# BiPer: Binary Neural Networks using a Periodic Function

**Edwin Vargas** [1]   **Claudia Correa** [2]   **Carlos Hinojosa** [3]   **Henry Arguello** [2]

## Abstract

Quantized neural networks employ reduced precision representations for both weights and activations. This quantization process significantly reduces the memory requirements and computational complexity of the network. Binary Neural Networks (BNNs) are the extreme quantization case, representing values with just one bit. Since the sign function is typically used to map real values to binary values, smooth approximations are introduced to mimic the gradients during error backpropagation. Thus, the mismatch between the forward and backward models corrupts the direction of the gradient causing training inconsistency problems and performance degradation. In contrast to current BNN approaches, we propose to employ a binary periodic (BiPer) function for the forward pass to obtain the binary values and employ the trigonometric sine function with the same period of the square wave function as a differentiable surrogate during the backward pass. We demonstrate that this approach can control the quantization error by using the frequency of the periodic function and improves network performance. Numerical experiments validate the effectiveness of BiPer in the classification task over ImageNet, with improvements of $0.69\%$.

## 1. Introduction

Deep Neural Networks (DNN) have achieved unprecedented results in many high-level tasks, such as classification, segmentation, and detection, with a tremendous concurrent impact in computer vision, natural language processing, information retrieval, and many others [1]. Typically, DNNs rely on full-precision (32 bit) weights and activation functions.

Accurate and precise models, however, become expensive in terms of computation, storage and number of parameters. For this reason, DNN deployment is usually prohibited for devices with limited resources, such as mobile, hand-held or wearables. Different approaches to reduce computation requirements include efficient neural network architecture design [2]–[4], network pruning [5], knowledge distillation [6], hashing [7], and network quantization [8], [9]. Among them, network quantization has become one of the most promising techniques, aiming at compressing large models usually stored as floating-point weights with low bitwidth numbers. Binary Neural Networks (BNNs) are the extreme quantization case, where weights and activation functions are constrained to just one bit, i.e., binary values, typically +1 or -1. In contrast to DNNs, BNNs replace heavy matrix computations by bit-wise operations, yielding to $32\times$ memory compression, and $58\times$ speed-up on CPUs [10]. Thus, this approach drastically reduces the computational requirements and accelerates inference, making BNNs particularly appealing for resource-constrained environments such as edge devices and mobile applications.

Despite significant advantages for efficient BNN deployment in hardware with limited capabilities, the binarization of full-precision models severely degrades the accuracy performance in high-level tasks such as object detection, classification, segmentation, and others [9]. For instance, in large datasets such as ImageNet, one of the earliest BNN models, the XNOR-Net [10], achieved an accuracy degradation of around $18\%$ compared to the full precision ResNet-18 architecture. Recent efforts have been devoted to close the performance gap of BNN with respect to their real-valued counterparts. Nonetheless, state-of-the-art approaches still exhibit accuracy degradations of approximately $8\%$ [11].

Binarization of real-valued weights and activations is generally performed using the sign function during the feed-forward procedure. A relevant limitation of the sign function is that its gradient is null everywhere except in zero, which makes it incompatible with error back-propagation methods, due to the non-differentiability of binary operations. To overcome this issue, various techniques like the straight-through estimator (STE) and relaxed training approaches have been adopted [12]. STE essentially substitutes the sign function for the identity function to calculate the gradients during the backwards process. Since there exists a mismatch

---

[1]Department of Electrical and Computer Engineering, Rice University, USA [2]Department of Computer Science, Universidad Industrial de Santander, Colombia [3]Department of Computer Science, KAUST, Saudi Arabia. Correspondence to: Edwin Vargas <edwin.vargas@rice.edu>.

between the forward and backward pass caused by the STE approximation, research efforts have focused on designing better smooth and differentiable functions to estimate the gradient of the sign function [13], [14]. Although these approaches have improved the accuracy of BNNs, gradient instability persists when the quantization error is minimized.

