# Layerwise Noise Maximisation
# to Train Low-Energy Deep Neural Networks

Sébastien Henwood, François Leduc-Primeau, and Yvon Savaria
Department of Electrical Engineering, École Polytechnique de Montréal
Montréal (QC) Canada
Email: {firstname.lastname}@polymtl.ca

*Abstract*—**Deep neural networks (DNNs) depend on the storage of a large number of parameters, which consumes an important portion of the energy used during inference. This paper considers the case where the energy usage of memory elements can be reduced at the cost of reduced reliability. A training algorithm is proposed to optimize the reliability of the storage separately for each layer of the network, while incurring a negligible complexity overhead compared to a conventional stochastic gradient descent training. For an exponential energy-reliability model, the proposed training approach can decrease the memory energy consumption of a DNN with binary parameters by 3.3× at isoaccuracy, compared to a reliable implementation.**

## I. INTRODUCTION

Deep learning [1] has attracted a lot of interest since 2012 and AlexNet's achievements on ImageNet [2], sparking a rapid improvement of the state-of-the-art in diverse fields. However the trade-off for such results is an antagonizing growth in resource requirements. We now see neural networks such as GPT-2 [3] with as much as 1.5 billion parameters. Such resource requirements make it difficult to implement high performance networks directly in IoT/embedded devices and imply a large energy usage. The energy consumed by memory cells can represent up to 60% of the overall energy consumption [4]. One natural countermeasure is to forsake a bit of accuracy and downscale the network thus reducing the number of parameters and operations. Another approach is to reduce the precision used to store the parameters: numerous quantization schemes have been proposed such as Binary Connect [5] (BC) that represents parameters using only 1 bit. Quantizing the parameters can reduce energy consumption at isoaccuracy, as shown for instance in [6]. Hence, we consider networks binarized with BC as a starting solution on which we want to improve the energy consumption at isoaccuracy.

To reach that goal, the supply voltage of memory cells can be reduced to decrease the energy consumption. However, this can result in a higher likelihood of faults occurring at write or read time and causing data corruption. Several previous works have proposed accepting memory faults to reduce the energy consumption of DNNs. A possible approach is to circumvent faults during inference using advanced fault detection schemes [7]. Sustaining faults through adapted training is also an option: for example, [8] explored the use of a faulty memory on a fully binary neural network and reported a tolerance to a fault rate up to 4% of the memory cells. A training approach for networks with higher precision weights has also been proposed [9], making use of fault detection schemes available in hardware.

Another approach for energy savings in neural networks is to supply less energy to the computation units instead of the memory units, as for example in [10]. While the assumed main type of error is slightly different (timing errors), the analysis is still focused on which weights are the most sensitive: authors propose to decide for the relative importance of weights by using a Taylor expansion of the computations.

As suggested in [9], the memory fault rate can be further optimized per block of a ResNet-style neural network. This echoes with the work from [11], in which a layerwise sensitivity analysis is proposed. They find that while some layers may be critical for the network accuracy, others could be randomized without too much performance degradation. However, optimizing the fault rate for each layer is not a straightforward task since faults occurring in each layer jointly affect the final accuracy.

We consider the case where no fault-detection mechanism is available, and consider the problem of optimizing the memory fault rate separately for each layer of the neural network to create more possibilities for energy reduction. For this, we propose the Layerwise Noise Maximisation (LaNMax) algorithm to optimize the energy-accuracy tradeoff of the network automatically. Compared to a binarized neural network implemented using reliable memories, the LaNMax algorithm finds a fault-rate assignment that reduces memory energy by 3.3× on the CIFAR10 dataset [12], or 2.4× when compared to a uniform-noise approach, at isoaccuracy.

The outline of this paper is as follows: Section II introduces the memory fault model used in our experiments, Section III describes the LaNMax algorithm, Section IV presents the experimental results obtained, and Section V concludes the paper.

