapped to AIMC, inference is then gauged for noise and drift at various delays (1 second, 1 hour, 1 day, and 1 year) after programming the weights into the crossbar arrays. We also introduce a set of AIMC characteristics including input, output and weight-scales (see Fig. 3 and Methods), and introduce a new approach for optimizing these scaling-factors during HWA training for use during inference (see Sec. A.1).

As shown in Table 3, our HWA training approach significantly improves achievable accuracy for AIMC across the full set of benchmark DNNs results. The normalized accuracies (relative to the FP$_{32}$ model) at one hour after programming are all higher than 96% ($\mathcal{A}_*^{1h}$, towards right edge of Table 3). This represents a significant improvement over 'direct' weight mapping without retraining shown earlier (Table 2), while establishing a new state-of-the-art in HWA training, as revealed by detailed comparisons on ResNet-32 with CIFAR10 (see Table A.2).

Table 3 indicates that five out of the 11 AI workloads can be trained to reach the $\mathcal{A}_*^{1h} > 99\%$ iso-accuracy target, including the BERT transformer model as well as all workloads based on LSTMs (last 3 rows, see 'Type' column in Table 1). Most of the remaining workloads use CNNs and exhibit more-pronounced accuracy drops of up to 3.6% on AIMC, although one CNN does reach iso-accuracy (WideResNet-32 on Cifar100).

For some DNNs, we find that the regularization effect of the added AIMC nonidealities allows HWA-training to actually improve the attainable accuracy (compare test-errors at 1 sec

| HWA training | Test Error in % | | | | | Normalized accuracy | |
|---|---|---|---|---|---|---|---|
| DNN | $FP_{32}$ | 1 sec | 1 hour | 1 day | 1 year | $\mathcal{A}_*^{1h}$ | $\mathcal{A}_*^{1y}$ |
| ResNet-32 CF10 | 5.80 | 6.79 0.02 | 7.08 0.03 | 7.50 0.04 | 8.36 0.06 | 98.5 | 97.0 |
| WideRN-16 CF100 | 20.00 | 19.75 0.02 | 20.03 0.02 | 20.27 0.02 | 21.31 0.04 | **100.0** | 98.3 |
| ResNet-18$^\dagger$ ImNet | 30.50 | 31.57 0.06 | 32.21 0.10 | 32.94 0.11 | 35.31 0.27 | 97.5 | 93.1 |
| ResNet-50$^\dagger$ ImNet | 23.87 | 24.52 0.02 | 24.82 0.05 | 25.24 0.06 | 27.16 0.13 | 98.8 | 95.7 |
| DenseNet-121$^\dagger$ ImNet | 25.57 | 27.39 0.07 | 28.22 0.10 | 29.36 0.15 | 36.09 0.31 | 96.4 | 85.8 |
| WideRN-50$^\dagger$ ImNet | 21.53 | 23.38 0.05 | 23.82 0.04 | 24.42 0.07 | 27.87 0.14 | 97.1 | 91.9 |
| BERT-base GLUE8 | 17.47 | 17.43 0.09 | 17.55 0.12 | 17.58 0.12 | 17.99 0.12 | **99.8** | 98.9 |
| Albert-base GLUE8 | 19.46 | 20.52 0.18 | 20.45 0.16 | 21.08 0.18 | 22.18 0.21 | 97.8 | 94.0 |
| Speech$^\square$ SWB300 | 14.05 | 14.26 0.03 | 14.29 0.03 | 14.35 0.03 | 14.52 0.03 | **99.7** | **99.4** |
| LSTM PTB | 72.90 | 72.94 0.01 | 72.94 0.01 | 72.96 0.01 | 73.04 0.01 | **99.8** | **99.5** |
| RNN-T SWB300 | 11.80 | 12.24 0.15 | 12.36 0.09 | 12.56 0.13 | 13.28 0.28 | **99.4** | 98.3 |

Table 3: Test error in % ± standard error of mean (across 15-25 inference repeats per training trial and up to 3 training trails) for DNN deployment on AIMC crossbars after HWA training. Rightmost two columns show the normalized accuracy, scaled between the FP reference and chance error, at 1 hour and 1 year after weight programming. Note that PCM drift is a post-programming physical effect that is initially rapid but then slows down logarithmically in time [45]. This means that the multiplicative conductance-changes induced by drift between 1 second and 1 hour (time-since-programming increased 3600×), and between 1 hour and 1 year (time-since-programming increased 8760×) are actually quite similar. HWA training hyper-parameters (injected noise strength, etc.) were chosen to produce the best average accuracy across the four widely-spaced time-points shown here. Other choices could be made to focus just on performance in either longer or shorter periods of drift. Models deemed iso-accurate ($\mathcal{A}_*^{1h}, \mathcal{A}_*^{1y} > 99\%$) are marked in bold. BERT and Albert results are averaged across eight GLUE tasks, as evaluated on validation datasets; SWB300 results are averaged over two benchmark tasks; results for Speech–SWB300 use HWA with distilling (see Methods).

after programming for WideRN-16 and BERT). Both RNNs and transformers are quite robust when subject to PCM conductance drift over longer periods as well. The rightmost column of Table 3 shows the long-term relative accuracy of the DNNs, $\mathcal{A}_*^{1y}$, for an hypothetical 1 year after programming without weight refresh.

While the RNNs and transformers remain near iso-accuracy over time, larger CNNs with higher resolution ImageNet inputs show the largest drop in accuracy. The deep DenseNet-121 (121 layers), as well as the large WideResNet-50 (69M parameters) models are clearly the most challenging for AIMC. That said, the resiliency to long-term drift is greatly improved by HWA-training as compared to "direct" deployment without retraining. For instance, the HWA-trained models for both the Speech-SWB300 and LSTM-PTB models remain iso-accurate out to a year, unlike the directly-mapped models (Table 2).

In general, we find that CNNs are more difficult to train to iso-accuracy for AIMC deployment compared to RNNs and transformers. In terms of AIMC workload execution latency and system mapping [32], CNNs are already less well-suited for resistive crossbar arrays due to the uneven temporal re-use between layers and spatial under-utilization of the large analog tiles by the small kernel matrices (see Table 1), although some optimization and mapping tricks [65] are available. Our results here indicate that AIMC noise-robustness issues will pose additional challenges when implementing CNNs onto AIMC systems.

## 2.4 Sensitivity of HWA trained models to various AIMC nonidealities

To determine which nonidealities are particularly problematic for analog inference across DNNs, we "stress test" our HWA-trained models. For each individual nonideality, such as PCM programming-error or IR-drop, we vary its strength and evaluate the resulting inference accuracy across DNNs using our base HWA-trained model. Our standard AIMC MVM model exhibits $\epsilon_M^* \approx 15\%$ (see Fig. 3 and Eq. 20), but combines many nonidealities. To estimate the
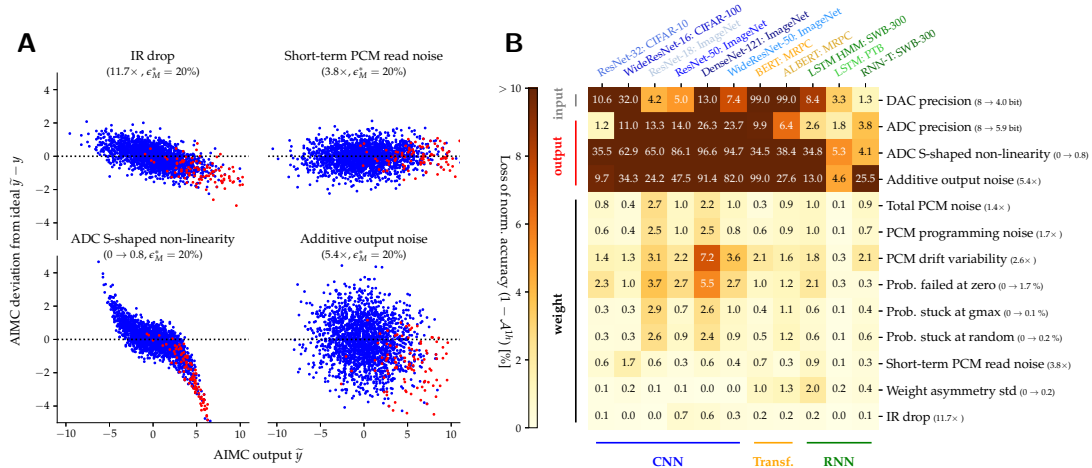
Figure 4: Comparison of the relative impact of various AIMC nonidealities on DNN accuracy. (A) AIMC deviations ($\tilde{y} - y$) from the ideal MVM output ($y$) are shown, for uncorrelated (blue dots) and weakly-correlated ($\rho = 0.05$) random activations (red dots), as a single nonideality is increased until standard MVM error reaches $\epsilon_M^* = 20\%$. All other parameters remain fixed to our standard AIMC crossbar-model ($\epsilon_M^* = 15\%$, Fig. 3). For instance, IR-drop needs to be scaled $11.7\times$ to incur $\epsilon_M^* = 20\%$. Even at constant $\epsilon_M^* = 20\%$, MVM deviations are structured differently and thus the impact on DNN accuracy can vary significantly. (B) Grid shows loss in normalized accuracy ($\mathcal{A}^{1h}$) over the base HWA-trained model at 1 hour after programming when boosting a given nonideality to $\epsilon_M^* = 20\%$. Thus 0% means no accuracy impact despite the amplified nonideality, whereas 100% means a drop to chance level. For the HMM LSTM sensitivity is reported for a portion of the training set (instead of the benchmark set) directly on the LSTM output without the Hidden-Markov-Model to speed up computations. For the transformer models, only one GLUE task is evaluated (MRPC).

relative accuracy impact due to each individual nonideality, we boost only that parameter value until MVM error increases to $\epsilon_M^* = 20\%$, and then re-measure DNN accuracy.

Even at constant MVM error, each parameter changes a different aspect of the AIMC compute. For instance, output noise is applied at each MVM, whereas PCM programming errors are only applied during programming and then persist throughout inference. Other nonidealities such as IR-drop or ADC "S-shaped" nonlinearity change the shape of the MVM deviations (Fig. 4A), causing large outputs to incur very significant MVM error. As a result, even at an identical average MVM error of $\epsilon_M^* = 20\%$, the impact on DNN accuracy can be much more pronounced. Such nonidealities are particularly detrimental for DNN inference, and thus deserve additional attention in future hardware designs or HWA training methods.

To gauge the relative impact of each individually-boosted nonideality parameter, Fig. 4B shows the loss in normalized accuracy ($\mathcal{A}^{1h}$), defined not with respect to the FP$_{32}$ model error ($\mathcal{A}_*^{1h}$ (Eq. 1)), but with respect to our standard AIMC crossbar-model (at 1 hour drift). A value of 0% means that boosting this particular nonideality has no impact on accuracy, as compared to our standard AIMC crossbar-model. A value of 100% means that simply boosting this nonideality to the same MVM error of $\epsilon_M^* = 20\%$ has degraded DNN accuracy to the level of random guessing.

Clearly, DNN accuracy reacts vastly differently to individual nonidealities. We observe that nonidealities that effectively add noise to the inputs or outputs – such as ADC and DAC resolution, additive output noise, and S-shaped non-linearity of the ADC – have the largest impact on the DNN accuracy, as normalized to impact on average MVM error. CNNs are the most sensitive DNN topology, while RNNs are the least sensitive (in particular the PTB-LSTM network).

Nonidealities that mostly affect weight-precision (all other nonidealities listed in Fig. 4B), have a much less severe impact on the DNN accuracy. In contrast to additive output noise, such

weight-related nonidealities all scale with the input norm, and thus disappear when no inputs are given. Since it arises from large currents, IR-drop becomes negligible when either inputs or weights are reduced (in either amplitude or occurrence). Such weight-related nonidealities impact CNNs slightly more than RNNs or transformers. In particular, DenseNet-121 with small kernel matrices and a high tile re-use factor seems the most affected by weight disturbances. Fig. 4 shows it is not enough to focus only on weight-related nonidealities, as most previous studies have done, when investigating AIMC.

We use this sensitivity analysis to assess additional nonidealities which our standard AIMC crossbar-model assumes to be perfect. For instance, imperfect device yield – where some fraction of the weight conductances are "stuck" either at zero (PCM reset), at $\hat{g}_{\mathrm{max}}$ (PCM set), or at some intermediate random value – is shown to have the same modest effect on DNN accuracy as other weight-related parameters. Weight asymmetry – a systematic difference in conductance for positive versus negative inputs such that $-w(-|x|) \neq w(|x|)$ – is shown to have only modest impact on DNN accuracy. Interestingly, RNNs and transformers are the models impacted by such polarity-dependent device response, since the ReLU activations used in CNNs cannot create negative inputs. Finally, systematic PCM programming errors – applied once to the conductance values and then remaining constant through repeated MVMs – are shown to have a slightly larger effect than the cycle-to-cycle short-term PCM read-noise that gets redrawn for every MVM.

## 2.5   AIMC robustness of DNN topologies

To extract the specific sensitivities of each individual DNN, we find the threshold value $x^*$ at which each nonideality degrades accuracy to $\mathcal{A}^{1h}(x) = 99\%$, with respect to the standard AIMC crossbar-model. From scans of $\mathcal{A}^{1h}$ as each nonideality is increased (Fig. 5A), we use linear interpolation to identify $x^*$ from the intersection with the dotted line at $\mathcal{A}^{1h} = 99\%$.