Instead of using the Sign function, in this work, we propose to address the aforementioned issues of extreme 1-bit quantization by using a binary periodic (BiPer) function or square wave function to promote binary weight values. Thus, opposite to the sign function which is always negative for negative values or positive for positive values, the proposed periodic function can reach positive and negative values in the whole domain of the latent weights. Since the gradient of the periodic function still faces the problem of being zero almost everywhere, it cannot be directly integrated within a back-propagation algorithm based on gradient descent. We solved this problem by employing a sinusoidal function with the same fundamental frequency of the periodic function as a differentiable surrogate during the backward pass. The continuity and differentiable characteristics of the sine function, make it suitable for stochastic gradient methods. In contrast to existing BNN methods that smoothly and progressively approximate the sign function to reduce the quantization error (QE), we will show that in the proposed BiPer approach the QE can be controlled by the frequency of the periodic function. We further leverage this property to provide an initialization of the weights that better balances the trade-off between the estimation error and performance accuracy. Experimental results demonstrate that BiPer provides the best network performance for the classification task, with respect to state-of-the-art BNN approaches on the CIFAR-10 and ImageNet data sets. The contributions of our work are summarized as follows:

- We propose a simple yet powerful and effective modification in the binarization process, by including a binary periodic function.

- We introduce a continuous, periodic sinusoidal function as a differentiable surrogate of the binary periodic function during the back-propagation process, suitable for stochastic gradient methods.

- We mathematically analyze the quantization error of BiPer and show that it can be controlled by the frequency of the periodic function.

## 2. BiPer

To overcome the gradient and quantization error challenges from existing binarization methods and their gradient approximation functions, we propose to use a binary periodic function or square wave function (see Fig. 1) instead of just

the sign function to model the binary weights. In contrast to the sign function depicted in Fig. 1(a), which is always negative for negative values of $w$, the proposed periodic function (Fig. 1(b)) can reach positive and negative values in the whole domain of the latent weights.
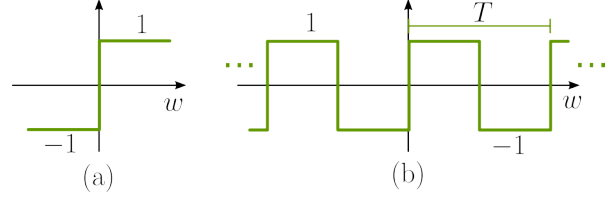


*Figure 1.* (a) Sign function. (b) Binary periodic function.

It should be pointed out that the gradient of the periodic function still faces the problem of being zero almost everywhere, therefore, it cannot be directly integrated within a back-propagation algorithm based on gradient descent. To solve this problem, we first rewrite the square wave function as

$$w^q = \text{Sign}\left(\sin(\omega_0 w)\right), \tag{1}$$

where $\omega_0 = \frac{2\pi}{T}$ is the angular frequency. We note that this corresponds to applying the sign function to the first harmonic of the periodic function. Based on (1), we can approximate the gradient with respect to the weights as

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial w^q}\frac{\partial w^q}{\partial w} \approx \frac{\partial \mathcal{L}}{\partial w^q}\frac{\partial \hat{w}}{\partial w} \tag{2}$$

where $\hat{w} = \sin(\omega_0 w)$. Note that the last differential term in (2) corresponds to the gradient of a continuous differentiable sinusoidal function, which is also a smooth periodic function, and proportional to the frequency $\omega_0$.