## II. MEMORY-FAULT MODEL

In this work, we consider networks binarized with BC: weights $\theta$ are taken in the set $\{-1, 1\}$. It is assumed that supplying less voltage to the memory cells will create faults at rate $p$ when the parameters will be read, with an energy consumption $\eta$ following the formula

$$\eta(p) = -\frac{\ln(p)}{a}, \tag{1}$$

where $a$ is a technology dependant parameter. Such an exponential energy-reliability tradeoff is observed for instance in CMOS SRAM circuits [13]. A one-bit parameter $\tilde{\theta}$ read back from the faulty memory can be described with a binary symmetric channel (BSC) model, defined as

$$BSC(\theta, p) = \begin{cases} -\theta & \text{with probability } p, \\ \theta & \text{with probability } 1-p. \end{cases} \quad (2)$$

We further assume that stored bits are affected independently by faults, which can be approximated in practice even in the presence of correlated bit-cell faults through interleaving.

Since the memory can be modeled as a communication channel, we refer to $p$ as the noise level. We define $\boldsymbol{p}$ a vector of noise levels for each layer $\ell$ of the network. We are interested in optimizing $\boldsymbol{p}$ layerwise to provide an additional degree of freedom when searching for an optimal energy-accuracy trade-off. We can thus express the overall memory energy consumption of the network as

$$E(\boldsymbol{p}) = \zeta * \sum_\ell \eta(p_\ell) n_\ell, \quad (3)$$

where $n_\ell$ is the number of parameters in layer $\ell$ of the network, $\zeta$ the number of bits per parameter and $\eta(p_\ell)$ the energy per bit given in (1).

## III. TRAINING LOW-ENERGY NETWORKS

### A. Uniform noise approach

A first approach for improving the energy efficiency of DNNs consists in assigning a common fault rate (or memory configuration) to all layers, and to train the network under this fault rate. For the network forward pass, weights $\theta$ are sampled through the BSC to obtain $\tilde{\theta}$. For the backpropagation, real weights $\theta$ are updated with the gradient of the corrupted weights $\tilde{\theta}$. It is possible to manually tune the training noise $p_t$ to obtain a desired accuracy level.

Fig. 1 shows the average accuracy achieved on CIFAR10 by networks trained at a fault rate $p_t$ but tested at a rate $p$. We can observe a noise overfitting phenomenon when a network achieves the best test-set accuracy at $p = p_t$, which we observe when $p_t$ becomes sufficiently large. This effect was also reported in [8]. Another interesting observation about the uniform noise case is that a small $p_t > 0$ can provide a regularization that improves accuracy for all choices of $p$, which happens for $p_t = 10^{-4}$ in this example. It is important to keep in mind these phenomena when developing an automated approach for choosing $p_t$.

Starting from a uniform-noise solution that provides a good accuracy-energy trade-off, we can perform a sensitivity analysis to determine whether the reliability of some layers can be decreased further. To do so, we tested the performance of the network when the parameters of layer $\ell$ are replaced with random values. As shown in Fig. 2, some layers are critical for performance (accuracy close to 0 when randomized), but we see that some are less critical, motivating us to propose the LaNMax algorithm to systematically optimize the reliability of each layer.
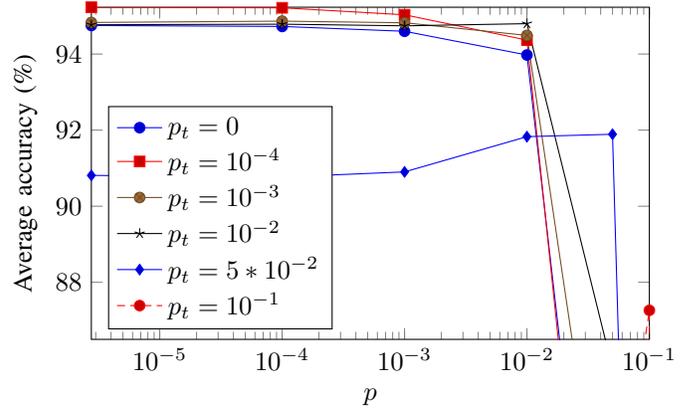


