
N Multipliers for N Bits: Learning Bit Multipliers for Non-Uniform Quantization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Effective resource management is critical for deploying Deep Neural Networks
2 (DNNs) in resource-constrained environments, highlighting the importance of
3 low-bit quantization to optimize memory and speed. In this paper, we introduce
4 N-Multipliers-for-N-Bits, a novel method for non-linear quantization designed for
5 efficient hardware implementation. Our method uses N parameters, distinct for
6 every layer and corresponding to the N quantization bits, whose linear combinations
7 span the set of allowed weights (and activations). Furthermore, we learn these
8 parameters in parallel with the weights ensuring exceptional flexibility in the
9 quantizer model with minimal hardware overhead. We validate our method on
10 CIFAR10 and ImageNet, achieving competitive results with 3- and 4-bit quantized
11 models. We demonstrate strong performance on 4-bit quantized Spiking Neural
12 Networks (SNNs), evaluated on the CIFAR10-DVS and N-Caltech 101 datasets.
13 Further, we address the issue of stuck-at faults in hardware, and demonstrate
14 robustness to up to 30% faulty bits.

15 1 Introduction

16 Deep learning dominates computer vision and broader AI applications, with cloud-based models
17 performing inference by transferring data to servers. While effective, this approach is inefficient in
18 terms of data transfer and power consumption. A more efficient alternative, especially for simple
19 tasks, is edge inference using low-power accelerators with fixed-point arithmetic and in-memory or
20 near-memory computing architectures [1, 2, 3, 4, 5]. These architectures, such as crossbar arrays,
21 perform matrix-vector multiplication by accumulating parallel operations. They can be implemented
22 using analog components or digital ones, but both approaches encounter a trade-off between energy
23 efficiency and performance [6, 7]. Quantization, while improving efficiency, often degrades accuracy
24 and is further impacted by hardware faults such as stuck-at (SA) faults [8] where certain weight
25 bits get stuck at either 0 or 1 and become unprogrammable. Addressing both quantization errors
26 and hardware faults is crucial for optimizing edge inference. Low-bit quantization for weights and
27 activations has been explored extensively through quantization-aware training (QAT) and methods
28 such as uniform and non-uniform quantization [9, 10, 11, 12]. Non-uniform methods, such as learning
29 quantization levels or companding functions [13, 14], offer flexibility by learning key parameters.
30 In this work, we propose a QAT scheme that optimizes bit multipliers for each quantization level,
31 balancing performance and hardware efficiency. Unlike prior approaches, our method provides
32 maximum flexibility during learning while still being hardware-friendly, and avoids the inaccuracies
33 introduced by using gradient estimators.

34 Spiking Neural Networks (SNNs) and neuromorphic hardware present a promising solution to the
35 challenge of energy-efficient edge inference. SNNs mimic the behavior of biological neurons by

36 processing information through discrete spikes, making them inherently event-driven and power-
 37 efficient [15, 16]. Neuromorphic systems, such as TrueNorth and Loihi [17, 18], are designed to
 38 leverage the sparse, asynchronous nature of SNNs, enabling real-time inference with significantly
 39 lower power consumption compared to conventional Artificial Neural Network (ANN) accelerators.
 40 When combined with low-bit quantization, SNNs offer further energy reductions without sacrificing
 41 performing, especially when using temporal datasets. We attain excellent performance on 4-bit SNNs,
 42 which can enable extreme low-power inference when ported on neuromorphic hardware.

43 Energy efficiency in edge devices often comes at the cost of circuit non-idealities such as line
 44 resistance and device variability [19, 20, 21], with SA faults introducing more significant challenges.
 45 Existing solutions [22, 8] attempt to handle SA faults via variable encoding or fault-aware training. We
 46 extend fault-aware training by incorporating faulty weights into QAT, modifying the regularization
 47 loss to prevent invalid weight configurations caused by SA faults. Our approach enables robust
 48 training for low-bit quantized models even with a high rate of hardware faults.

49 Our key contributions are summarized as follows:

- 50 • We introduce a novel, flexible, and hardware-compatible quantization framework that learns
 51 N bit multipliers per layer alongside network weights, enabling adaptable precision with
 52 minimal hardware overhead, while spanning a rich set of quantization levels.
- 53 • We show our method’s effectiveness across multiple networks and datasets, achieving
 54 comparable state-of-the-art results for 3- and 4-bit DNNs on CIFAR10 [23] and ImageNet
 55 [24], and 4-bit SNNs on event-based datasets: CIFAR10-DVS [25] and N-Caltech 101 [26].
- 56 • We propose a fault-tolerant quantization method that enables low-bit models to maintain
 57 performance up to 30% faulty bits, as demonstrated on CIFAR10, enhancing robustness.
- 58 • We propose a custom implementation of bit-level multipliers for analog/digital crossbars,
 59 optimized for our quantization scheme and directly portable to neuromorphic hardware.