The grid in Fig. 5B shows this threshold value $x^*$, for each nonideality and each DNN. For example, considering just total PCM noise, even small increases beyond the current hardware-calibrated values markedly degrade ResNet-18 ($x^* = 1.2\times$ for $\mathcal{A}^{1h} = 99\%$ ), while LSTM-PTB is not affected until this particular nonideality is significantly larger ($x^* = 3\times$). The colors ranging from red to green in Fig. 5 illustrate the relative sensitivity among the DNNs, obtained by scaling $x^*$ linearly between the minimal and maximal values across the 11 DNNs. For many of these nonidealities, yet again RNNs tend to be the most robust, followed by small CNNs on the CIFAR dataset.

Some nonideality parameters can be increased quite dramatically with respect to our standard AIMC crossbar-model baseline. For instance, DAC precision can be lowered from 8bit to 6bit without any retraining, with little accuracy impact across all DNNs – this could produce considerable energy savings and throughput improvement for AIMC designs. Also, IR-drop can be increased beyond than the baseline before becoming problematic, and short-term weight noise could be up to $3\times$ larger, similarly informing future AIMC designs, both with and without PCM devices. While direct examination of Fig. 5 might suggest that IR-drop could be increased by $15\times$ without issue, note that the assumptions inherent in our IR-drop calculations, concerning average rather than instantaneous currents, imply a small safety margin of perhaps $3\times$ (see Methods).

We also estimated the effect of imperfect PCM device yield. Even the least robust model can tolerate 0.27% failed-at-zero devices (stuck in the reset state, at random locations), rising to 3-4% for some of the RNNs. However, DNN accuracies are more sensitive to devices stuck either at random intermediate conductance values or at $\hat{g}_{\mathrm{max}}$ (in the set state). As few as 0.02% such failed devices would already cause a noticeable accuracy drop in some large CNNs. However, our analysis only assumes one pair of conductances per weight — since many existing AIMC designs provide multiple pairs of PCM devices per weight [39, 57], such additional redundancy can potentially counteract such stringent device yield requirements.
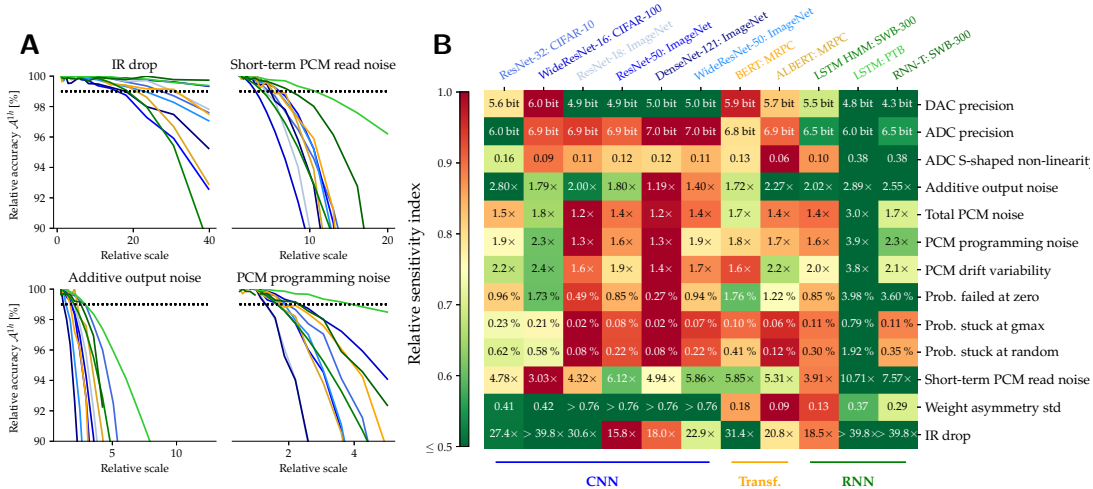
Figure 5: Sensitivities of individual AIMC nonidealities across DNNs. (A) As a single nonideality parameter is increased from the our standard setting, accuracy $\mathcal{A}^{1h}$ eventually drops to 99% (compared to accuracy of the standard AIMC crossbar-model). Four nonidealities are shown, with DNN line-colors matching the textlabel color in (B). (B) Grid shows $x^*$, the threshold value at which that particular nonideality produces $\mathcal{A}^{1h} = 99\%$ (DNN's curve crosses dotted-line in (A)). For instance, reducing DAC precision from 8bit down to 4.8bit, while maintaining all other parameters from the standard AIMC crossbar-model, causes exactly 1% additional accuracy loss in the LSTM-PTB model. Text-label colors at top match the lines in (A); grid colors reflect relative sensitivity index $r_S = \frac{x^* - \min x^*}{\max x^* - \min x^*}$, with min and max values taken across all DNNs. $r_s = 1$ (red) indicates the most sensitive and $r_s = 0$ (green) the least sensitive DNN. RNNs are generally observed to be more robust to AIMC nonidealities than CNNs, even with the limited hyper-parameter tuning available for RNN-T due to its large number of MVM FLOPS and parameters.

## 2.6 Impact of weight distributions on AIMC MVM fidelity

The MVM error of each AIMC tile is affected by the shape of the weight distributions in interesting ways. While weight-clipping might seem disadvantageous, directly programming a very "long-tailed" weight distribution by mapping its largest outlying weight-value to $\hat{g}_{\max}$ can cause even larger problems. Such mappings tend to produce low average output currents which fail to employ the available ADC range, leading to larger MVM errors thanks to ADC quantization, output noise, and other nonidealities that remain stubbornly independent of the reduced output signal-levels.

To show this effect, we calculate the MVM error for different arbitrarily-constructed weight distribution shapes, obtained by sampling the generalized normal distribution,

$$p(x|\mu, \alpha, \beta) = \frac{\beta}{2\alpha\Gamma(1/\beta)} \, e^{-(|x-\mu|/\alpha)^\beta}, \tag{2}$$

where we use $\alpha = 1$ and $\mu = 0$. As $\beta$ increases, this distribution becomes more compact, moving through the Laplace ($\beta = 1$) and normal distributions ($\beta = 2$) along the way (see red curves above Fig. 6A). Fig. 6A shows the MVM error $\epsilon_M$ at 1 hour drift, for weight-values sampled from Eq. 2 as $\beta$ increases from long-tailed ($\beta \leq 1$) to compact (high $\beta$) weight distributions. Here we map weights directly to conductance values, with the maximum weight assigned to $\hat{g}_{\max}$; inputs are uniformly distributed between $(-1, 1)$. MVM error increases rapidly for longer-tailed distributions ($\beta \leq 1$).

One simple measure of a distribution's shape is the kurtosis, obtained by dividing the fourth moment ($\langle (x - \mu)^4 \rangle$) of the distribution by its variance squared ($[\langle (x - \mu)^2 \rangle]^2$). In the plots and the remainder of this section we use the excess kurtosis — defined as the kurtosis minus 3, so that its value is 0 for normal distributions. Since kurtosis increases for long-tailed distributions,
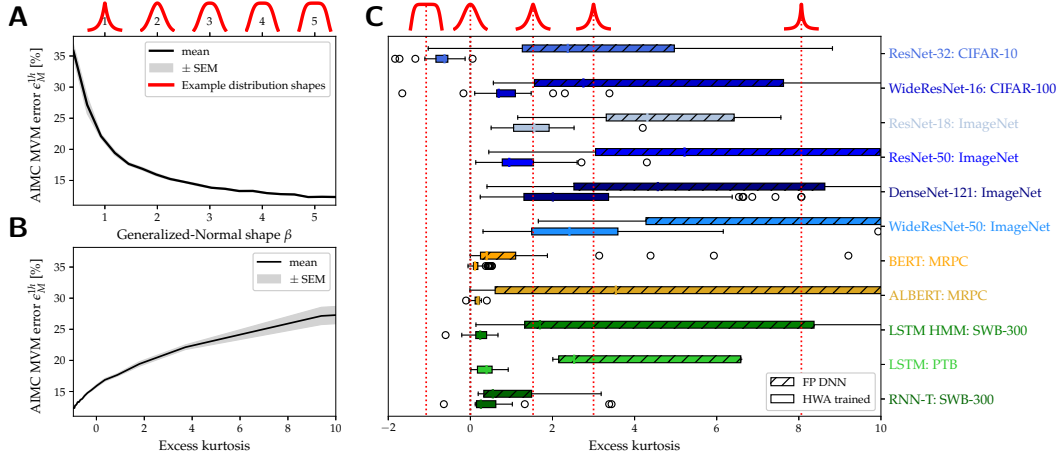
Figure 6: HWA training reduces MVM error by creating more compact conductance distributions. (A) MVM error decreases as constructed conductance distributions, produced by a generalized normal distribution (Eq. 2), are made more compact by increasing $\beta$. Example distributions in red at top show $\beta = 1$ (Laplace distribution), $\beta = 2$ (normal distribution), and even more compact distributions for higher $\beta$. 'SEM' indicates standard-error of the mean. (B) Data from (A) is replotted as a function of the (excess) kurtosis of the distribution. According to the definition of excess kurtosis, a normal distribution (that is $\beta = 2$ in (A)) has a value of 0, and positive or negative values for longer tail distributions (ie. $\beta < 2$) or more compact distributions (ie. $\beta > 2$), respectively. Note that longer tail distributions (large kurtosis) lead to higher MVM error, while more compact distributions (lower kurtosis) reduce MVM error (C) Kurtosis of the conductance values per layer, comparing HWA trained models (solid bars), to FP$_{32}$ weight data scaled by the overall absolute maximum weight (hashed bars). Column-wise scaling, and the tuning of both weights and scaling-parameters during HWA training, help lead to significantly more compact distributions with smaller kurtosis values.

we find that lower kurtosis – and thus more compact weight distributions – means lower MVM error (Fig. 6B).

Fortunately, our HWA training and conductance mapping approach tends to inherently produce more compact conductance distributions, for several different reasons. First, the individual digital scales $\gamma_i$ available for each MVM output (see Eq. 3) are initialized to scale conductances by the absolute maximal value of each weight matrix-column rather than by the overall maximum across the entire weight matrix. With each column individually scaled, the overall conductance distribution becomes more compact than the original weight distribution. During HWA training, these digital scales are optimized – which may lead the system to choose to clip some output columns – and any large weight deviations and outliers created during training are also clipped. Finally, since the AIMC nonidealities cause large weights and outputs to increase the errors that SGD is attempting to correct, HWA training should be expected to drive towards more compact weight distributions during retraining.

Indeed, we find that our HWA training and mapping scheme greatly increases the compactness of the conductance distributions for each layer, as indicated by the kurtosis values shown for our 11 DNN models in Fig. 6C. Hashed bars show kurtosis for direct mapping of the FP$_{32}$ model without HWA training, using a single global digital scale-factor per layer. Solid bars illustrate that our columnwise-scaled and HWA-trained models get mapped into conductance distributions that are significantly more compact, which helps reduce both MVM and DNN error.
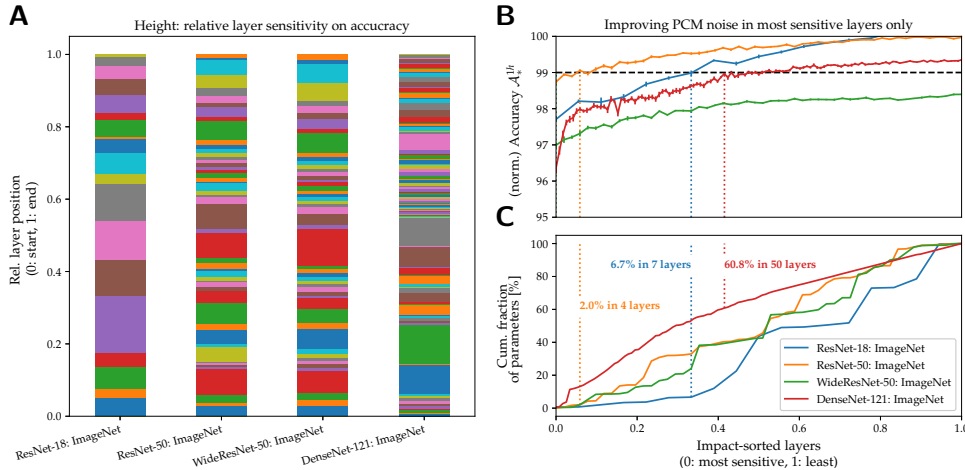
Figure 7: Layer-wise breakdown for ImageNet CNNs. (A) Bar-charts reveal the relative impact that different DNN layers have on AIMC accuracy when the PCM conductances in just that layer are made very noisy (overall PCM noise scale set to 15), while all other layers see only minimal PCM noise (overall noise scale set to 0). The height of each bar-segment, arranged in sequential DNN layer order, corresponds to the relative impact of that layer; colors simply delineate layer boundaries. Note that ResNet-50 and WideResNet-50 have very similar graphs since their layers only differ in width. (B) Accuracy $\mathcal{A}_*^{1h}$ as the most critical layers in these four CNNs are exempted from PCM noise, plotted as the fraction of noise-exempt layers is increased, in order from most-sensitive to least-sensitive. (C) Corresponding cumulative fraction of weight-parameters that are PCM-noise-exempted. For ResNet-50 and ResNet-18, reducing PCM noise in just a few layers (dotted vertical lines) allows an AIMC crossbar-model to achieve iso-accuracy (dashed horizontal line). For DenseNet-121, more than half of the network parameters need to be exempted, while for WideResNet, even the non-PCM nonidealities would need to be improved.