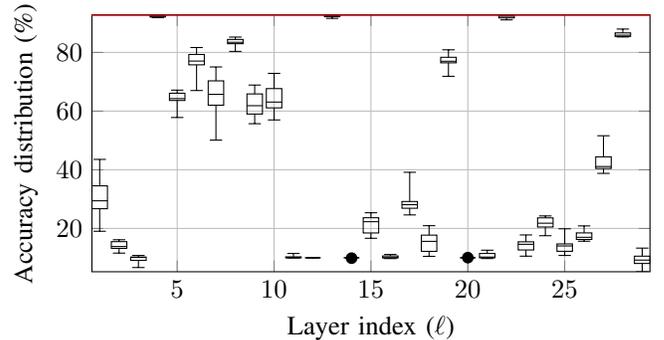Fig. 1. Accuracy VS $p$ for a uniform training fault rate $p_t$.



Fig. 2. Sensitivity analysis under uniform noise $p = 1\%$

### B. LaNMax algorithm: finding the best layerwise $p$

We model the accuracy-energy trade-off with an *outer* loss function $L_O$ made up from the usual loss expectation $A$ and from the energy consumption $E(\boldsymbol{p})$. The usual loss expectation can be written as $A = \mathbb{E}(L_I(f(x|\tilde{\boldsymbol{\theta}}), y))$ with $\tilde{\boldsymbol{\theta}}$ the noisy weights, and $L_I(\cdot, \cdot)$ the *inner* loss of the network output $f(x|\tilde{\boldsymbol{\theta}})$ for an input $x$ and a desired output $y$ (e.g. cross entropy). Note that the loss expectation is taken over the inputs $x$, but also over the noisy weights $\tilde{\boldsymbol{\theta}}$. We define the outer loss as

$$L_O(A, \boldsymbol{p}) = A + \alpha E(\boldsymbol{p}) + \lambda \boldsymbol{p}. \quad (4)$$

Parameter $\alpha$ allows us to set the desired trade-off between energy and accuracy: decreasing $\alpha$ will lead to solutions with better accuracy and larger energy. Because of the possibility of noise overfitting mentioned in Sec. III-A, it could be possible for an optimization scheme to remain stuck in a local minimum. A regularizer $\lambda \boldsymbol{p}$, which we call *noise decay*, is used to correct this problem by adjusting $\lambda$.

We want to find the best reliability for each layer

$$\boldsymbol{p}^\star = \underset{\boldsymbol{p}}{\operatorname{argmin}} \, L_O(A, \boldsymbol{p}), \quad (5)$$

subject to

$$0 < p_\ell \leq \frac{1}{2} \quad \forall \ell. \quad (6)$$

We propose to optimise $L_O$ with a Stochastic Gradient Descent [14] (SGD)-like algorithm (SGD$_O$). Clearly, $L_I$ and $L_O$ cannot be optimized separately since $L_I$ depends indirectly on $\boldsymbol{p}$. We thus need to adapt the network training so that $\theta$ and $\boldsymbol{p}$ are optimized jointly. A possible approach would be to use Block Coordinate Descent [15] with alternating updates based on $L_I$ and $L_O$. Instead, to speed up the training process, we make use of the usual training mini-batch updates of $L_I$ to measure $L_O$ and approximate its gradient.

At each mini-batch, we apply a small random perturbation on $\boldsymbol{p}$: we thus obtain $\tilde{\boldsymbol{p}}$. We can then first estimate $A$ conditioned on $\tilde{\boldsymbol{\theta}}$ for the current mini-batch (of size $m$), which yields $\hat{A} = \frac{1}{m}\sum_i L_I(f(x^{(i)};\tilde{\boldsymbol{\theta}}), y^{(i)})$ and then estimate $L_O$. When all mini-batches have been seen for an epoch, a linear regression is performed on the estimates of $L_O$ of the current epoch, and the slope of this regression is used as an estimate of the outer gradient. An additional advantage of this approach is that the randomness added to $\boldsymbol{p}$ reduces the likelihood of noise overfitting.