60 2 Methodology

61 **Preliminaries:** Quantization aims to replace floating-point weights and activations in DNNs with low-
 62 bit representations to reduce memory usage and speed up computations. A general N-bit quantizer
 63 function will have 2^N levels, say $l_1, l_2, \dots, l_{2^N}, 2^N - 1$ transition thresholds, say $t_1, t_2, \dots, t_{2^N-1}$,
 64 and is defined as follows:

$$Q(x) = \begin{cases} l_1 & \text{if } x < t_1 \\ l_i & \text{if } t_{i-1} \leq x < t_i, \quad i = 2, 3, \dots, 2^N - 1 \\ l_{2^N} & \text{if } x \geq t_{2^N-1} \end{cases} \quad (1)$$

65 **Quantizer Model:** We introduce an N -dimensional learnable vector $r \in \mathbb{R}^N$, which defines the
 66 N bit multipliers, alongside a scalar offset c in our quantizer model. The set of allowed quantized
 67 weights or activations is given by:

$$W_r = \{ \langle r, b \rangle + c \mid b \in \{0, 1\}^N \} \quad (2)$$

68 The quantization function maps each full-precision weight to its nearest quantized counterpart:

$$\hat{x} = Q(x, r) = \arg \min_{w_q \in W_r} |x - w_q| \quad (3)$$

69 This design enables a flexible non-uniform quantizer with multiple step sizes, offering hardware
 70 efficiency while preserving the structure of N-bit quantization. Although learning all 2^N quantization
 71 levels would offer maximum flexibility, it would undermine hardware efficiency and the core benefits
 72 of N-bit quantization. Figure 1a illustrates a sample quantizer function. Drawing parallels between a
 73 general N-bit quantizer and the one introduced above, we can see that the elements of the set W_r
 74 serve as the levels, l_1, l_2, \dots, l_{2^N} , and the transition thresholds are defined as $t_i = (l_i + l_{i+1})/2$

75 **Loss and Learning:** We jointly optimize the bit multipliers, offsets, and weights by introducing
 76 an additional quantization-aware loss alongside the standard cross-entropy loss. This allows the
 77 model parameters to be optimized through backpropagation within the usual training pipeline. During
 78 training, the weights remain in full precision but gradually align with their quantized counterparts

79 due to the influence of the quantization-aware loss. The actual quantization is applied post-training,
 80 where the full-precision weights are mapped to their nearest quantized values.

81 **Quantization-Aware Loss:** We define a regularization loss that minimizes the squared error between
 82 each weight and its nearest quantized value. To balance gradient contributions across layers, we
 83 introduce a layer-specific scaling factor. The total loss is formulated as:

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda \sum_{l=1}^L \alpha_l \sum_{i=1}^{n_l} \min_{w_q \in W_r^l} |w_i - w_q|^2 \quad (4)$$

84 where \mathcal{L}_{CE} is the cross-entropy loss and W_r^l represents the set of quantized weights for layer l ,
 85 defined by parameters r^l and c^l . The term α_l is a layer-wise scaling factor, and λ controls the
 86 regularization strength. Following other works[10], we set α_l as $1/\sqrt{N \cdot Q_P}$, where Q_P is $2^b - 1$ for
 87 activations (unsigned data) and $2^{b-1} - 1$ for weights (signed data), respectively; b denotes the number
 88 of bits. Figure 1b illustrates the regularization loss for a sample weight using an arbitrary vector r to
 89 define the quantized weight set. Equivalently, the loss can be expressed as a function of the weights
 90 and bit multipliers. This formulation jointly optimizes the overall objective and the quantization
 91 parameters, including the bit multipliers and offsets that define the quantization function itself.

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda \sum_{l=1}^L \alpha_l \sum_{i=1}^{n_l} |w_i - Q(w_i, r^l)|^2 \quad (5)$$

Gradient Calculation: The gradient calculation for the weights and quantizer parameters is