## 2.7 Improving AIMC fidelity of selected layers to reach iso-accuracy in large CNNs

Our results show that larger CNNs, particularly those using the ImageNet database, are the most challenging for AIMC. Even with HWA training, our standard AIMC crossbar-model cannot achieve iso-accuracy for these DNN models (Table 3). Clearly the fidelity of the MVMs must be further improved, either through better materials or through hardware design-choices. For instance, designers could dedicate multiple conductance pairs per weight [44] to reduce PCM programming errors, but at the cost of larger tile area and energy. Or designers could average the results from multiple passes through the tile to reduce the effects of cycle-to-cycle PCM read and additive output noise, but at significant cost to latency, throughput, and energy-efficiency. Given these unpleasant tradeoffs, such approaches should be used as infrequently as possible, ideally only on a small set of DNN layers that really require these extra resources, which can then allow the entire model to achieve iso-accuracy.

Thus we need to determine which of the layers in ImageNet CNNs are the most sensitive to AIMC nonidealities, and then assess whether improving just a small subset of these layers would have sufficient impact. To do this, we sequentially introduce AIMC nonidealities at each layer of the HWA-trained DNNs individually, while turning off all nonidealities in all other layers (using FP$_{32}$ operations on their HWA-trained weight matrices). By repeating this process over the $L$ layers with different overall PCM noise settings, we can determine the sensitivity and relative importance of single layers.

We first rank the layers according to accuracy impact for each DNN by exposing each layer to significant PCM noise with all other layers exempted from noise (Fig. 7A). Then, in order from most- to least-sensitive layer, we introduce this noise-exemption into multiple layers (Fig. 7B),

causing normalized accuracy at 1 hour drift $\mathcal{A}_*^{1h}$ with respect to the $FP_{32}$ model to increase as more and more model-parameters are made noise-exempt (Fig. 7C). Eventually the 99% iso-accuracy is achieved (dashed horizontal line) and then exceeded for most of these models. For Fig. 7A, the one layer being assessed sees $15\times$ the usual PCM noise; for Fig. 7B, the layers not yet PCM-noise-exempted see our standard AIMC crossbar-model. While PCM-noise-exempt layers experience no long-term conductance noise, programming errors, or drift, they still are subject to the same cycle-to-cycle read noise, additive output noise, and DAC/ADC quantization choices in our standard AIMC crossbar-model.

For ResNet-18 and ResNet-50, we find that improving just a few layers can help achieve iso-accuracy ($\mathcal{A}_*^{1h} \geq 99\%$, dashed line in Fig. 7B). This involves only 7% and 2% of the model parameters, respectively (Fig. 7C). Improving MVM fidelity for such a limited number of parameters should prove less costly than across the full DNN. However, for the other two ImageNet DNNs, a considerably larger fraction of the model ($> 60\%$) would need to be improved to reach iso-accuracy. Therefore, these DNNs will either need further advances in either HWA-training or the overall AIMC specifications, in order to support AIMC deployment without significant accuracy drop.

# 3    Discussion

We have introduced a new approach for successfully deploying DNNs onto realistic AIMC inference hardware, at or near iso-accuracy. Our standard AIMC crossbar-model incorporates well-known but hardware-calibrated nonidealities caused by the analog devices, such as read noise, programming errors, and conductance drift. Going well beyond previous studies, our model also includes nonidealities due to MVM circuit integration, such as system noise, DAC and ADC quantization, and dynamically-computed IR-drop. Finally, our model fully addresses the fixed dynamic-range constraints on inputs, weights, and outputs found in all AIMC systems, but previously neglected.

While a few aspects of this study are not directly applicable to hardware designed around non-PCM devices, our standard AIMC crossbar-model and our carefully-designed inference protocols can readily serve as the basic core for studying such systems. The intuition we have developed in terms of how various types of noise affect different families of DNN models is also readily transferable. As such, the present work establishes a baseline that can both guide – and be compared against – future AIMC simulation studies. To help make this even more straight-forward, our standard AIMC crossbar-model has now been incorporated into our open source AIHWKit [66], which is based on the popular ML framework PyTorch [62].

Our AIMC crossbar-model is not a replacement for the detailed circuit simulations essential to hardware verification. We use many simplifications and abstractions of the various AIMC nonidealities, since our goal is quick and relatively realistic functional verification of larger DNN workloads. For instance, we assume noise sources are Gaussian, avoiding physically-modeled distributions that would be more accurate but significantly slower. We also devised a method to rapidly approximate IR-drop which can adjust dynamically with the input. We chose to intentionally ignore static crossbar effects that would change the conductance value systematically [31, 69], since read-write-verify conductance programming can readily adapt to such effects.

Some prior works proposed using on-chip or chip-in-the-loop training methods [36, 83, 31, 88, 85], which can greatly increase the attainable accuracy by addressing the specific fabrication variations found on that particular chip. However, we strongly believe that the time and cost of such individualized preparation is likely to be untenable for widespread deployment. Thus in this paper, we have focused on HWA training that can be general enough to be performed once per model per AIMC chip-family, greatly simplifying the deployment onto individual chips. That said, our HWA training approach could readily be combined with more sophisticated online compensation methods, with on-chip or chip-in-the-loop training, or with more than one device pair used per weight, including optimization of how weights are assigned across these

conductances [50].

Since HWA training is performed in software before deployment, it has no first-order impact on the latency, throughput or energy-efficiency of AIMC hardware. However, as we have shown, HWA training is essential to understanding the tradeoffs between accuracy and these important system performance metrics. For instance, because of the sequential nature of layers of a deep network, shallower but wider layers should generally be preferable for AIMC, since higher utilization of large matrices stored on the crossbar arrays does not significantly change the runtime [25, 65] and helps improve energy-efficiency. In terms of noise robustness, excessively deep DNNs have disadvantages. Among the ImageNet CNNs tested, DenseNet-121 showed the worst accuracy drop from its $FP_{32}$ model, while WideResNet-50 offered the best raw test-error at 1 hour drift (23.82%, versus 24.82% for ResNet-50, Table 3). Together with information on the latency, throughput, and energy-efficiency, this kind of information on available accuracy gain is critical when trying to decide which DNN model to deploy.

A few previous studies have attempted to improve the robustness of DNNs to nonidealities by noise-aware training, where multiplicative or additive Gaussian noise [38, 36]) is added to weights or activations during training. Similarly, other studies seeking to prevent overfitting or to enhance robustness to adversarial attacks have injected noise into standard floating-point training as a regularization technique ([74, 81, 26, 37, 58, 63, 48]). While all these methods qualitatively increase the noise robustness of DNNs, the quantitative benefits on real AIMC can neither be accurately reported nor fully optimized by these studies. Since our HWA approach keeps weights mapped in conductance units, a variety of realistic hardware-relevant constraints can be incorporated in a straightforward manner. These include the complexities of PCM programming, and the shallow input-output ranges, IR-drop and quantization affecting the MVM compute – aspects neglected in most previous studies.

We have tried distilling with the FP model as a teacher (similar to [90]) and found some benefits when HWA training time is limited. However, since the improvements offered by distilling disappeared at longer training times for most DNN models, we mostly report results without distilling. We did find that accuracy with distilling is significantly higher for the HMM Speech LSTM, and these results are shown in Table 3, implying that distilling can be helpful for some DNNs.

Rather than simple Gaussian weight noise [36], we use the expected weight noise distribution characterized from PCM measurements [56]. We find that applying readout-noise on the weights – together with the correct incorporation of injected longer-term programming-noise when modifying the weight matrix during the backward pass – is crucial for achieving AIMC robustness. One drawback of our approach is that this type of noise injection is currently applied only once per mini-batch, which reduces the effectivity of the noise as batch-size increases. One possible improvement would be to sample the weight noise sources multiple times per mini-batch. Such an extension of our methods should further improve the noise robustness of the HWA-trained DNNs.

# 4    Conclusion

In summary, we show that comprehensive hardware-aware (HWA) training can greatly enhance the robustness of a variety of deep neural networks (DNNs) – including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformers – to the unavoidable device and circuit nonidealities of emerging analog in-memory computing (AIMC) accelerators. In five of the 11 models studied, the techniques we introduce lead to software-equivalent accuracy, defined as 99% of the accuracy-performance offered by the original DNN model beyond random guessing. Across all models, HWA-training reduces the worst-case gap in raw model accuracy from 21.81% down to just 2.65%.

Through a systematic sensitivity analysis, we identify the nonidealities that are most critical for maintaining accuracy in future system designs. For instance, we observe that nonidealities that effectively add noise to the inputs or outputs – such as ADC and DAC resolution, additive

output noise, and S-shaped non-linearity of the ADC – have the largest impact on DNN accuracy. We also show that certain DNN topologies, such as RNNs or shallower CNNs, can tolerate more AIMC nonidealities than others.

By making this standard AIMC crossbar-model available in the open-source AIHWKit [66], we make it possible for future advances in HWA training techniques to be readily compared to these results. By pinpointing the measures needed to compensate for imperfect AIMC hardware, the tools we have introduced here enable better understanding and optimization of the tradeoffs between model accuracy and desirable performance characteristics such as latency, throughput, and energy-efficiency. Continued coordination between HWA-training and architectural assessments may even lead to brand-new DNN topologies, specifically designed to maximize the benefits of AIMC hardware — accurate inference at high speed and low power.

# Methods

## AIMC standardized evaluation model

### Affine transform in tile periphery

We assume that each output column of the analog crossbar has a floating-point scale $\alpha_i$ and offset $\beta_i$ available which implement together an affine transformation. We assume that conductances can be linearly mapped to weight values, so that we can normalize the analog weight values from -1 to 1, corresponding to $-\hat{g}_{\max}, \ldots, \hat{g}_{\max}$ (see conductance programming below). This affine transform then maps the column's physical output (e.g. current), as quantized using an ADC into integers within a certain range, to the value expected by the DNN for the next layer (e.g. activation). Note that such ADC conversion using a scale and bias per column is already available in prototypes [39] but has not previously been incorporated into studies on HWA training.

This digital periphery of an analog MVM can thus summarized at

$$\widetilde{y}_i = \beta_i + \alpha\gamma_i \ \text{quant}^q_{b_{\text{out}}} \left( \check{\mathbf{F}} \left( \text{quant}^q_1 \left( \mathbf{x}/\alpha \right) \right) \right), \tag{3}$$

where $\check{\mathbf{F}}$ describes the analog aspects of the AIMC MVM (see "Analog MVM model"), and

$$\text{quant}^q_b(z) \equiv \text{clip}^b_{-b} \left( \text{round} \left( \frac{(2^q - 2)z}{2b} \right) \right), \tag{4}$$

describes linear quantization to $2^q - 1$ values in $-b, \ldots, b$ centered around 0. One bin is discarded to force an odd number of bins on either side of zero. Here $\text{clip}^b_a(x)$ constrains $z$ between minimum $a$ and maximum $b$,

$$\text{clip}^b_a(z) = \begin{cases} z, & \text{if } a < z < b \\ a, & \text{if } z \leq a \\ b, & \text{if } z \geq b. \end{cases} \tag{5}$$

$\alpha$ is a scalar, per-crossbar value which determines the usable input range. This can either be a learned parameter which is then held fixed during inference (static input range), or can depend dynamically on the current input vector $\mathbf{x}$ (dynamic input range). While main results assume a static input range, we examine performance improvements for the dynamic option (Supplemental Material, Sec. A.1).

The scales $\gamma_i$ determine the mapping of conductances to weight values, individually for each crossbar column $i$. During HWA we allow SGD to optimize this parameter, starting from values initialized as described below. $\beta_i$ is used to implement the bias of the MVM, which we implement in digital (FP) precision here. We assume $q = 8$bit quantization, and investigate lower precision as part of our sensitivity analysis.

16

## Dynamic MVM range

A critical feature of our crossbar-model is that it fully encompasses the finite dynamic-range constraints on inputs, weights and outputs that will be present and unavoidable in any real AIMC implementation. Since both input and weights are normalized within $-1, \ldots 1$ (in analog units), our output-bound setting of $b_{\text{out}} = 10$ means that just 10 fully-on inputs, applied to rows containing maximal-value weights, would fully saturate the output. This is a conservative choice that works for modest-size crossbars and for our assumption that positive current contributions (produced by weight and activation pairs of the same sign) and negative contributions (weights and activations have opposite signs) cancel within the array. This mode is energy-efficient and minimizes IR-drops, but requires the ADC to be capable of measuring bipolar currents [39]. If the crossbar is made much larger, or the positive and negative terms are integrated separately, this may increase energy usage and exacerbate IR-drops, but simplify the ADC design. Furthermore, such choices will likely alter the overall dynamic range limitations, calling for a reoptimization of $b_{\text{out}}$.