To try and achieve a better convergence for SGD$_O$, we rely on a surrogate loss

$$\bar{L}_O(A, \boldsymbol{p}) = \frac{\hat{A}}{\hat{A}^\star} + \alpha\sqrt{E(\boldsymbol{p})} + \lambda\boldsymbol{p}, \tag{7}$$

where $\hat{A}^\star$ is the best loss observed on a mini-batch up to this point in the training of the neural network. This formulation allows to stabilize the value of the first term (accuracy), which would vary due to the dynamics of the neural network training. The gradient of the second term (energy) has large variations depending on the $p_\ell$ values. Taking the square root of the energy reduces these variations and we observe that this improves the SGD convergence.

The procedure to train a neural network using LaNMax is shown in Algorithm 1. Here, LaNMax is applied to a regular SGD training of a neural network's weights $\theta$ (adapted from the SGD update rule in [16]) but it could be integrated to other SGD variants. The network is trained with the usual stopping criterion (line 2), e.g. the number of epochs $K$. While iterating over the $M$ mini-batches (line 3) and if LaNMax is activated (line 5), a random experiment is performed in the neighbourhood of $\boldsymbol{p}$ as follows. First, random perturbation amounts $\delta_\ell$ are sampled uniformly from the set $\{-h, 0, h\}$ (line 6). These values are assembled into a perturbation vector that is added to $\boldsymbol{p}$ (line 7) in order to evaluate the loss function around $\boldsymbol{p}$. After sampling the random weights (line 8), the loss function $L_O$ is computed and stored (line 9). Repeating this experiment across all the mini-batches allows to iteratively refine the numeric gradient value.

If the network is training without LaNMax, the weights are instead sampled from the BSC at the fixed $\boldsymbol{p}$ (line 11). The usual weight update for SGD$_I$ is performed either way at the end of a mini-batch (lines 12, 13). When all the mini-batches have been seen and if LaNMax is still active (line 15), we perform SGD$_O$. First, a numerical gradient is computed via an ordinary least square (OLS) regression (line 16). Then the

---

1   $K \leftarrow 1$
2   **while** *stopping criterion not met* **do**
3     **for** $k \in \{1, 2, \ldots, M\}$ **do**
4       Sample a minibatch $\{x^{(1)}, \ldots, x^{(m)}\}$ from the training set with targets $y^{(i)}$
5       **if** $K \leq s$ **then**
6         $\delta_\ell \sim \mathcal{U}\{-h, 0, h\} \, \forall\ell$
7         $\tilde{\boldsymbol{p}}_k \leftarrow \boldsymbol{p} + \boldsymbol{\delta}$
8         $\tilde{\theta}_i \leftarrow BSC(\theta_i, \tilde{p}_{k,\ell(i)}) \, \forall i$
9         $L_k \leftarrow \bar{L}_O(\hat{A}, \tilde{\boldsymbol{p}}_k)$
10      **else**
11        $\tilde{\theta}_i \leftarrow BSC(\theta_i, p_{\ell(i)}) \, \forall i$
12      $\hat{\boldsymbol{g}} \leftarrow \frac{1}{m}\nabla_\theta \sum_i L_I(f(x^{(i)};\tilde{\boldsymbol{\theta}}), y^{(i)})$
13      $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon\hat{\boldsymbol{g}}$
14     **end**
15     **if** $K \leq s$ **then**
16      $\nabla_p L_O \leftarrow OLS_p(\{L_1, \ldots, L_M\})$
17      $\nabla_p L_O \leftarrow \frac{\nabla_p L_O}{||\nabla_p L_O||}$
18      $\boldsymbol{p} \leftarrow \boldsymbol{p} - \epsilon\nabla_p L_O$
19    $K \leftarrow K + 1$
20 **end**
21 **return** $\boldsymbol{p}, \boldsymbol{\theta}$

**Algorithm 1:** Layerwise Noise Maximisation (LaNMax) update

gradient is normalized to ensure stable updates to $\boldsymbol{p}$ (line 17). Finally, $\boldsymbol{p}$ is updated (line 18).