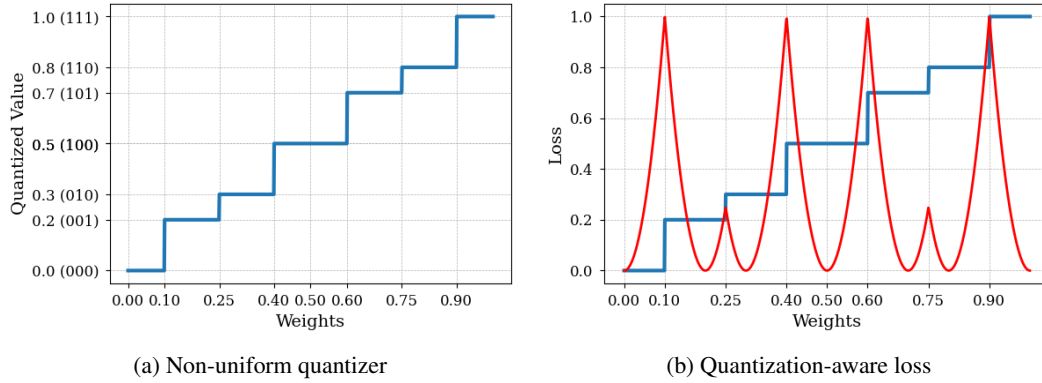


Figure 1: Our quantizer model is non-uniform and learnable. The quantization-aware loss forces weights towards their allowed quantized levels, and the levels towards the weights. The weights are kept in FP during training, and are quantized to their closest allowed level during inference.

92 straightforward. Since we use full precision weights throughout the training, we can simply define
 93 $\frac{\partial Q(w, r)}{\partial w} = 0$, thereby eliminating the need of any gradient approximation techniques. For the
 94 quantizer parameters, $\frac{\partial Q(w, r)}{\partial c} = 1$ and $\nabla_r Q(w, r) = B_r(Q(w, r))$, where B_r is an inverse map
 95 defined as $B_r : W_r \rightarrow \{0, 1\}^N$, providing the bit representation vector of the quantized weights.
 96 This encoding function satisfies $w_q = \langle r, B_r(w_q) \rangle + c \forall w_q \in W_r$. The gradients for the weights, bit
 97 multipliers, and offsets are calculated as follows:
 98

$$\frac{\partial L}{\partial w} = \frac{\partial L_{CE}}{\partial w} + 2\lambda \cdot \alpha_l \cdot (w - Q(w, r^l)) \quad (6)$$

$$\frac{\partial L}{\partial r^l} = 2\lambda \cdot \alpha_l \sum_{i=1}^{n_l} (w_i - Q(w_i, r^l)) \cdot B_r(Q(w_i, r^l)) \quad (7)$$

$$\frac{\partial L}{\partial c^l} = 2\lambda \cdot \alpha_l \sum_{i=1}^{n_l} (w_i - Q(w_i, r^l)) \cdot \left(-\frac{\partial Q(w_i, r^l)}{\partial c^l}\right) = 2\lambda \cdot \alpha_l \sum_{i=1}^{n_l} (Q(w_i, r^l) - w_i) \quad (8)$$

99 **SNN Training:** SNNs inherently produce *quantized activations* in the form of *spike trains*, we thus
 100 need to solely quantize the weights of the network. We use a Leaky Integrate-and-Fire (LIF) model

101 [27] for the spiking neuron in our SNN models. These discrete-time equations describe its dynamics:

$$H[t] = V[t - 1] + \beta(X[t] - (V[t - 1] - V_{reset})) \quad (9)$$

$$S[t] = \Theta(H[t] - V_{th}) \quad (10)$$

$$V[t] = H[t] (1 - S[t]) + V_{reset} S[t] \quad (11)$$

102 where $X[t]$ denotes the input current at time step t . $H[t]$ denotes the membrane potential following
 103 neural dynamics and $V[t]$ denotes the membrane potential after a spike at step t , respectively. The
 104 model uses a firing threshold V_{th} and utilizes the Heaviside step function $\Theta(x)$ to determine spike
 105 generation. The output spike at step t is denoted by $S[t]$, while V_{reset} represents the reset potential
 106 following a spike. The membrane decay constant is denoted by β . To facilitate error backpropagation,
 107 we use the surrogate gradient method [28], defining $\Theta'(x) \triangleq \sigma'(x)$, where $\sigma(x)$ is the arctan
 108 surrogate function [29]. The remaining part of the training and quantization follows that of the
 109 non-spiking networks described earlier.

110 **Fault-Aware Modification:** We propose a two-pronged approach to address SA faults in quantized
 111 neural networks. Firstly, we enhance fault awareness during training by periodically (every 4 epochs)
 112 loading faulty weights onto the model. Secondly, we introduce a fault-aware modification to our
 113 algorithm, designed to avoid weight configurations rendered impossible by SA faults. We introduce a
 114 *validity* term that constrains weights to only those quantization levels that are achievable, avoiding
 115 those rendered unreachable by faulty bits. The *validity* term is defined for each layer as a binary
 116 map that indicates whether a specific weight can attain a given quantization level (1 if achievable, 0
 117 otherwise). This allows us to modify the quantization-aware training loss in Equation 4 as follows:

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda \sum_{l=1}^L \alpha_l \sum_{i=1}^{n_l} \min_{w_q \in W_r^l} (val_{i,q}^l |w_i - w_q|^2 + (1 - val_{i,q}^l) \cdot \Delta) \quad (12)$$

118 Here, $val_{i,q}^l$ represents the *validity* term for weight w_i in layer l with respect to the quantization level
 119 $w_q \in W_r^l$. If w_i can reach w_q , then $val_{i,q}^l = 1$; otherwise, $val_{i,q}^l = 0$. The term Δ is a large constant
 120 that penalizes unreachable quantization levels, effectively excluding them from the optimization.