## Analog MVM model

Our basic model is illustrated in Fig. 3A. The analog MVM $\breve{\mathbf{y}} = \breve{\mathbf{F}}(\breve{\mathbf{x}})$ in Eq. 3 for the quantized, clipped and scaled input vector $\breve{\mathbf{x}} \equiv \text{quant}_1^q(\mathbf{x}/\alpha)$ takes the following general form

$$\breve{y}_i = \sigma^{\text{out}} \xi_i + f_i^{\text{NL}} \left( \Delta \breve{y}_i^{\text{IR-drop}} + \sum_j \left( \breve{w}_{ij}(t_{\text{eval}}) + \sigma^{\text{w}}(\breve{w}_{ij}) \xi_{ij} \right) \breve{x}_j \right), \tag{6}$$

where analog weights $\breve{w}_{ij}(t)$ represent normalized conductances with programming errors, drift, and long-term noise up to time $t_{\text{eval}}$ applied (see 'Weight programming' below). We include a point-wise non-linear functions $f_i^{\text{NL}}(x)$ to support special cases such as ADC nonlinearities; in our standard model, $f_i^{\text{NL}}(x) \equiv x$. Normal random numbers ($\xi_i, \xi_{ij} \sim \mathcal{N}(0, 1)$ are drawn for each MVM, representing additive output noise with standard deviation $\sigma^{\text{out}} = 0.04$, and short-term weight noise $\sigma^{\text{w}}(\breve{w})$ that depends on the current weight values (see 'Short-term PCM read noise'), respectively. Since the analog output values running from $-10, \ldots 10$ get quantized into digital values from $-127, \ldots 127$ (8bit), this choice of $\sigma^{\text{out}} = 0.04$ corresponds to almost exactly half of one ADC quantization bin.

## Weight programming

We adopt a previously-described and -characterized weight-programming and drift model for PCM devices [56]. We assume that the crossbar provides 1 pair of conductances per weight, where the first (second) member of the device-pair is programmed to a conductance between reset (0) and set ($\hat{g}_{\text{max}}$) to handle positive (negative) weights, with the non-active conductance programmed to reset. Only the active conductance is considered in our model. Although recent prototypes support two pairs per weight [57, 39], having only one conductance-pair increases the weight density and thus compute efficiency, and poses a more difficult challenge in terms of accuracy and yield.

Each column $\mathbf{w}_i$ of each weight matrix is mapped to a column of target conductances $\hat{\mathbf{g}}_i$. We first initialize each affine scale coefficient using the maximum weight found in that column, $\gamma_i = \max_j |w_{ij}|$. This allows each weight to be mapped to a scaled target conductance, $\hat{g}_{ij} = \hat{g}_{\text{max}} \frac{w_{ij}}{\gamma_i}$. In our HWA training approach (described below), after this initialization of target conductance and affine scales based on the FP$_{32}$ model weights, we then use SGD to further optimize both the mapped target conductances and scales $\gamma_i$ separately. Table 3 uses this learned weight-to-conductance mapping when evaluating AIMC inference performance.

In a real AIMC system, a positive $\hat{g}$ value gets programmed onto a different physical device than if that particular $\hat{g}$ had been negative. We here assume that only one of the two devices are programmed to particular target conductance whereas the other device is always at reset

conductance ($\hat{g}_{ij} = 0$). In this case, one can simplify and compute the MVM directly with signed conductances as done in our model. The programmed conductances $g_{ij}^{\mathrm{P}}$ differ from the desired target values $\hat{g}_{ij}$ as $g_{ij}^{\mathrm{P}} = \hat{g}_{ij} + \sigma^{\mathrm{P}}(\hat{g}_{ij})\xi$ due to programming noise, assumed to be Gaussian ($\xi \in \mathcal{N}(0,1)$). In turn, the standard deviation of this programming noise depends on the target conductance as

$$\sigma^{\mathrm{P}}(\hat{g}) = c_0 + c_1 \frac{\hat{g}}{\hat{g}_{\mathrm{max}}} + c_2 \frac{\hat{g}^2}{\hat{g}_{\mathrm{max}}^2}, \tag{7}$$

where $c_0 = 0.26348\,\mu\mathrm{S}$, $c_1 = 1.9650\,\mu\mathrm{S}$, and $c_2 = -1.1731\,\mu\mathrm{S}$, as obtained by fitting to extensive PCM hardware data [56].

**Weight drift and read noise**

Once a PCM device is programmed, the device exhibits both conductance drift and $1/f$ (long-term) read noise. As briefly described below, both are modeled in a statistical manner based on measurements of doped-$Ge_2Sb_2Te_5$ (d-GST) mushroom PCMs from a large device array integrated in 90nm CMOS technology [56].

**PCM drift**   PCM conductance drift, attributed to post-programming structural relaxation, follows an empirical relation

$$g^{\mathrm{D}}(t_{\mathrm{eval}}) = g^{\mathrm{P}} \left( \frac{t_{\mathrm{eval}} + t_0}{t_0} \right)^{-\nu}, \tag{8}$$

where $g^{\mathrm{D}}(t_{\mathrm{eval}})$ is the conductance measured at time $t_{\mathrm{eval}}$ after the programming (assumed to complete at $t_0 = 20$s [55]) and $\nu$ is the drift coefficient.

The drift coefficients for each device are assumed to be normally distributed, that is $\nu_{ij} \in \mathcal{N}(\mu_\nu(\hat{g}_{ij}), \sigma_\nu(\hat{g}_{ij}))$, where the mean and standard deviation are empirically determined by fitting to experimental data. The data fits are expressed by a clipped linear function in log-space, that is (with Eq. 5)

$$L(x|a, b, y_{\mathrm{min}}, y_{\mathrm{max}}) \equiv \mathrm{clip}_{y_{\mathrm{min}}}^{y_{\mathrm{max}}} (a \ln x + b) \tag{9}$$

where here $x \equiv \frac{\hat{g}}{\hat{g}_{\mathrm{max}}}$. The parameters for $\mu_\nu$ are given by $a = -0.0155$, $b = 0.0244$, $y_{\mathrm{min}} = 0.049$, and $y_{\mathrm{max}} = 0.1$. For $\sigma_\nu$ the parameter are $a = -0.0125$, $b = -0.0059$, $y_{\mathrm{min}} = 0.008$, and $y_{\mathrm{max}} = 0.045$. The drift coefficient $\nu_{ij}$ thus determined for each device are used to model the conductance at any time $t_{\mathrm{eval}}$ using Eq. 8.

**PCM read noise**   PCM is also known to demonstrate low frequency noise such as random telegraph noise (RTN) and $1/f^\gamma$ noise with $\gamma \in [0.9, 1.1]$. We follow the empirical noise model of [56], which assumes $\gamma = 1$ and arrives at a read noise standard deviation at time $t_{\mathrm{eval}}$ of ([56])

$$\sigma_{\mathrm{read}}(t_{\mathrm{eval}}) = \hat{g}\, Q_s(\hat{g}) \sqrt{\ln \left( \frac{t_{\mathrm{eval}} + T_{\mathrm{read}}}{2\, T_{\mathrm{read}}} \right)}, \tag{10}$$

where $Q_s(\hat{g})$ is measured to be

$$Q_s(\hat{g}) = \mathrm{clip}_0^{c_3} \left( c_1 \left( \frac{\hat{g}}{\hat{g}_{\mathrm{max}}} \right)^{c_2} \right), \tag{11}$$

with $c_1 = 0.0088$, $c_2 = -0.65$, $c_3 = 0.2$.

This read noise is added to the post-drift conductance $g^{\mathrm{D}}(t_{\mathrm{eval}})$ to arrive at the final PCM conductance

$$\tilde{g} = \mathrm{clip}_{\hat{g}_{\mathrm{min}}}^{\infty} \left( g^{\mathrm{D}}(t_{\mathrm{eval}}) + \sigma_{\mathrm{read}}(t_{\mathrm{eval}})\xi \right) \tag{12}$$

where we set $\hat{g}_{\min} = 0$ here and $\xi \sim \mathcal{N}(0,1)$. The weight values $\breve{w}_{ij}$ of the crossbar array for Eq. 6 are then obtained by scaling and combining positive and negative parts

$$\breve{w}_{ij} = \frac{\tilde{g}_{ij}}{\hat{g}_{\max}} \operatorname{sgn} w_{ij} \tag{13}$$

These long-term PCM effects are applied to all weights prior to the evaluation at time $t_{\text{eval}}$ and the weights are subsequently fixed during the evaluation of the test set. Short-term weight noise, redrawn for each MVM, is included separately in Eq. 6 as described in the following paragraph.

**Short-term PCM read noise**   When evaluating the AIMC DNN at a time $t_{\text{eval}}$, the analog weights $\breve{W}$ are established as described in Eq. 13. However, weights are often re-used multiple times during a single input, say across image-pixels in a CNN image or sequence-tokens in an RNNs or transformer model. Here short-term weight noise can cause small but perceptible cycle-to-cycle variations (Fig. 3B).

Modifying the weight matrix at each MVM would be highly inefficient for our HWA training software running on GPUs. To efficiently model such short-term read noise, we use the read noise definition Eq. 10 to set $\sigma^{\text{w}}$ in Eq. 6, but refer the resulting noise to the output $\breve{y}_i$. Assuming zero-mean independent normal distributions, we can sum the variances as

$$\tilde{\sigma}_i^{\text{w}} = \sigma_0^{\text{w}} \sqrt{\sum_j |\breve{w}_{ij}| \, |\breve{x}_j|^2}, \tag{14}$$

implying that the weight dependence of the read noise can be approximated as $\propto \sqrt{|\breve{w}|}$. Thus weight-noise $\sigma^{\text{w}}$ in Eq. 6 effectively adds $\xi_i \tilde{\sigma}_i^{\text{w}}$ (with $\xi_i \sim \mathcal{N}(0,1)$) to the analog output $\breve{y}_i$. The parameter $\sigma_0^{\text{w}}$ can be identified with $c_1 \sqrt{\ln(\frac{\Delta t + t_{\text{r}}}{2 t_{\text{r}}})}$ for read noise accumulated over time-period $\Delta t$ (Eq. 10, [56]). Assuming a read duration of $t_{\text{r}} = 250$ns and approximate waiting time between two consecutive MVMs ($\Delta t$) to be $100\times$ longer, we find $\sigma_0^{\text{w}} \approx 0.0175$.

### Drift compensation

For evaluation times $t_{\text{eval}}$ long after NVM programming, the conductance drift Eq. 8 can be compensated in the digital domain without any expensive re-programming [46, 4]. This can be done by running a number of analog MVMs on some known test inputs $\{\mathbf{x}^k\}$ immediately after weight programming and recording the overall output magnitude as $s_{\text{ref}} = \sum_{ik} |y_i^{(k)}|$. At time $t_{\text{eval}}$, just before beginning inference, the same inputs can be applied to measure $s_{t_{\text{eval}}}$. We then correct the MVM outputs by adjusting the digital $\gamma_i$ (see Eq. 3) by $\frac{s_{\text{ref}}}{s_{t_{\text{eval}}}}$ to accommodate the average conductance decrease due to drift. We assume one global drift compensation applied to all columns, although this could be done individually at each column if $s_{\text{ref}}|_i$ can be measured sufficiently accurately. Other more sophisticated drift compensation and adaptive refresh methods including in-memory re-training could potentially be applied as well (e.g. [36]).

### Crossbar tile size

The NVM crossbars available on an AIMC chip are of finite size, typically ranging from $256 \times 256$ (e.g. [39]) to $512 \times 512$ (e.g. [57]). We assume a tile-size of $512 \times 512$, and assume that enough crossbars are available to support separate crossbars for each weight matrix. Any weight matrix with input dimension $> 512$ is divided into roughly equal parts for programming on as many tiles necessary. Partially-used tiles have weights are situated at the bottom of the crossbar, to minimize interference and potential IR drop, and unused inputs are clamped to zero.

Each tile computes an MVM (Eq. 6) using its own periphery (Eq. 3). Inter-tile summation is performed at FP precision (FP16), after affine-scaling but before being passed to subsequent digital compute such as activation functions. Because our AIMC nonidealities have no dependencies across output columns, the HWA training code does not need to explicitly break the

compute along the output-dimension into tile-sized chunks. This helps the simulations run more efficiently on GPUs.

## IR-drop

Ideally, the voltage along each long bitline in the crossbar would remain constant, so that conductances with the same value could contribute the same current, whether in the farthest or nearest row from where peripheral circuitry is holding the bitline voltage and measuring currents. In a physical crossbar, however, IR-drops imposed by finite wire resistance cause the bitline voltage to vary [13], especially as instantaneous currents get large. To keep the simulation-time reasonable, we make a number of approximations when modeling this effect. IR-drop is modeled independently for each crossbar-column, because any column-to-column differences will be implicitly corrected (to first order) when programming the weight with an appropriate read-write-verify scheme.

However, within each crossbar column, the current contributed by each weight depends on the local bitline voltage, which in turn depends on the other currents being generated elsewhere along the column by that particular input vector. This situation will evolve throughout the integration period due to the pulse-length modulation of those inputs as well as any resulting transients, including the response of the peripheral circuit establishing the bitline voltage. Here, for simplicity and speed of computation for large DNNs, we only consider the average integration current.