When training with LaNMax, it is advisable to stop updating $\boldsymbol{p}$ before the end of training through the parameter $s$ (lines 5 and 15) so that the weights can be fine-tuned with the final $\boldsymbol{p}$ in the remaining training epochs. Finally, as in a traditional SGD, a learning rate $\epsilon$ must be specified, but thanks to the gradient normalization in the SGD$_O$, it has a small impact and for all experiments that were performed it was sufficient to use the same learning rate for SGD$_I$ and SGD$_O$.

## IV. EXPERIMENTS

We evaluate the energy reductions obtained by LaNMax by using it to train a WideResNet network [17] that has been shown to work well with binary weights [18] on the CIFAR10 dataset [12]. As a point of comparison, we also evaluate the energy gains obtained by simply reducing the size of the network. For all results, we fix the number of layers to 28, and use the layer-width parameter of the WideResNet architecture to vary the number of parameters. We use a default width parameter of 10, and denote by $\rho$ the number of parameters normalized to this default network. In addition, we also consider the simpler alternative of training the network with a uniform noise level applied to all layers.

For all networks, the energy metric is defined by (3). Note that (1) is not defined at $p = 0$, since the fault probability can never be zero in a physical system. However, for simplicity, we consider that standard implementations have $p = 0$, and an associated energy consumption $\eta(0) \triangleq 1$. Finally, we set the model parameter $a = 12.8$ as suggested in [9].

For training, $SGD_I$ is configured as suggested in [17]. Both $SGD_I$ and $SGD_O$ are configured with 200 epochs, a learning rate decayed every 60 epochs by a factor of 0.2 starting at 0.1, and a batch size of 128. $SGD_I$ uses a Nesterov momentum of 0.9, a weight decay of $5 * 10^{-4}$ and a batch size of 128. $SGD_O$ uses a Nesterov momentum $\beta = 0.2$ and a noise decay $\lambda = 5 \cdot 10^{-4}$. The training set is augmented with padding, crops, and horizontal flips. The network weights are initialized following the state-of-the-art practices [19].

When $p_\ell > 0$, the network weights become stochastic at inference time, and therefore the network output is also stochastic. To properly measure the average accuracy, we evaluate the test set multiple times, each time sampling new random weights, until the confidence interval on the average reaches a magnitude of five percents at 95% confidence level.

We perform several training runs using LaNMax while varying the accuracy-energy tradeoff parameter $\alpha \in \{0.1, 0.05, 0.03, 0.02, 0.01, 0.001, 0\}$, with LaNMax hyperparameters set to $h = 0.01$ and $s = 160$. Since we know from the uniform noise experiments that the network achieves a higher accuracy when $p = 10^{-4}$ than when $p = 0$, we also restrict $p \in [10^{-4}, 0.5]$, and $\boldsymbol{p}$ is initialized for all layers at $p_\ell = 0.01$.