121 3 Experiments

122 We initialize quantized networks with weights from a trained full-precision model of the same
 123 architecture, then fine-tune in the quantized space, which has been proven to improve performance
 124 [30, 31, 32]. We quantize input activations and weights to 3- or 4-bits for all matrix multiplication
 125 layers except the first and last, which use 8-bits. This approach is commonly used for quantizing deep
 126 networks, and has been proven to increase effectiveness at the cost of minimal overhead [10]. The
 127 weights and the quantization parameters: bit multipliers and the offset values, are trained using SGD
 128 with a momentum of 0.9 and a cosine learning rate decay schedule [33]. We sweep over different
 129 values of the regularization hyperparameter λ and chose $\lambda = 100$ for our results.

130 **ANN Training Details.** We use the ResNet-18 [34] architecture for experiments on CIFAR10 [23]
 131 and ImageNet [24] datasets. Models are trained for 200 epochs on CIFAR10 and 90 epochs on
 132 ImageNet with the weights having a learning rate of 0.01 and 0.1 respectively. The other parameters
 133 are trained with a learning rate of 0.001. For ImageNet, we preprocess images by resizing them to
 134 256×256 pixels. During training, we apply random 224×224 crops and horizontal flips half the
 135 time. At inference, we use a center crop of 224×224 . For CIFAR-10, we augment the training
 136 data by padding images with 4 pixels on each side, then taking random 32×32 crops. We also apply
 137 random horizontal flips half the time. The results are shown in Table 1 and 2.

138 **SNN Training Details.** We use the ResNet-19 [35] and VGG-11 [36] models, after adapting them to
 139 SNNs. Specifically, we replace all ReLU activation functions with LIF modules and substitute max-
 140 pooling layers with average pooling operations. We follow the implementation and data augmentation
 141 technique used in NDA [37] as our baseline training method. The weights and the other parameters
 142 are trained with a learning rate of 0.01 and 0.001 respectively. We evaluate on the N-Caltech 101
 143 and CIFAR10-DVS benchmarks. N-Caltech 101 consists of 8,831 DVS images converted from the
 144 original Caltech 101 dataset, while CIFAR10-DVS comprises 10,000 DVS images derived from the
 145 original CIFAR10 dataset. For both these datasets, we apply a 9:1 train-validation split and resize
 146 all images to 48×48 . Each sample is temporally integrated into 10 frames using spikingjelly [38].
 147 V_{reset} is set to 0 and the membrane decay β is 0.25. Our results are presented in Table 3.

148 **Fault-Aware Training.** We evaluate our method on the VGG-13 architecture, training with 3-bit and
 149 4-bit precision for both weights and activations on the CIFAR10 dataset. Our experiments consider
 150 varying levels of SA fault density. Figures 2a and 2b illustrate the efficacy of our approach for 4-bit
 151 and 3-bit quantization, respectively.

152 4 Results and Analysis

153 **Comparison with Baselines.** Tables 1 and 2 present our quantized ANN results for CIFAR10 and
 154 ImageNet, respectively. Our method outperforms existing approaches, with 4-bit ResNet-18 achieving
 155 a 0.24% accuracy increase over full-precision (FP) on CIFAR10 and matching FP performance
 156 on ImageNet. For 4-bit quantized SNNs (Table 3), we observe performance gains on N-Caltech
 157 101 and marginal losses on CIFAR10-DVS compared to FP. We attribute occasional performance
 158 improvements in both 4-bit ANNs and SNNs to the regularization effect of our quantization loss.

Table 1: Accuracy (%) for 3- and 4- bit quantized ResNet-18 models on CIFAR10. FP denotes full-precision accuracy, Δ FP denotes difference in performance compared to the corresponding FP network. **Best/second best** relative performances for each bit-width are marked in **bold/underlined**.