The steady-state bitline voltages $\bar{v}_i$ can be computed by solving the equation system

$$(\bar{v}_{i+1} - \bar{v}_i)\, g_w + g_i^+(v_i^+ - \bar{v}_i) = (\bar{v}_i - \bar{v}_{i-1})\, g_w + g_i^-\left(\bar{v}_i - v_i^-\right) \tag{15}$$

where $g_w$ is the wire conductance between the crosspoint nodes and $g_i^{+/-}$ the weight programmed onto either the positive or negative conductance (with the other programmed into the reset condition, $g = 0$). The individual input voltages, $v_i^-$ and $v_i^+$ of spatially-ordered inputs $i$, are linearly prorated from the supply voltages $(v_{\mathrm{ref}} \pm V_{\mathrm{read}})$ to represent the time-averaged current. The analog output current $\breve{y}$ located at location $i = 0$ is given by $g_w\,(\bar{v}_0 - v_{\mathrm{ref}})$, with $V_{\mathrm{read}} = 0.2\mathrm{V}$.

This linear system (Eq. 15) can be solved by inverting the unique coefficient matrix produced by a given input vector. To speed up the simulation and avoid inverting a $512 \times 512$ matrix for each MVM, we further approximate the solution with a quadratic equation. Thus, in our analog MVM (Eq. 6), the IR-drop amount is computed from the normalized weights and inputs by

$$a_i \;\equiv\; \gamma n \sum_j |\breve{w}_{ij}||\breve{x}_j| \tag{16}$$

$$c_i \;\equiv\; 0.05\, a_i^3 - 0.2 a_i^2 + 0.5 a_i \tag{17}$$

$$\Delta \breve{y}_i^{\mathrm{IR\text{-}drop}} \;\equiv\; -c_i \sum_j \breve{w}_{ij}\breve{x}_j \left(1 - (1 - \frac{j}{n})^2\right), \tag{18}$$

where $\gamma$ is the unitless product of the wire resistance between adjacent cross-points (assumed $0.35\ \Omega$) and the maximal (set) conductance of the device ($g_{\mathrm{max}} = 5\mu\mathrm{S}$), and $n$ is the number of cross-points occupied by the weight matrix. We assume that smaller weight matrices are located at the lower edge of the crossbar to avoid excess IR-drop. We use Eq. 18 to dynamically approximate the IR-drop across the 512 input channels in Eq. 6 when computing normalized MVM outputs $\widetilde{y}$ in all our results. Multiplying these normalized outputs by $g_{\mathrm{max}}V_{\mathrm{read}}$ produces the (time-averaged) physical output currents. To amplify these IR-effects for the sensitivity analysis (Fig. 4), we simply multiply the IR-drop error $\Delta \breve{y}_i^{\mathrm{IR\text{-}drop}}$ by a varying scaling factor.

For large inputs where current is flowing throughout the integration window, our estimations using time-averaged current are quite accurate. However, for small inputs where much of the current-flow occurs in a small portion of the integration window, instantaneous and average

currents differ strongly, and IR-drop will be underestimated. We find that for a Normal distributed weight matrix and random but correlated inputs (as in Fig. 3E), IR-drop deviations are underestimated by roughly a factor of 5. Unfortunately, similar conditions arise across many of our DNNs. Fortunately, our sensitivity analysis (Fig. 4) finds that scaling our time-averaged IR-drop approximation by a factor of $> 10\times$ does not significantly impact the accuracy of the DNNs, so we can still conclude that DNNs are reasonably robust to IR-drop, albeit by a modest rather than large safety margin. Since IR-drop depends heavily on both on the hardware design (crossbar size, wire resistances, and absolute device conductances) and on the input and weight distributions, detailed circuit-based simulations using the intended workload(s) will remain a critical part of assessing new hardware designs.

## Additional nonlinearities for sensitivity analysis

### PCM device yield

Emerging memory devices such as PCM exhibit imperfect yield, and some fraction of the devices in a given crossbar array will simply not switch properly [40, 14]. PCM devices can end up stuck-at-set ($\hat{g}_{\mathrm{max}}$), stuck-at-reset (conductance set to 0) and stuck-at-random (stuck somewhere between 0 and $\hat{g}_{\mathrm{max}}$). In our sensitivity analysis (Fig. 4), we vary the fraction of failed devices and randomly-select their locations.

### S-shaped ADC output non-linearity

The output level might gradually saturate more gradually than the desired linear response due to non-linearity in the ADC [80, 39]. To estimate the impact of this for our sensitivity analysis (Fig. 4), we define $f_i^{\mathrm{NL}}$ in Eq. 6 with

$$f_i^{\mathrm{NL}}(z) \equiv \left( 1 + \frac{2}{d_{\mathrm{out}}} \sum_{k=1}^{d_{\mathrm{out}}} |\zeta_k| \right)^2 \frac{z}{1 + |\zeta_i z|}, \tag{19}$$

which models a S-shaped saturation with variable slope scaled to approximately cover the full output range. Each of the $d_{\mathrm{out}}$ outputs has an independent ADC and thus a slightly different (pre-determined) shape, $\zeta_i = \mu_\zeta (1 + \sigma_\zeta \xi)$ with $\xi \sim \mathcal{N}(0, 1)$ and $\mu_\zeta = \frac{1}{4}$. $\sigma_\zeta$ is only varied in the sensitivity analysis ("ADC S-shaped nonlinearity"); for our standard AIMC crossbar-model, $\mu_\zeta$ and $\sigma_\zeta$ are both set to 0, causing $f_i^{\mathrm{NL}}(z) = z$.

### PCM polarity

Depending on the hardware and unit-cell design, positive and negative inputs might not create perfectly symmetric read-currents. The measured conductance of a PCM-device can depend on whether read-current passes from top to bottom electrode, or vice-versa. This read-polarity dependence can cause weights to appear systematically altered for negative inputs as compared to positive inputs. Although the average effect can be corrected by adjusting read voltages, device-to-device or conductance-dependent variations can remain. To model this effect in our sensitivity analysis, we separate positive and negative inputs into two phases (setting a negative input to 0 in the positive phase and vice versa), and scale each weight in the negative phase by $(1 + a_{ij})$ where $a_{ij} \sim \mathcal{N}(0, \sigma_a)$. We then vary this new nonideality parameter $\sigma_a$ as "weight asymmetry std."

## MVM error calculation

To quantify the fidelity of the analog MVM, we calculate the expected deviation of the analog MVM as compared to the ideal MVM as MVM-error $\epsilon_M$, defined by the relative normalized

deviations (see e.g. [8])

$$\epsilon_M(W, \{\mathbf{x}_k\}) = \frac{\langle ||\mathbf{y}_k - \widetilde{\mathbf{y}}_k||_2 \rangle_k}{\langle ||\mathbf{y}_k||_2 \rangle_k}, \tag{20}$$

where $\mathbf{y}_k = W\mathbf{x}_k$ is the ideal MVM output to input vector $\mathbf{x}_k$ using matrix $W$, and $\widetilde{\mathbf{y}}$ is the actual AIMC output considering all hardware-related nonidealities as defined in Eq. 3.

The MVM-error is obviously zero if the AIMC is equal to the ideal outcome, but otherwise it depends on both the particular weight matrix $W$ and set of input vectors $\mathbf{x}_k$ used to estimate Eq. 20. To best reflect the impact of the nonidealities on the DNN, inputs $\mathbf{x}_k$ should ideally be taken from the distribution of actual input activation vectors, and $W$ should be the target weight matrix, for the specific DNN layer in question.

However, to quantify the MVM-error independent of the DNN in question, we calculate the standard MVM-error $\epsilon_M^*$ by using normal distributed weights, $w_{ij} \sim \mathcal{N}(0, 0.246)$ and uniform inputs $x_i \sim \mathcal{U}(-1, 1)$ with a tile size of $512 \times 512$. For our standard AIMC crossbar-model as described in 'AIMC standardized evaluation model', the standard MVM error is $\epsilon_M^* = 15\%$ (not considering drift).

## AIMC hardware-aware DNN training

Robustness to the nonidealities of AIMC inference hardware can be improved by hardware-aware (HWA) training — a DNN re-training method that applies expected nonidealities to the forward pass of the SGD, with the backward pass performed using regular FP precision.

Our HWA training approach is to use the general form of the expected analog MVM nonidealities as described in Eq. 6 together with the injection of the expected programming errors (but without any conductance drift). Further, we use the HWA training step to also establish the digital peripheral parameters of Eq. 3, in particular the static input range $\alpha$ (see 'Learning the input range') and the weight-to-conductance mapping $\gamma_i$ (see 'Learning of weight-to-conductance conversion factors'). Additionally, we found that ramping up the injected programming error strength (see 'Re-training with weight noise injection'), fixed scales and individual learning rates per tile (see 'Learning of weight-to-conductance conversion factors'), weight clipping (see 'Weight mapping and clipping') and distilling (see 'Distilling with floating-point teacher') improved the robustness and achievable accuracy in the presence of AIMC nonidealities.

In general, the HWA training starts from an already FP-trained DNN, and hyper-parameters (learning rate, injected noise strength) are optimized. We verified the effectiveness of our new HWA training approach on the very same DNNs used in a previous study [36] and found, on average, a $> 10\%$ decrease in AIMC test error for long $t_{\text{eval}}$ times. This directly indicates the improvement of our approach over previous methods (see Table A.2).

In the following paragraphs, our new HWA training methods are presented in more detail.

### Re-training with weight noise injection

Injecting noise to improve robustness to nonidealities was suggested by a number of studies (e.g. [36, 24, 38]), and has been one of the hallmarks of HWA training for AIMC. In previous studies, noise has been injected in multiple ways, such as output [24, 36], input [36], or weight noise [36, 38]. Different types of weight noise distributions have been used, such as additive (scaled by the current maximal weight [36]) or multiplicative [38] Gaussian.

Methods for injecting weight noise have differed across previous studies. For instance, Joshi et al. [36] added newly drawn Gaussian weight noise to the weight matrix reversibly for each image input (not mini-batch) only during the forward pass (and not during backward pass which was done with the actual weight matrix). However, it is more mathematically-correct to also apply these same weight perturbations during the backward pass (but not to the reference weights to which updates are applied), as is commonly done for weight regularization techniques such as drop-connect [82]. Furthermore, although the exact noise injection method (input, output, or weight noise) does not seem to matter much [36], generic additive Gaussian noise does not

conform with the expected AIMC noise structure. For instance, PCM programming errors are actually conductance-value dependent and not just additive.

Here, we improve on the earlier approaches in the following ways: First, rather than just a generic noise term, we apply all expected nonidealities and hardware design choices (given by Eq. 6) into the HWA retraining. This includes dynamic range limitations, system noise, and analog-digital conversions — all previously ignored. We inject weight noise in a mathematically-consistent way to both forward and backward passes, redrawing from random distributions once per mini-batch. We draw the weight noise from the (scaled) expected programming error distribution (see Eq. 7) instead of using generic additive or multiplicative Gaussian distributions. Finally, the scale of the injected weight noise is ramped up linearly over a number of epochs, which we found to improve the HWA training. See Sec. A.2 (Supplemental Material) for the detailed hyper-parameters and noise settings used for each DNN.

## Learning of weight-to-conductance conversion factors

To achieve a good weight-to-conductance conversion, we train the $\gamma_i$ scale-factors in Eq. 3 using SGD. To improve the training of CNN models, it is beneficial to represent these scale-factors by $\gamma_i = \tilde{\gamma}_i \kappa$, where both the column-wise $\tilde{\gamma}_i$ and per-tile $\kappa$ factors can be learned. We treat the learning of either factor as a hyper-parameter for a particular DNN. In case of not learning, $\gamma_i$ is initialized by the weight mapping described below (see 'Weight mapping and clipping') and $\kappa$ is set to 1.

In case of CNNs, where the matrix-sizes vary widely, the learned values $\tilde{\gamma}_i$ are uniquely scaled for each weight matrix by a fixed $c_{\text{aws}}$ value, which re-scales the learning rates per tile so that the trained parameters can all have similar magnitude $\approx 1$. This auto-weight scaling factor, $c_{\text{aws}}$, is set to the value suggested by the Xavier weight initialization [22, 64], $c_{\text{aws}} = \sqrt{\frac{3}{n}}$, where $n$ is the input dimension of the weight matrix.

If $\kappa$ is learned, we encourage the learning of larger outputs and weights by down-scaling the output range to $[-1, 1]$ which typically improves the signal-to-noise ratio, thus $\kappa = \frac{\tilde{\kappa}}{b_{\text{out}}}$. Here $b_{\text{out}}$ is the fixed output bound of Eq. 3, and $\tilde{\kappa}$ is a per-tile learnable scalar which is initialized to $b_{\text{out}}$ (and is subject to weight decay).

Note that during inference evaluation, the digital periphery can simply apply one scale-factor per output-column, since the various scale-factors described above can be re-combined after the completion of HWA training.