The accuracy results in terms of the energy metric are shown in Fig. 3. The results obtained by a 16-bit floating-point (FP16) implementation are also shown as a point of comparison. Different curves use different strategies for trading off accuracy for reduced energy. For noiseless curves, energy is reduced by decreasing the number of parameters in the architecture. For the uniform noise BC, energy is reduced by increasing the unique $p$ parameter. Finally, for LaNMax results, the algorithm's $\alpha$ parameter is modified. Note that for most data points of the uniform BC result, the test-time $p$ is the same as the training-time $p_t$, but since for $p = 0$ it is better to train at $p_t = 10^{-4}$, this result is shown instead. We can first see that a network trained with uniform noise for all layers provides a limited improvement in energy efficiency. For example the BC with uniform noise decreases energy by 28% at an accuracy of 95.2%, compared to the noiseless BC. On the other hand, the more fine-grained optimization performed by LaNMax manages to find solutions with much smaller energy consumption at equal accuracy. At an accuracy of 95.3%, the LaNMax solution decreases energy by $2.4\times$ with respect to the uniform-noise BC solution, and by $3.3\times$ with respect to the noiseless BC.

It is interesting to note that for the energy-reliability model considered, changes in reliability alone do not allow to cover a wide range of accuracy values with good energy efficiency. Indeed, Fig. 3 shows that it is essential to also modify the number of parameters to move along the Pareto front.

Fig. 4 shows the reliability assignment generated by LaNMax for $\rho = 1$ and $\alpha = 0.001$. This solution achieves an accuracy of 95.4% for $E = 35\%$. As a point of comparison, the uniform-noise BC achieves 94.8% for $E = 36\%$, using $p = 0.01$. To improve the accuracy, the reliability of the first layers was increased to $p_\ell = 10^{-4}$, while the reliability of the last layers was decreased. In the WideResNet architecture,
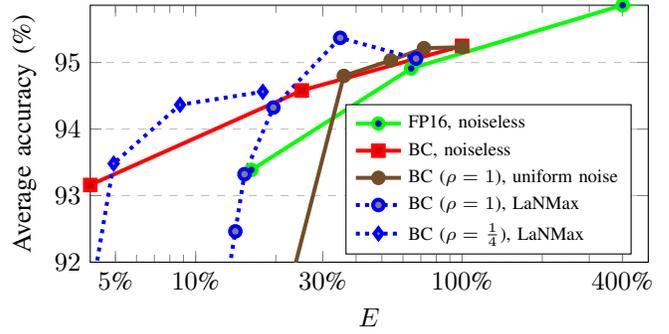


Fig. 3. Average classification accuracy in terms of the normalized energy consumption for a network with 16-bit floating-point parameters and for BC networks with various noise configurations. For the FP16 and noiseless BC curves, energy reductions are obtained by reducing the number of parameters. For the other curves, the normalized number of parameters ($\rho$) is given.
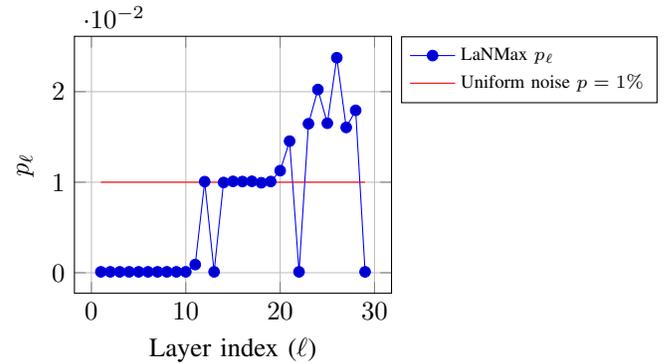


Fig. 4. The layerwise $p_\ell$ for the best result when varying $\alpha$ when $\rho = 1$.

deeper layers have more weights, making it more interesting to reduce the reliability of these layers.

## V. CONCLUSION

In this paper we proposed the LaNMax algorithm to optimise the storage energy of a neural network, with the possibility to control the energy-accuracy trade-off. We considered that the reliability of memory cells decreases exponentially as their energy usage is reduced. LaNMax provides a way to automatically find the best weight-storage reliability for each layer of the neural network while retaining the accuracy. This was achieved by adding an outer SGD-like optimization to the training, while re-using the results generated by the standard training algorithm to reduce overhead.

For a binarized WideResNet on the CIFAR10 dataset, a LaNMax-optimized network achieves the same accuracy as the reliable network with more than three times less energy.