Method	FP	W4/A4 (Δ FP)	W3/A3 (Δ FP)
L1 Reg [39]	93.54	89.98 (-3.56)	-
BASQ [40]	91.7	90.21 (-1.49)	-
LTS [41]	91.56	91.7 (+0.1)	90.58 (-0.98)
PACT [9]	91.7	91.3 (-0.4)	91.1 (-0.6)
LQ-Nets [13]	92.1	-	91.6 (-0.5)
LCQ [14]	93.4	93.2 (-0.2)	92.8 (-0.6)
Ours	93.26	93.50 (+0.24)	92.84 (-0.42)

Table 2: Accuracy (%) for 4- bit quantized ResNet-18 models on ImageNet. FP denotes full-precision accuracy, Δ FP denotes difference in performance compared to the corresponding FP network. **Best/second best** relative performances for each bit-width are marked in **bold/underlined**.

Method	FP	W4/A4 (Δ FP)
L1 Reg [39]	69.7	57.5 (-12.5)
SinReQ [42]	70.5	64.6 (-5.9)
LTS [41]	69.6	68.3 (-1.3)
PACT [9]	69.7	69.2 (-0.5)
LQ-Nets [13]	70.3	69.3 (-1.0)
QIL [43]	70.2	70.1 (-0.1)
QSin [44]	69.8	69.7 (-0.1)
LCQ [14]	70.4	71.5 (+1.1)
Ours	69.6	69.6 (-0.0)

Table 3: Accuracy (%) for 4- bit quantized SNNs on CIFAR10-DVS and N-Caltech 101. FP denotes full-precision accuracy, Δ FP denotes difference in performance compared to the FP network.

Dataset	Model	FP	W4 (Δ FP)
CIFAR10-DVS	Spiking VGG-11	71.92	71.84 (-0.08)
CIFAR10-DVS	Spiking ResNet-19	72.91	72.14 (-0.77)
N-Caltech 101	Spiking VGG-11	73.19	74.18 (+0.99)
N-Caltech 101	Spiking ResNet-19	75.27	75.93 (+0.66)

159 **Robustness to Faults.** SA faults represent extreme non-idealities in hardware, with each faulty
 160 bit halving the range of possible weight values. High device variability in conductance states can
 161 similarly cause significant discrepancies between expected and realized weights. Our approach,
 162 combining periodic loading of faulty weights during training with a fault-aware modified QAT
 163 algorithm, demonstrates robust performance even under high SA fault densities.

164 **Hardware Compatibility.** Figures 3a and 3b illustrate implementations of custom bit multipliers in
 165 analog and digital crossbar arrays, respectively, compatible with our quantization method. For analog

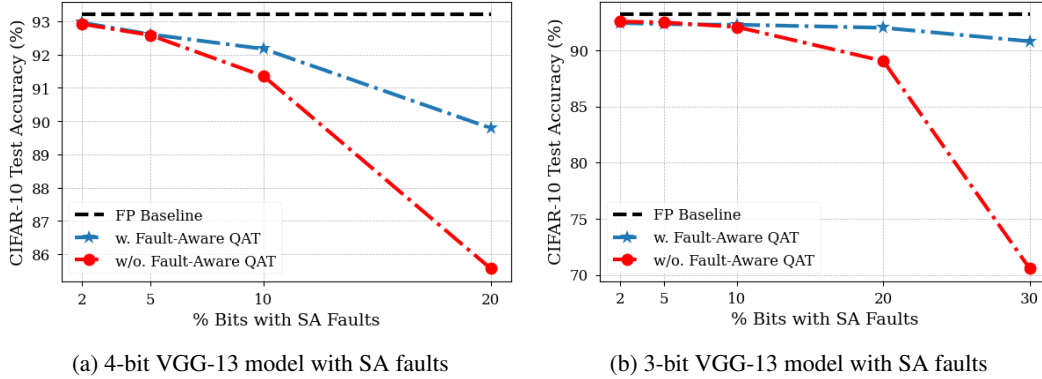


Figure 2: Performance preservation with SA faults: periodic faulty weight loading maintains accuracy for low fault densities; our fault-aware modified QAT extends robustness to high fault fractions.

166 arrays, the implementation incurs no additional cost, requiring only adjustment of bit-multiplier
 167 conductance values from power-of-2 proportions to custom values. In digital arrays, the multiply-
 168 accumulate operation remains unchanged, but peripheral circuitry must be modified to convert right-
 169 shift operations to multiplications, introducing a modest overhead. Learning custom bit multipliers
 170 within QAT can enable highly effective low-bit quantization models, which are compatible with the
 standard in-memory computing architectures.