## Weight mapping and clipping

Since we use the output scales $\gamma_i$ to keep the analog weights $\breve{w}_{ij}$ of Eq. 6 mapped in (normalized) conductance units (within $-1, \ldots, 1$), the FP weights $w_{ij}$ of the trained DNN need to be mapped to conductances before initiating HWA training. For that we set initially

$$\breve{w}_{ij} \quad \leftarrow \quad \frac{w_{ij}}{\max_j |w_{ij}|} \tag{21}$$

$$\gamma_i \quad \leftarrow \quad \max_j |w_{ij}| \tag{22}$$

so that $\gamma_i \breve{w}_{ij} = w_{ij}$.

We keep training from creating excessively large analog weights. $\breve{w}$, by clipping after each update to this same range. In some cases (see Supplemental Information), we encourage learning of larger analog weights to maintain signal-to-noise ratio by remapping weights according to Eq. 21 once every epoch.

## Learning the input range

The input range clipping bound $c_{\text{input}}$ in Eq. 3 is learned during HWA training. To encourage a smaller clipping value (and thus a more compact input distribution), a decay is introduced. To

augment the gradient update for the clipping bound, we scale gradient updates by the current bound value. For small data sets (such as for transformer fine-tuning tasks), the HWA training is too short to learn the clipping bound value from scratch. In such cases, we initialize $c_{\text{input}}$ to the average absolute maximal value of the input vectors over a number of mini-batches before starting HWA training, subject to a cap (nominally $\max(c_{\text{input}}) = 10$).

### Distilling with floating-point teacher

If the model output dimension is large, such as for the LSTM models with large vocabulary size, the HWA training greatly benefits from distilling with the FP model. In knowledge distillation [29], an already trained "teacher" model augments the usual one-hot labels with expected class probabilities, which can drive a "student" model to a good solution more rapidly than when training only with the one-hot label vectors. We use the distilling applied at the last layer, with the FP model without any AIMC nonidealities as the teacher and the HWA training as described above as the student. The temperature controlling the distribution of pseudo-probabilities was fixed to 10, and training loss was weighted by a mixture of 75% from the distillation and 25% from the regular loss.

## HWA training experiments

We applied and optimized the HWA training process described in this section to a variety of AI workloads – including text prediction, speech-to-text translation, and image classification – as listed in Table 1. In general, our HWA training approach addressed these DNNs similarly, since a too DNN-specific re-training approach would be impractical. In Sec. A.2, we detail any specific differences used in the HWA training of these DNNs, including learning-rates and injected noise strength. We select the last available rather than the best checkpoint, and we repeat experiments multiple times and average the results to obtain repeatable results.

# 5 Acknowledgment

# References

[1] A. Agrawal, S. K. Lee, J. Silberman, M. Ziegler, M. Kang, S. Venkataramani, N. Cao, B. Fleischer, M. Guillorn, M. Cohen, S. Mueller, J. Oh, M. Lutz, J. Jung, S. Koswatta, C. Zhou, V. Zalani, J. Bonanno, R. Casatuta, C.-Y. Chen, J. Choi, H. Haynie, A. Herbert, R. Jain, M. Kar, K.-H. Kim, Y. Li, Z. Ren, S. Rider, M. Schaal, K. Schelm, M. Scheuermann, X. Sun, H. Tran, N. Wang, W. Wang, X. Zhang, V. Shah, B. Curran, V. Srinivasan, P.-F. Lu, S. Shukla, L. Chang, and K. Gopalakrishnan. A 7nm 4-core AI chip with 25.6 TFLOPS hybrid FP8 training, 102.4 TOPS INT4 inference and workload-aware throttling. In *IEEE International Solid-State Circuits Conference (ISSCC)*, volume 64, pages 144–146, 2021.

[2] A. Agrawal, S. M. Mueller, B. M. Fleischer, X. Sun, N. Wang, J. Choi, and K. Gopalakrishnan. Dlfloat: A 16-b floating point format designed for deep learning training and inference. In *2019 IEEE 26th Symposium on Computer Arithmetic (ARITH)*, pages 92–95, 2019.

[3] J. Albericio, P. Judd, T. Hetherington, T. Aamodt, N. E. Jerger, and A. Moshovos. Cnvlutin: Ineffectual-neuron-free deep neural network computing. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pages 1–13, 2016.

[4] S. Ambrogio, M. Gallot, K. Spoon, H. Tsai, C. Mackin, M. Wesson, S. Kariyappa, P. Narayanan, C.-C. Liu, A. Kumar, A. Chen, and G. W. Burr. Reducing the impact

of phase-change memory conductance drift on the inference of large-scale hardware neural networks. In *IEEE International Electron Devices Meeting*, pages 1–4, 2019.

[5] S. Ambrogio, P. Narayanan, H. Tsai, R. M. Shelby, I. Boybat, C. di Nolfo, S. Sidler, M. Giordano, M. Bodini, N. C. P. Farinha, B. Killeen, C. Cheng, Y. Jaoudi, and G. W. Burr. Equivalent-accuracy neuromorphic hardware acceleration of neural network training using analog memory. *Nature*, 558(7708):60–67, 2018.

[6] M. Boniardi, D. Ielmini, S. Lavizzari, A. L. Lacaita, A. Redaelli, and A. Pirovano. Statistics of resistance drift due to structural relaxation in phase-change memory arrays. *IEEE Transactions on Electron Devices*, 57(10):2690–2696, 2010.

[7] R. L. Bruce, S. Ghazi Sarwat, I. Boybat, C. W. Cheng, W. Kim, S. R. Nandakumar, C. Mackin, T. Philip, Z. Liu, K. Brew, N. Gong, I. Ok, P. Adusumilli, K. Spoon, S. Ambrogio, B. Kersting, T. Bohnstingl, M. Le Gallo, A. Simon, N. Li, I. Saraf, J. P. Han, L. Gignac, J. M. Papalia, T. Yamashita, N. Saulnier, G. W. Burr, H. Tsai, A. Sebastian, V. Narayanan, and M. Brightsky. Mushroom-Type phase change memory with projection liner: An array-level demonstration of conductance drift and noise mitigation. In *IEEE International Reliability Physics Symposium Proceedings*, volume 2021-March, pages 1–6, 2021.

[8] J. Büchel, A. Vasilopoulos, B. Kersting, F. Odermatt, K. Brew, I. Ok, S. Choi, I. Saraf, V. Chan, T. Philip, et al. Gradient descent-based programming of analog in-memory computing cores. In *2022 International Electron Devices Meeting (IEDM)*, pages 33–1. IEEE, 2022.

[9] G. W. Burr, M. J. BrightSky, A. Sebastian, H.-Y. Cheng, J.-Y. Wu, S. Kim, N. E. Sosa, N. Papandreou, H.-L. Lung, H. Pozidis, E. Eleftheriou, and C. H. Lam. Recent Progress in Phase-Change Memory Technology. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 6(2):146–162, 2016.

[10] G. W. Burr, A. Sebastian, T. Ando, and W. Haensch. Ohm's law plus Kirchhoff's current law equals better AI. *IEEE Spectrum*, 58(12):44–49, 2021.

[11] G. W. Burr, R. M. Shelby, A. Sebastian, S. Kim, S. Kim, S. Sidler, K. Virwani, M. Ishii, P. Narayanan, A. Fumarola, L. L. Sanches, I. Boybat, M. Le Gallo, K. Moon, J. Woo, H. Hwang, and Y. Leblebici. Neuromorphic computing using non-volatile memory. *Advances in Physics: X*, 2(1):89–124, 2017.

[12] H.-Y. Chang, P. Narayanan, S. C. Lewis, N. C. P. Farinha, K. Hosokawa, C. Mackin, H. Tsai, S. Ambrogio, A. Chen, and G. W. Burr. AI hardware acceleration with analog memory: micro-architectures for low energy at high speed. *IBM Journal of Research and Development*, 63(6):1–14, 2019.

[13] A. Chen. A comprehensive crossbar array model with solutions for line resistance and nonlinear device characteristics. *IEEE Transactions on Electron Devices*, 60(4):1318–1326, 2013.

[14] L. Chen, J. Li, Y. Chen, Q. Deng, J. Shen, X. Liang, and L. Jiang. Accelerator-friendly neural-network training: Learning variations and defects in rram crossbar. In *Design, Automation Test in Europe Conference Exhibition (DATE), 2017*, pages 19–24, 2017.

[15] J. Choi, Z. Wang, S. Venkataramani, P. I.-J. Chuang, V. Srinivasan, and K. Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural networks, 2018.

[16] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1, 2016.

[17] X. Cui, V. Goel, and G. Saon. Embedding-Based Speaker Adaptive Training of Deep Neural Networks. In *Proc. Interspeech 2017*, pages 122–126, 2017.

[18] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet. Front-end factor analysis for speaker verification. In *IEEE Transactions on Audio, Speech, and Language Processing*, volume 19, pages 788–798, 2011.

[19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[20] L. Fick, S. Skrzyniarz, M. Parikh, M. B. Henry, and D. Fick. Analog matrix processor for edge ai real-time video analytics. In *2022 IEEE International Solid-State Circuits Conference (ISSCC)*, volume 65, pages 260–262. IEEE, 2022.

[21] E. J. Fuller, S. T. Keene, A. Melianas, Z. Wang, S. Agarwal, Y. Li, Y. Tuchman, C. D. James, M. J. Marinella, J. J. Yang, A. Salleo, and A. A. Talin. Parallel programming of an ionic floating-gate memory array for scalable neuromorphic computing. *Science*, 364(6440):570–574, 2019.

[22] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010.

[23] T. Gokmen, M. Onen, and W. Haensch. Training deep convolutional neural networks with resistive cross-point devices. *Frontiers in neuroscience*, 11:538, 2017.

[24] T. Gokmen, M. J. Rasch, and W. Haensch. The marriage of training and inference for scaled deep learning analog hardware. In *2019 IEEE International Electron Devices Meeting (IEDM)*, pages 22–3. IEEE, 2019.

[25] T. Gokmen and Y. Vlasov. Acceleration of deep neural network training with resistive cross-point devices: Design considerations. *Frontiers in neuroscience*, 10:333, 2016.

[26] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks, 2013.

[27] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, 2016.

[28] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[29] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015.

[30] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

[31] S. Jain, A. Sengupta, K. Roy, and A. Raghunathan. Rxnn: A framework for evaluating deep neural networks on resistive crossbars. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 40(2):326–338, 2021.

[32] S. Jain, H. Tsai, C.-T. Chen, R. Muralidhar, I. Boybat, M. M. Frank, S. Woźniak, M. Stanisavljevic, P. Adusumilli, P. Narayanan, et al. A heterogeneous and programmable compute-in-memory accelerator architecture for analog-ai using dense 2-d mesh. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 31(1):114–127, 2023.

[33] J.-W. Jang, S. Park, G. W. Burr, H. Hwang, and Y.-H. Jeong. Optimization of Conductance Change in $Pr_{1-x}Ca_xMnO_3$-Based Synaptic Devices for Neuromorphic Systems. *IEEE Elec. Dev. Lett.*, 36(5):457–459, 2015.

[34] J.-W. Jang, S. Park, Y.-H. Jeong, and H. Hwang. ReRAM-based Synaptic Device for Neuromorphic Computing. In *IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 1054–1057, 2014.

[35] H. Jia, H. Valavi, Y. Tang, J. Zhang, and N. Verma. A programmable heterogeneous microprocessor based on bit-scalable in-memory computing. *IEEE Journal of Solid-State Circuits*, 55(9):2609–2621, 2020.

[36] V. Joshi, M. Le Gallo, S. Haefeli, I. Boybat, S. R. Nandakumar, C. Piveteau, M. Dazzi, B. Rajendran, A. Sebastian, and E. Eleftheriou. Accurate deep neural network inference using computational phase-change memory. *Nature Communications*, 11(2473):1–13, 2020.

[37] G. Kang, J. Li, and D. Tao. Shakeout: A new regularized deep neural network training scheme. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI'16, page 1751–1757. AAAI Press, 2016.

[38] S. Kariyappa, H. Tsai, K. Spoon, S. Ambrogio, P. Narayanan, C. Mackin, A. Chen, M. Qureshi, and G. W. Burr. Noise-Resilient DNN: Tolerating Noise in PCM-Based AI Accelerators via Noise-Aware Training. *IEEE Transactions on Electron Devices*, 68(9):1–7, 2021.

[39] R. Khaddam-Aljameh, M. Stanisavljevic, J. F. Mas, G. Karunaratne, M. Braendli, F. Liu, A. Singh, S. M. Müller, U. Egger, A. Petropoulos, T. Antonakopoulos, K. Brew, S. Choi, I. Ok, F. L. Lie, N. Saulnier, V. Chan, I. Ahsan, V. Narayanan, S. R. Nandakumar, M. L. Gallo, P. A. Francese, A. Sebastian, and E. Eleftheriou. HERMES Core – A 14nm CMOS and PCM-based In-Memory Compute Core using an array of 300ps/LSB Linearized CCO-based ADCs and local digital processing. In *Symposium on VLSI Circuits*, 2021.

[40] W. Kim, M. BrightSky, T. Masuda, N. Sosa, S. Kim, R. Bruce, F. Carta, G. Fraczak, H. Y. Cheng, A. Ray, Y. Zhu, H. L. Lung, K. Suu, and C. Lam. Ald-based confined pcm with a metallic liner toward unlimited endurance. In *2016 IEEE International Electron Devices Meeting (IEDM)*, pages 4.2.1–4.2.4, 2016.