The proposed training approach can help achieve ultra-low energy deep-learning inference by enabling implementations based on near-threshold CMOS circuits as well as emerging technologies.

REFERENCES

[1] Y. Lecun, Y. Bengio, and G. Hinton, "Deep learning," pp. 436–444, May 2015.

[2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[3] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language Models are Unsupervised Multitask Learners," Open AI, Tech. Rep., 2018.

[4] S. Kim, P. Howe, T. Moreau, A. Alaghi, L. Ceze, and V. S. Sathe, "Energy-Efficient Neural Network Acceleration in the Presence of Bit-Level Memory Errors," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 65, no. 12, pp. 4285–4298, Dec. 2018.

[5] M. Courbariaux, Y. Bengio, and J. P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in *Advances in Neural Information Processing Systems*, 2015, pp. 3123–3131.

[6] B. Moons, K. Goetschalckx, N. Van Berckelaer, and M. Verhelst, "Minimum energy quantized neural networks," in *Conference Record of 51st Asilomar Conference on Signals, Systems and Computers, ACSSC 2017*, vol. 2017-October. Institute of Electrical and Electronics Engineers Inc., Apr. 2018, pp. 1921–1925.

[7] B. Reagen, P. Whatmough, R. Adolf, S. Rama, H. Lee, S. K. Lee, J. M. Hernandez-Lobato, G. Y. Wei, and D. Brooks, "Minerva: Enabling Low-Power, Highly-Accurate Deep Neural Network Accelerators," in *Proceedings - 2016 43rd International Symposium on Computer Architecture, ISCA 2016*, 2016, pp. 267–278.

[8] T. Hirtzlin, M. Bocquet, J.-O. Klein, E. Nowak, E. Vianello, J.-M. Portal, and D. Querlioz, "Outstanding Bit Error Tolerance of Resistive RAM-Based Binarized Neural Networks," in *2019 IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, 2019, pp. 288–292.

[9] G. B. Hacene, F. Leduc-Primeau, A. B. Soussia, V. Gripon, and F. Gagnon, "Training modern deep neural networks for memory-fault robustness," in *Proceedings - IEEE International Symposium on Circuits and Systems*, 2019.

[10] D. Shin, W. Choi, J. Park, and S. Ghosh, "Sensitivity based Error Resilient Techniques with Heterogeneous Multiply-Accumulate Unit for Voltage Scalable Deep Neural Network Accelerators," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, pp. 1–1, 2019.

[11] C. Zhang, S. Bengio, and Y. Singer, "Are All Layers Created Equal?" *arXiv preprint arXiv:1902.01996*, Feb. 2019.

[12] A. Krizhevsky, "Learning Multiple Layers of Features from Tiny Images," *Science Department, University of Toronto, Tech.*, pp. 1–60, 2009.

[13] R. G. Dreslinski, M. Wieckowski, D. Blaauw, D. Sylvester, and T. Mudge, "Near-Threshold Computing: Reclaiming Moore's Law Through Energy Efficient Integrated Circuits," *Proceedings of the IEEE*, vol. 98, no. 2, pp. 253–266, Feb. 2010.

[14] L. Bottou, "Online learning and stochastic approximations," *On-line learning in neural networks*, vol. 17, no. 9, p. 142, 1998.

[15] S. J. Wright, "Coordinate Descent Algorithms," Feb. 2015. [Online]. Available: http://arxiv.org/abs/1502.04759

[16] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.

[17] S. Zagoruyko and N. Komodakis, "Wide Residual Networks," in *British Machine Vision Conference 2016, BMVC 2016*, 2016, pp. 1–87.

[18] M. D. McDonnell, "Training wide residual networks for deployment using a single bit for each weight," Feb. 2018. [Online]. Available: http://arxiv.org/abs/1802.08530

[19] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE International Conference on Computer Vision*, Feb. 2015, pp. 1026–1034.