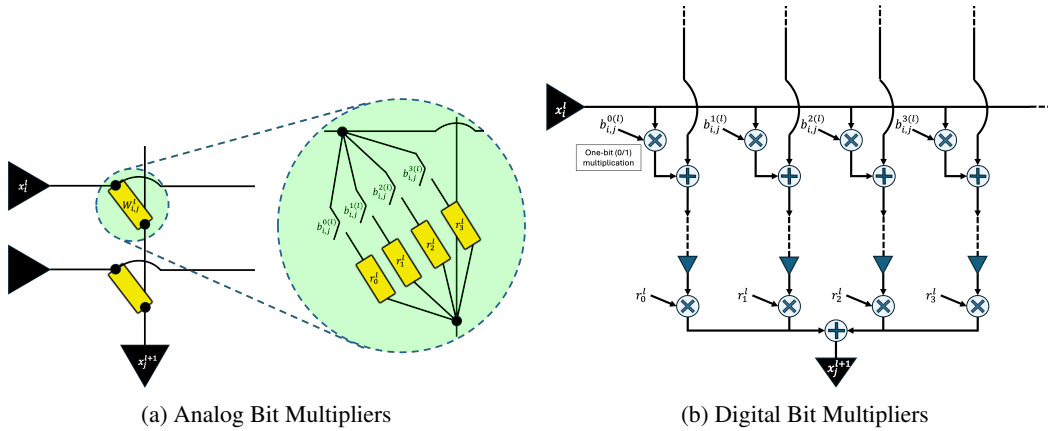


Figure 3: Custom bit multipliers in analog/digital crossbar arrays, compatible with our quantizer.

171

172 5 Conclusion and Future Work

173 We introduce a novel algorithm for learning bit multipliers within QAT, enabling efficient low-bit
 174 quantization models with learnable, non-uniform levels compatible with in-memory computing
 175 architectures. Our approach demonstrates minimal accuracy drops for 3- and 4-bit models compared
 176 to FP baselines across various datasets and architectures, including CIFAR10 and ImageNet using
 177 ResNet-18, and CIFAR10-DVS and N-Caltech 101 using spiking VGG-11 and ResNet-19. Notably,
 178 our quantized models occasionally outperform their FP counterparts. We further extend our method
 179 to address SA faults, maintaining performance with up to 30% faulty bits. Future directions include
 180 extending the method to channel-specific quantizers, conducting fault-aware training experiments
 181 on additional benchmarks, expanding ANN and SNN model evaluations, and exploring sub-3-bit
 182 quantization. These advancements aim to enhance the efficiency and robustness of quantized neural
 183 networks for resource-constrained environments and hardware non-idealities.

References

- 184
- 185 [1] Navjot Kukreja, Alena Shilova, Olivier Beaumont, Jan Huckelheim, Nicola Ferrier, Paul
186 Hovland, and Gerard Gorman. Training on the edge: The why and the how. In *2019 IEEE*
187 *International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages
188 899–903. IEEE, 2019.
- 189 [2] Yu-Hsin Chen, Tien-Ju Yang, Joel Emer, and Vivienne Sze. Eyeriss v2: A flexible accelerator
190 for emerging deep neural networks on mobile devices. *IEEE Journal on Emerging and Selected*
191 *Topics in Circuits and Systems*, 9(2):292–308, 2019.
- 192 [3] Yu-Der Chih, Po-Hao Lee, Hidehiro Fujiwara, Yi-Chun Shih, Chia-Fu Lee, Rawan Naous,
193 Yu-Lin Chen, Chieh-Pu Lo, Cheng-Han Lu, Haruki Mori, et al. 16.4 an 89tops/w and 16.3
194 tops/mm² all-digital sram-based full-precision compute-in memory macro in 22nm for machine-
195 learning edge applications. In *2021 IEEE International Solid-State Circuits Conference (ISSCC)*,
196 volume 64, pages 252–254. IEEE, 2021.
- 197 [4] Hongyang Jia, Hossein Valavi, Yinqi Tang, Jintao Zhang, and Naveen Verma. A programmable
198 heterogeneous microprocessor based on bit-scalable in-memory computing. *IEEE Journal of*
199 *Solid-State Circuits*, 55(9):2609–2621, 2020.
- 200 [5] Jae-sun Seo, Jyotishman Saikia, Jian Meng, Wangxin He, Han-sok Suh, Yuan Liao, Ahmed
201 Hasssan, Injune Yeo, et al. Digital versus analog artificial intelligence accelerators: Advances,
202 trends, and emerging designs. *IEEE Solid-State Circuits Magazine*, 14(3):65–79, 2022.
- 203 [6] Lixia Han, Renjie Pan, Zheng Zhou, Hairuo Lu, Yiyang Chen, Haozhang Yang, Peng Huang,
204 Guangyu Sun, Xiaoyan Liu, and Jinfeng Kang. Comn: Algorithm-hardware co-design platform
205 for non-volatile memory based convolutional neural network accelerators. *IEEE Transactions*
206 *on Computer-Aided Design of Integrated Circuits and Systems*, 2024.
- 207 [7] Hanbo Sun, Zhenhua Zhu, Chenyu Wang, Xuefei Ning, Guohao Dai, Huazhong Yang, and
208 Yu Wang. Gibbon: An efficient co-exploration framework of nn model and processing-in-
209 memory architecture. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and*
210 *Systems*, 2023.
- 211 [8] Muhammad Abdullah Hanif and Muhammad Shafique. Faq: Mitigating the impact of faults
212 in the weight memory of dnn accelerators through fault-aware quantization. *arXiv preprint*
213 *arXiv:2305.12590*, 2023.
- 214 [9] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi
215 Srinivasan, and Kailash Gopalakrishnan. Pact: Parameterized clipping activation for quantized
216 neural networks. *arXiv preprint arXiv:1805.06085*, 2018.
- 217 [10] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dhar-
218 mendra S Modha. Learned step size quantization. *arXiv preprint arXiv:1902.08153*, 2019.
- 219 [11] Chen Tang, Kai Ouyang, Zhi Wang, Yifei Zhu, Wen Ji, Yaowei Wang, and Wenwu Zhu.
220 Mixed-precision neural network quantization via learned layer-wise importance. In *European*
221 *Conference on Computer Vision*, pages 259–275. Springer, 2022.
- 222 [12] Matthias Wess, Sai Manoj Pudukotai Dinakarrao, and Axel Jantsch. Weighted quantization-
223 regularization in dnns for weight memory minimization toward hw implementation. *IEEE*
224 *Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 37(11):2929–2939,
225 2018.
- 226 [13] Dongqing Zhang, Jiaolong Yang, Dongqiangzi Ye, and Gang Hua. Lq-nets: Learned quantiza-
227 tion for highly accurate and compact deep neural networks. In *Proceedings of the European*
228 *conference on computer vision (ECCV)*, pages 365–382, 2018.
- 229 [14] Kohei Yamamoto. Learnable companding quantization for accurate low-bit neural networks. In
230 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
231 5029–5038, 2021.