[41] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur. Audio augmentation for speech recognition. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3586–3589, 2015.

[42] R. Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.

[43] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images, 2009.

[44] M. Le Gallo, S. Nandakumar, L. Ciric, I. Boybat, R. Khaddam-Aljameh, C. Mackin, and A. Sebastian. Precision of bit slicing with in-memory computing based on analog phase-change memory crossbars. *Neuromorphic Computing and Engineering*, 2(1):014009, 2022.

[45] M. Le Gallo and A. Sebastian. An overview of phase-change memory device physics. *Journal of Physics D: Applied Physics*, 53(21):213002, 2020.

[46] M. Le Gallo, A. Sebastian, G. Cherubini, H. Giefers, and E. Eleftheriou. Compressed sensing with approximate message passing using in-memory computing. *IEEE Transactions on Electron Devices*, 65(10):4304–4312, 2018.

[47] Y. Li, S. Kim, X. Sun, P. Solomon, T. Gokmen, H. Tsai, S. Koswatta, Z. Ren, R. Mo, C. C. Yeh, et al. Capacitor-based cross-point array for analog neural network with record symmetry and linearity. In *2018 IEEE Symposium on VLSI Technology*, pages 25–26. IEEE, 2018.

[48] Y. Li and F. Liu. Adaptive gaussian noise injection regularization for neural networks. *Lecture Notes in Computer Science*, page 176–189, 2020.

[49] S. Lim, M. Kwak, and H. Hwang. Improved Synaptic Behavior of CBRAM Using Internal Voltage Divider for Neuromorphic Systems. *IEEE Transactions on Electron Devices*, 65(9):3976–3981, 2018.

[50] C. Mackin, M. J. Rasch, A. Chen, J. Timcheck, R. L. Bruce, N. Li, P. Narayanan, S. Ambrogio, M. Le Gallo, S. Nandakumar, et al. Optimised weight programming for analogue memory-based deep neural networks. *Nature communications*, 13(1):1–12, 2022.

[51] C. Mackin, H. Tsai, S. Ambrogio, P. Narayanan, A. Chen, and G. W. Burr. Weight Programming in DNN Analog Hardware Accelerators in the Presence of NVM Variability. *Advanced Electronic Materials*, 5(9):1900026, 2019.

[52] F. Merrikh-Bayat, X. Guo, M. Klachko, M. Prezioso, K. K. Likharev, and D. B. Strukov. High-performance mixed-signal neurocomputing with nanoscale floating-gate memory cell arrays. *IEEE Trans. Neur. Netw. Learn. Sys.*, 2017.

[53] B. Murmann. Mixed-signal computing for deep neural network inference. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 29(1):3–13, 2020.

[54] M. Nagel, M. Fournarakis, R. A. Amjad, Y. Bondarenko, M. van Baalen, and T. Blankevoort. A white paper on neural network quantization. *arXiv preprint arXiv:2106.08295*, 2021.

[55] S. R. Nandakumar, I. Boybat, J.-P. Han, S. Ambrogio, P. Adusumilli, R. L. Bruce, M. BrightSky, M. J. Rasch, M. Le Gallo, and A. Sebastian. Precision of synaptic weights programmed in phase-change memory devices for deep learning inference. In *IEEE International Electron Devices Meeting (IEDM)*, pages 1–4, 2020.

[56] S. R. Nandakumar, I. Boybat, V. Joshi, C. Piveteau, M. Le Gallo, B. Rajendran, A. Sebastian, and E. Eleftheriou. Phase-change memory models for deep learning training and inference. *IEEE International Conference on Electronics, Circuits and Systems*, pages 727–730, 2019.

[57] P. Narayanan, S. Ambrogio, A. Okazaki, K. Hosokawa, H. Tsai, A. Nomura, T. Yasuda, C. Mackin, S. C. Lewis, A. Friz, M. Ishii, Y. Kohda, H. Mori, K. Spoon, R. Khaddam-Aljameh, N. Saulnier, M. Bergendahl, J. Demarest, K. W. Brew, V. Chan, S. Choi, I. Ok, I. Ahsan, F. L. Lie, W. Haensch, V. Narayanan, and G. W. Burr. Fully on-chip mac at 14nm enabled by accurate row-wise programming of pcm-based weights and parallel vector-transport in duration-format. In *2021 Symposium on VLSI Technology*, pages 1–2, 2021.

[58] H. Noh, T. You, J. Mun, and B. Han. Regularizing deep neural networks by noise: Its interpretation and optimization. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 5115–5124, Red Hook, NY, USA, 2017. Curran Associates Inc.

[59] M. Onen, N. Emond, B. Wang, D. Zhang, F. M. Ross, J. Li, B. Yildiz, and J. A. Del Alamo. Nanosecond protonic programmable resistors for analog deep learning. *Science*, 377(6605):539–543, 2022.

[60] N. Papandreou, H. Pozidis, A. Pantazi, A. Sebastian, M. Breitwisch, C. Lam, and E. Eleftheriou. Programming algorithms for multilevel phase-change memory. *IEEE International Symposium on Circuits and Systems*, pages 329–332, 2011.

[61] D. S. Park, W. Chan, Y. Zhang, C.-c. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le. SpecAugment: A simple data augmentation method for automatic speech recognition. In *Proc. Interspeech*, pages 2613–2617, 2019.

[62] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.

[63] A. S. Rakin, Z. He, and D. Fan. Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack, 2018.

[64] M. J. Rasch, T. Gokmen, and W. Haensch. Training large-scale artificial neural networks on simulated resistive crossbar arrays. *IEEE Design & Test*, 37(2):19–29, 2019.

[65] M. J. Rasch, T. Gokmen, M. Rigotti, and W. Haensch. RAPA-ConvNets: Modified convolutional networks for accelerated training on architectures with analog arrays. *Frontiers in Neuroscience*, 13:753, 2019.

[66] M. J. Rasch, D. Moreda, T. Gokmen, M. Le Gallo, F. Carta, C. Goldberg, K. E. Maghraoui, A. Sebastian, and V. Narayanan. A flexible and fast pytorch toolkit for simulating training and inference on analog crossbar arrays. *IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, pages 1–4, 2021.

[67] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks, 2016.

[68] A. Reuther, P. Michaleas, M. Jones, V. Gadepally, S. Samsi, and J. Kepner. Ai accelerator survey and trends. In *2021 IEEE High Performance Extreme Computing Conference (HPEC)*, pages 1–9. IEEE, 2021.

[69] S. Roy, S. Sridharan, S. Jain, and A. Raghunathan. Txsim: Modeling training of deep neural networks on resistive crossbar systems. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 29(4):730–738, 2021.

[70] G. Saon, Z. Tüske, K. Audhkhasi, and B. Kingsbury. Sequence noise injected training for end-to-end speech recognition. In *2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6261–6265, 2019.

[71] A. Sebastian, M. Le Gallo, R. Khaddam-Aljameh, and E. Eleftheriou. Memory devices and applications for in-memory computing. *Nature Nanotechnology*, 15:529–544, 2020.

[72] J. Sevilla, L. Heim, A. Ho, T. Besiroglu, M. Hobbhahn, and P. Villalobos. Compute trends across three eras of machine learning. *arXiv preprint arXiv:2202.05924*, 2022.

[73] K. Spoon, H. Tsai, A. Chen, M. J. Rasch, S. Ambrogio, C. Mackin, A. Fasoli, A. M. Friz, P. Narayanan, M. Stanisavljevic, and G. W. Burr. Toward Software-Equivalent Accuracy on Transformer-Based Deep Neural Networks With Analog Memory Devices. *Frontiers in Computational Neuroscience*, 15(July):1–9, 2021.

[74] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15(1):1929–1958, jan 2014.

[75] X. Sun, J. Choi, C.-Y. Chen, N. Wang, S. Venkataramani, V. V. Srinivasan, X. Cui, W. Zhang, and K. Gopalakrishnan. Hybrid 8-bit floating point (hfp8) training and inference for deep neural networks. *Advances in neural information processing systems*, 32, 2019.

[76] X. Sun, N. Wang, C.-Y. Chen, J. Ni, A. Agrawal, X. Cui, S. Venkataramani, K. El Maghraoui, V. V. Srinivasan, and K. Gopalakrishnan. Ultra-low precision 4-bit training of deep neural networks. In *Advances in Neural Information Processing Systems*, volume 33, pages 1796–1807, 2020.

[77] V. Sze, Y. H. Chen, T. J. Yang, and J. S. Emer. Efficient Processing of Deep Neural Networks: A Tutorial and Survey. *Proceedings of the IEEE*, 105(12):2295–2329, 2017.

[78] A. Taylor, M. Marcus, and B. Santorini. The Penn Treebank: An overview, 2003.

[79] H. Tsai, S. Ambrogio, C. Mackin, P. Narayanan, R. M. Shelby, K. Rocki, A. Chen, and G. W. Burr. Inference of long-short term memory networks at software-equivalent accuracy using 2.5m analog phase change memory devices. In *2019 Symposium on VLSI Technology*, pages T82–T83, 2019.

[80] J.-H. Tsai, Y.-C. Chen, and Y.-T. Liao. A power-efficient bidirectional potentiostat-based readout ic for wide-range electrochemical sensing. In *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 1–5. IEEE, 2018.

[81] S. Wager, S. Wang, and P. Liang. Dropout training as adaptive regularization, 2013.

[82] L. Wan, M. Zeiler, S. Zhang, Y. Le Cun, and R. Fergus. Regularization of neural networks using dropconnect. In S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1058–1066. PMLR, Atlanta, Georgia, USA, 17–19 Jun 2013.

[83] W. Wan, R. Kubendran, C. Schaefer, S. B. Eryilmaz, W. Zhang, D. Wu, S. Deiss, P. Raina, H. Qian, B. Gao, et al. A compute-in-memory chip based on resistive random-access memory. *Nature*, 608(7923):504–512, 2022.

[84] S. Wang and P. Kanwar. BFloat16: The secret to high performance on Cloud TPUs, 2019.

[85] L. G. Wright, T. Onodera, M. M. Stein, T. Wang, D. T. Schachter, Z. Hu, and P. L. McMahon. Deep physical neural networks trained with backpropagation. *Nature*, 601(7894):549–555, 2022.

[86] C.-X. Xue, Y.-C. Chiu, T.-W. Liu, T.-Y. Huang, J.-S. Liu, T.-W. Chang, H.-Y. Kao, J.-H. Wang, S.-Y. Wei, C.-Y. Lee, et al. A cmos-integrated compute-in-memory macro based on resistive random-access memory for ai edge devices. *Nature Electronics*, 4(1):81–90, 2021.

[87] X. Yang, C. Wu, M. Li, and Y. Chen. Tolerating noise effects in processing-in-memory systems for neural networks: A hardware–software codesign perspective. *Advanced Intelligent Systems*, page 2200029, 2022.

[88] P. Yao, H. Wu, B. Gao, J. Tang, Q. Zhang, W. Zhang, J. J. Yang, and H. Qian. Fully hardware-implemented memristor convolutional neural network. *Nature*, 577(7792):641–646, 2020.

[89] S. Zagoruyko and N. Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

[90] C. Zhou, P. Kadambi, M. Mattina, and P. N. Whatmough. Noisy machines: Understanding noisy neural networks and enhancing robustness to analog hardware errors using distillation, 2020.

# A  Supplemental Information

## A.1  Static input range learning versus dynamic scaling

In general it is difficult, if not impossible, to simultaneously achieve perfect utilization of the input activation range, weight range, and output activation range for every input and across every tile in the network.

In our hardware model, we assume that input ranges are statically clipped (and scaled) by a learnable input range $\alpha$ (see Sec. A.1). This input range is thus fixed and re-determined during inference, and is therefore not dynamically adjusted based on the actual values of each input vector. If dynamically adjusted, the input range might get optimized (eg. very small inputs might not be buried in the noise floor because of a dynamic scaling of the input range), and this could improve the signal-to-noise ratio of the AIMC MVM. For instance, it was suggested to scale the input range dynamically by the absolute maximum input value for each input vector to maximize the output range (e.g. so-called noise management in [23]). One drawback of dynamic scaling, however, is that determining the absolute max value of each input vector requires more on-chip digital computation than simply fixing a predetermined scaling factor. Another caveat of enlarging the input range too much is that larger inputs might saturate the output range. In our hardware model, we assume a relatively shallow input-to-output dynamic ratio (IO-ratio), that is we assume 10 fully-on inputs together with 10 fully-on weights (ie. at $g_{max}$) would saturate the output bound so that larger output values are clipped (compare to Fig. 3).

To test whether the trained input ranges in our static approach are sufficiently learned, we turn on the described noise management additionally so that too small inputs would have an improved SNR. To avoid a potential bound saturation when testing the effect of the dynamic input scaling, we also turned on the so-called bound management (which uses an dynamic down-scaling of the input only if the output was clipped, see [23]). We found that across DNNs, the improvement in performance (as measured with $\mathcal{A}^{1s}$) is less than 0.1% for all DNNs except for the transformers, where the improvement is more pronounced but still only on the order of 1%. Altogether, we thus conclude that our HWA training method learned the static input range well and could adapt the weight and input code to the shallow IO-ratio assumption across DNNs. A dynamic input scaling mechanism seems not necessary to reach good accuracy for inference when using HWA training approach. However, such an approach might be beneficial to increase the flexibility of models that must be deployed with less or no HWA training.

## A.2  Details on the HWA training simulations

### A.2.1  LSTM on Penn Treebank Dataset

Here the DNN model is a two-layer LSTM with a hidden size of 650 for word-based prediction on the Penn Treebank Dataset (PTB) [78] using cross entropy loss. The encoder layer is implemented digitally, whereas the decoder layer is assumed to be on AIMC and consists of 10K classes – each corresponding to a word in the dictionary. We performed HWA training generally as described in 'AIMC hardware aware DNN training' with the following particularities. The HWA model is trained for 60 epochs using a conventionally FP trained model as a starting point. We scan several hyper-parameters and find the best HWA-trained model to have a base learning rate of 0.01, a learning rate decay of 0.95 (applied after each epoch if the test error is not improving on the validation set), dropout ratio of 0.5, injection of PCM programming error at $5\times$ the nominal scaling, SGD momentum of 0.9, and maximum gradient norm of 10. We use a weight decay (L2 regularization) factor of $10^{-5}$. Finally, we randomly select 1% of devices to be stuck at $\hat{g}_{min} = 0$ throughout training to provide some added robustness (i.e. drop-connect [82]).

| Parameter description | value | AIHWKIT configuration name |
|---|---|---|
| ADC precision | 8 bit | `forward.out_res` |
| DAC precision | 8 bit | `forward.inp_res` |
| System noise referred to output | 0.04 | `forward.out_noise` |
| Output bound | 10.0 | `forward.out_bound` |
| Input bound | learned | `forward.inp_bound` |
| G-max | $25\mu$ S | `noise_model.g_max` |
| (Max) tile size | $512{\times}512$ | `mapping.max_input_size` |
| Wire-conductance to $\hat{g}_{\mathrm{max}}$ ratio | 571428.57 | `forward.ir_drop_g_ratio` |
| Short-term weight noise | 0.0175 | `PCM_READ`, `forward.w_noise` |
| IR-drop scale | 1.0 | `forward.ir_drop` |
| PCM noise and drift | data-calibrated [56] | `PcmLikeNoiseModel` |
| Drift compensation | global | `GlobalDriftCompensation` |
| Layer bias | digital | `mapping.digital_bias` |
| Digital output scale | column-wise | `mapping.out_scale_columnwise` |

Table A.1: Parameter of our Standard PCM AIMC model for inference (see also Fig. 2 and Fig. 3). The parameter names of RPU configuration settings ( `InferenceRPUConfig`) of the open-source package AIHWKit are indicated. Noise management and bound management techniques [23] are turned off and instead the learned input range is taken.

| Hardware-aware (HWA) training method | 1 sec | 1 hour | 1 day | 1 year |
|---|---|---|---|---|
| HWA method from [36] adapted | $6.79_{\,0.01}$ | $7.27_{\,0.02}$ | $7.80_{\,0.02}$ | $9.60_{0.05}$ |
| our HWA training (avg 3 training trials) | $6.79_{\,0.02}$ | $7.08_{\,0.03}$ | $7.50_{\,0.04}$ | $8.36_{0.06}$ |
| FP$_{32}$ mapped to analog (no HWA training) | $10.03_{\,0.18}$ | $11.94_{\,0.35}$ | $14.27_{\,0.39}$ | $22.63_{1.07}$ |
| only output noise | $7.04_{\,0.05}$ | $7.40_{\,0.07}$ | $8.03_{\,0.09}$ | $10.33_{0.25}$ |
| only short-term weight noise | $7.01_{\,0.04}$ | $7.62_{\,0.09}$ | $8.15_{\,0.12}$ | $10.47_{0.38}$ |
| no noise (but other nonidealities) | $7.36_{\,0.05}$ | $8.04_{\,0.06}$ | $8.67_{\,0.18}$ | $11.38_{0.25}$ |

Table A.2: Validation of our HWA training methods and a number of ablations for a ResNet-32 model on CIFAR10 [36]. Average inference test errors are in percent (mean over 25 inference trials, and standard error of the mean). For comparison, we obtained the same weights as in [36], but mapped them to our modified hardware-model with additional noise sources (see Eq. 3 and Eq. 6) Our method performs results in a $> 10\%$ more robust model (for the most challenging $t_{\mathrm{eval}}$) validating our approach. Note that we here do not use the learning rate schedule nor the Gaussian weight clipping suggested by [36], nor any best checkpoint selection beyond hyper-parameter tuning (added noise strength and learning rate). Note that it is crucial to add weight noise, however, as was also done in [36].

### A.2.2 Speech-to-text LSTM with HMM

An automatic speech recognition DNN based on an hidden Markov model (HMM) acoustic model on the Switchboard (SWB) 300 dataset [17] is used. The training set consists of 262 hours of SWB 1 audio with transcripts provided by the Mississippi State University. The test set is the Hub5 2000 evaluation set composed of two parts: 2.1-hour SWB data from 40 speakers and 1.6-hour call-home (CH) data from 40 speakers. The acoustic model is a 4-layer bidirectional LSTM network with input size of 140 and 512 hidden units per direction. On top of the LSTM layers, there is a linear projection layer with 256 hidden units followed by a softmax output layer with 32K units corresponding to context-dependent HMM states. The acoustic model has a total of 30M trainable weights.

We train the acoustic model in a HWA manner for 20 epochs while monitoring the validation loss, using an 8-bit FP trained model [75] as a starting point. We use an SGD optimizer with momentum of 0.9, batch size of 256, gradient clipping of 10, dropout of 0.1, initial learning rate of 0.005 which drops by $\sqrt{10}$ at epochs 9 and 18, and $2\times$ scaling of the injected nominal PCM programming errors. These were found to be the best hyper-parameters for the HWA trained

network. We also use knowledge distillation, which significantly improves the accuracy of the HWA network. In addition, we remap all weights to the full conductance range every 2000 mini-batches (see Eq. 21). The testing is performed by pushing the input features from test-utterances through the HWA acoustic model, simulated as being implemented in AIMC, to get their posteriors. Thereafter, the HMM-based decoding network – assumed to be implemented in a conventional digital processor without any analog nonidealities – is run on the posteriors to compute the word error rate.

### A.2.3    BERT base: GLUE

We evaluate the BERT base model [19] on the GLUE benchmark, which nominally comprises 9 tasks. However, we excluded one task (WNLI) due to the unusual construction of the data set and the small size of the test set. Therefore, we evaluate 8 GLUE tasks: RTE, STS-B, MRPC, CoLA, SST-2, QNLI, QQP, and MNLI. All linear layers of the 12 layer transformer are assumed to be deployed on AIMC. In a previous study, we examined the effect of integer precision on the attention computation [73]. However, here we assume that the activation multiplications in the attention is done in FP accuracy in digital. We do not train the BERT model from scratch but instead only fine-tune the pre-trained BERT model using HWA training techniques. We use a maximum sequence length of 128, because the length of the vast majority of samples are within 128. During fine-tuning, we scan a variety of hyper-parameters, e.g., learning rate from $1 \times 10^{-5}$ to $2 \times 10^{-4}$. We use a batch size of 5 and fixed-value weight clipping. Dropout is not used because it was not found to benefit accuracy. The number of fine-tuning epochs varies from 5 to 20, depending on the size of the data-set for each GLUE task. By using a trained FP model as a starting point, the HWA model can achieve iso-accuracy with the results obtained using standard FP fine-tuning as a reference on most GLUE tasks. Our accuracy results are reported on the validation data set.

### A.2.4    Albert base: GLUE

The Albert base model is structurally identical to the BERT base model, however, it shares the weights across the 12 layers and thus has a much smaller number of total parameters. Similar to the BERT model, we evaluate the Albert model on the same 8 GLUE tasks, by fine tuning the pre-trained Albert model. Similarly hyper-parameters are scanned to identify the best parameters for the highest accuracy, e.g., learning rate in the range from $5 \times 10^{-6}$ to $1.5 \times 10^{-4}$. Batch size is 10 and no dropout is used.

### A.2.5    ResNet-32: CIFAR10

We used the very same ResNet-32 DNN with slightly non-standard channel setting as defined in the paper [36]. Additional floating point training was done starting from the trained model obtained via ONNX export from the authors of [36], reaching slightly improved accuracy (50 batch size; 2 steps learning rate (LR) schedule with final LR of 1% of the starting LR 0.025; first 2 epochs having 20 times reduced LR warmup; 600 epochs training with cutout augmentation; random cropping and mirroring). Using the same LR scheduling (but different base LR), our standard HWA training was performed starting from the same ONNX model, where we assumed that all layers are computing in AIMC. Weights were remapped roughly once per epoch and injected weight noise ramped up over the initial epochs. The best final injected programming weight noise scales was selected (from values of 3, 4, and 5). Training was repeated for 3 trials and average results are reported. LR was optimized among for 4 values in the range from 0.005 to 0.05. Weight and bias decay was set to 0.001 and auto weight scaling as well as constant down-scaling was used (see 'Learning of weight-to-conductance conversion factors').

### A.2.6  WideResNet-16: CIFAR100

We trained a WideResnet-16 [89] on the CIFAR-100 dataset [43] using standard floating point training (using random resize, flipping, and cutout augmentation). Then we used our standard HWA-training as described in 'AIMC hardware aware DNN training' assuming all layers in AIMC. HWA training was done for 80 epochs, with a brief warmup and 2 step-wise learning rate reductions (to 1% of initial LR). Weights were remapped roughly once per epoch  and injected weight noise ramped up over the initial epochs (best of final weight noise scale of 3 and 4). We averaged the results of 3 training runs. We used auto-weight scaling, constant down-scaling as described in 'Learning of weight-to-conductance conversion factors'. Additionally, we used bias and weight decay (0.001), and 1% drop connect [82], and selected the best of 6 learning rates (roughly log-spaced between 0.0005 and 0.01). Finally, we used distilling of the last activation layer with the floating point model.

### A.2.7  ImageNet DNNs

We used 4 standard architectures for the ImageNet data set: Resnet18 and Resnet50 [28], WideResnet-50 [89], and Densenet [30]. The main difference between the image classification DNNs are the various layer depths and widths. As a starting point, we used a fully trained DNN (using floating point) that is provided with the "torchvision" package [62]. We then re-trained the DNNs in HWA manner as described in 'AIMC hardware aware DNN training'. First and last layers were assumed to be computed in digital and all others in AIMC. The injected weight noise scale and the learning rate were treated as a hyper-parameters. We reduced the learning rate by a factor 10 of after the initial 7-10 epochs. Batch size was adjusted to fit into the GPU memory (around 20). Due to computing time constraints, HWA training was limited a few epochs (below[1] 12) and a couple of hyper-parameter settings (e.g. noise strength scales $1, 2, 3$ and learning rates e.g. $10^{-4}$, $5 \times 10^{-5}$, $10^{-5}$ and $5 \times 10^{-6}$; all combinations simulated). The best setting (typically, learning range $10^{-5}$and noise scale 3.0) was selected based on a good compromise between short PCM drift accuracy and long-term stability. Results in Table 3 are presented for the very same hyper-parameter setting for all drift times. Auto-weight scaling and constant down-scaling techniques were used, as well as ramping up of the weight noise strength of the initial epoch and occasional remapping of the affine scales during training. We also experimented with distilling, which in some cases performed well, however, the final results were obtained without distilling since results were found to be more consistent. Finally, we used a small weight decay ($10^{-5}$) and a small weight drop connect probability of 0.1% ([82]), and turned off the IR-drop during training (but not during inference) as its effect was minimal.

### A.2.8  RNN-T speech-to-text

The architecture consists of an acoustic encoder (transcriptor network), a decoder (prediction network), and a joint network. The encoder comprises 6 bi-directional LSTM layers and a FC layer. Input size to the RNN is 340 (including 100-dimensional i-vectors [18]), hidden size is 640. The prediction network consists of an embedding layer (dictionary size 46, embedding vector size 10), a single unidirectional LSTM layer with 768 cells, and a FC layer. Encoder and prediction network outputs are combined multiplicatively in the joint network, with a tanh activation function, and an additional FC layer and a final log-softmax over 46 characters.

HWA training performs fine-tuning from a model with the same architecture, pre-trained for 20 epochs on the audio and character-level transcripts from the Switchboard-1 Telephone Speech Corpus (262 h of audio), further augmented with speed and tempo perturbation [41], SpecAugment [61], and Sequence Noise Injection [70]. The pre-trained model achieves 11.8% average Word Error Rate (WER) on the Switchboard and Call-Home test sets of the NIST Hub5-2000 evaluation.

---

[1]In case of WideResnet-50, we extended the training to 75 epochs for the last LR reduction step, which however did not change the train error significantly beyond epoch 15 due to the small learning rate.

The HWA model is trained for few (1-5) epochs on audio sequences of maximum length capped at 500 frames, using SGD with momentum as optimizer, with learning rate 1e-3 and the methods described in (see 'AIMC hardware aware DNN training'). The injected weight noise scale was set to 1. Here, neither down-scaling nor auto-scaling was used, and IR-drop was set to 0 during training.