- 232 [15] Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models.
233 *Neural networks*, 10(9):1659–1671, 1997.
- 234 [16] Kaushik Roy, Akhilesh Jaiswal, and Priyadarshini Panda. Towards spike-based machine
235 intelligence with neuromorphic computing. *Nature*, 575(7784):607–617, 2019.
- 236 [17] Filipp Akopyan, Jun Sawada, Andrew Cassidy, Rodrigo Alvarez-Icaza, John Arthur, Paul
237 Merolla, Nabil Imam, Yutaka Nakamura, Pallab Datta, Gi-Joon Nam, et al. Truenorth: Design
238 and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip. *IEEE transactions*
239 *on computer-aided design of integrated circuits and systems*, 34(10):1537–1557, 2015.
- 240 [18] Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham China, Yongqiang Cao, Sri Harsha
241 Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic
242 manycore processor with on-chip learning. *Ieee Micro*, 38(1):82–99, 2018.
- 243 [19] Xiaochen Peng, Shanshi Huang, Hongwu Jiang, Anni Lu, and Shimeng Yu. Dnn+ neurosim v2.
244 0: An end-to-end benchmarking framework for compute-in-memory accelerators for on-chip
245 training. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*,
246 40(11):2306–2319, 2020.
- 247 [20] Malte J Rasch, Diego Moreda, Tayfun Gokmen, Manuel Le Gallo, Fabio Carta, Cindy Goldberg,
248 Kaoutar El Maghraoui, Abu Sebastian, and Vijay Narayanan. A flexible and fast pytorch toolkit
249 for simulating training and inference on analog crossbar arrays. In *2021 IEEE 3rd international*
250 *conference on artificial intelligence circuits and systems (AICAS)*, pages 1–4. IEEE, 2021.
- 251 [21] Corey Lammie, Wei Xiang, Bernabé Linares-Barranco, and Mostafa Rahimi Azghadi. Mem-
252 torch: An open-source simulation framework for memristive deep learning systems. *Neurocom-*
253 *puting*, 485:124–133, 2022.
- 254 [22] Tao Liu, Wujie Wen, Lei Jiang, Yanzhi Wang, Chengmo Yang, and Gang Quan. A fault-tolerant
255 neural network architecture. In *Proceedings of the 56th Annual Design Automation Conference*
256 *2019*, pages 1–6, 2019.
- 257 [23] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
258 2009.
- 259 [24] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
260 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
261 recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 262 [25] Hongmin Li, Hanchao Liu, Xiangyang Ji, Guoqi Li, and Luping Shi. Cifar10-dvs: an event-
263 stream dataset for object classification. *Frontiers in neuroscience*, 11:309, 2017.
- 264 [26] Garrick Orchard, Ajinkya Jayawant, Gregory K Cohen, and Nitish Thakor. Converting static
265 image datasets to spiking neuromorphic datasets using saccades. *Frontiers in neuroscience*,
266 9:437, 2015.
- 267 [27] Wulfram Gerstner and Werner M Kistler. *Spiking neuron models: Single neurons, populations,*
268 *plasticity*. Cambridge university press, 2002.
- 269 [28] Emre O Neftci, Hesham Mostafa, and Friedemann Zenke. Surrogate gradient learning in spiking
270 neural networks: Bringing the power of gradient-based optimization to spiking neural networks.
271 *IEEE Signal Processing Magazine*, 36(6):51–63, 2019.
- 272 [29] Wei Fang, Zhaofei Yu, Yanqi Chen, Timothée Masquelier, Tiejun Huang, and Yonghong Tian.
273 Incorporating learnable membrane time constant to enhance learning of spiking neural networks.
274 In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2661–2671,
275 2021.
- 276 [30] Jeffrey L McKinstry, Steven K Esser, Rathinakumar Appuswamy, Deepika Bablani, John V
277 Arthur, Izzet B Yildiz, and Dharmendra S Modha. Discovering low-precision networks close
278 to full-precision networks for efficient embedded inference. *arXiv preprint arXiv:1809.04191*,
279 2018.

- 280 [31] Wonyong Sung, Sungho Shin, and Kyuyeon Hwang. Resiliency of deep neural networks under
281 quantization. *arXiv preprint arXiv:1511.06488*, 2015.
- 282 [32] Asit Mishra and Debbie Marr. Apprentice: Using knowledge distillation techniques to improve
283 low-precision network accuracy. *arXiv preprint arXiv:1711.05852*, 2017.
- 284 [33] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv*
285 *preprint arXiv:1608.03983*, 2016.
- 286 [34] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
287 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
288 pages 770–778, 2016.
- 289 [35] Hanle Zheng, Yujie Wu, Lei Deng, Yifan Hu, and Guoqi Li. Going deeper with directly-trained
290 larger spiking neural networks. In *Proceedings of the AAAI conference on artificial intelligence*,
291 volume 35, pages 11062–11070, 2021.
- 292 [36] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
293 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 294 [37] Yuhang Li, Youngeun Kim, Hyoungeob Park, Tamar Geller, and Priyadarshini Panda. Neuro-
295 morphic data augmentation for training spiking neural networks. In *European Conference on*
296 *Computer Vision*, pages 631–649. Springer, 2022.
- 297 [38] Wei Fang, Yanqi Chen, Jianhao Ding, Zhaofei Yu, Timothée Masquelier, Ding Chen, Liwei
298 Huang, Huihui Zhou, Guoqi Li, and Yonghong Tian. Spikingjelly: An open-source machine
299 learning infrastructure platform for spike-based intelligence. *Science Advances*, 9(40):eadi1480,
300 2023.
- 301 [39] Milad Alizadeh, Arash Behboodi, Mart Van Baalen, Christos Louizos, Tijmen Blankevoort,
302 and Max Welling. Gradient l1 regularization for quantization robustness. *arXiv preprint*
303 *arXiv:2002.07520*, 2020.
- 304 [40] Han-Byul Kim, Eunhyeok Park, and Sungjoo Yoo. Basq: Branch-wise activation-clipping
305 search quantization for sub-4-bit neural networks. In *European Conference on Computer Vision*,
306 pages 17–33. Springer, 2022.
- 307 [41] Yunshan Zhong, Gongrui Nan, Yuxin Zhang, Fei Chao, and Rongrong Ji. Exploiting the partly
308 scratch-off lottery ticket for quantization-aware training. *arXiv preprint arXiv:2211.08544*,
309 2022.
- 310 [42] Ahmed T Elthakeb, Prannoy Pilligundla, and Hadi Esmaeilzadeh. Sinreq: Generalized sinu-
311 soidal regularization for low-bitwidth deep quantized training. *arXiv preprint arXiv:1905.01416*,
312 2019.
- 313 [43] Sangil Jung, Changyong Son, Seohyung Lee, Jinwoo Son, Jae-Joon Han, Youngjun Kwak,
314 Sung Ju Hwang, and Changkyu Choi. Learning to quantize deep networks by optimizing
315 quantization intervals with task loss. In *Proceedings of the IEEE/CVF conference on computer*
316 *vision and pattern recognition*, pages 4350–4359, 2019.
- 317 [44] Kirill Solodskikh, Vladimir Chikin, Ruslan Aydarkhanov, Dehua Song, Irina Zhelavskaya, and
318 Jiansheng Wei. Towards accurate network quantization with equivalent smooth regularizer. In
319 *European Conference on Computer Vision*, pages 727–742. Springer, 2022.