### 2.1. Quantization Error Analysis

This section shows that an additional advantage of using the periodic function is its flexibility to control the quantization error. In particular, we mathematically demonstrate how a lower quantization error can be achieved by setting the fundamental period of the wave function. To this end, let us first assume that the latent weights roughly follow the zero-mean Laplace distribution, i.e., $\mathcal{W} \sim La(0, b)$ [15]–[17]. Since the weights $\hat{w}$ before quantization are a function of a random variable, they are also a random variable $\hat{\mathcal{W}} \in [-1, 1]$. Computing the probability density function (pdf) of a random variable $\mathcal{Y} = g(\mathcal{X})$ from the pdf of $\mathcal{X}$ $(f_{\mathcal{X}}(x))$ can be easily done employing the method of transformation [18], if the function $g$ is differentiable and strictly increasing or decreasing, i.e., strictly monotonic. Thus, the pdf of $\mathcal{Y}$ can be computed as

$$f_{\mathcal{Y}}(y) = \begin{cases} \frac{f_{\mathcal{X}}(x_1)}{|g'(x_1)|} = f_X(x_1) \cdot \left|\frac{dx_1}{dy}\right| & \text{where } g(x_1) = y \\ 0 & \text{if } g(x) \neq y. \end{cases}$$

The more general case in which $g$ is not monotonic requires splitting the domain into $n$ intervals, so that $g$ is strictly monotonic and differentiable on each partition. Then, the pdf can be obtained as

$$f_{\mathcal{Y}}(y) = \sum_{k=1}^{n} \frac{f_{\mathcal{X}}(x_k)}{|g'(x_k)|} = \sum_{k=1}^{n} f_{\mathcal{X}}(x_k) \cdot \left| \frac{dx_k}{dy} \right|, \quad (3)$$

where $x_1, \cdots, x_n$ are real solutions to $g(x) = y$. For BiPer, since the sin function is not monotonic, we can use (3) to compute the pdf of $\hat{\mathcal{W}}$ using the pdf of $\mathcal{W}$. Letting $f_{\mathcal{W}}(w) = \frac{1}{b} exp(|w|/b)$ denote the pdf of $\mathcal{W}$, and setting $g$ as the sine function, we can divide the sinusoidal function into subsequent intervals of $T/2$ where it is strictly increasing or decreasing, alternately. The summation in (3) converges to the probability density function of the latent weights before binarization $\hat{w}$ for an arbitrary frequency $\omega_0$ given by

$$f_{\hat{\mathcal{W}}}(\hat{w}) = \frac{1}{2b\omega_0} \frac{1}{\sqrt{1 - \hat{w}^2}} \exp\left( \frac{-|\arcsin(\hat{w})|}{b\omega_0} \right)$$
$$+ \frac{1}{2b\omega_0} \frac{1}{\sqrt{1 - \hat{w}^2}} \cosh\left( \frac{\arcsin(\hat{w})}{b\omega_0} \right) \frac{1}{e^{\pi/b\omega_0} - 1}. \quad (4)$$
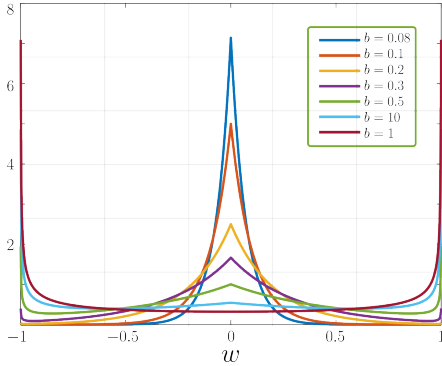


*Figure 2.* Probability density function of $\hat{w} = \sin(\omega_0 w)$ assuming that the random variable $w$ follows a Laplace distribution with parameter $b$ and a fixed value of $\omega_0 = 1$.

Figure 2 depicts the distribution of the weights for different values of the Laplace distribution parameter $b$ and a fixed frequency $\omega_0 = 1$. Note that different from the random variable $\mathcal{W}$ which can take any real value, the codomain of the random variable $\hat{\mathcal{W}}$ is $[-1, 1]$. From Fig. 2 we can observe that when the value of $b$ increases the pdf of $\hat{w}$ behaves as an arcsin distribution with values concentrated around -1 and 1. This reduces the quantization error in comparison to the Laplacian distribution. Also, a similar behavior occurs when the frequency value increases for a fixed $b$. To further analyze these observations consider the QE defined as

$$\text{QE} = \int_{-\infty}^{+\infty} f_{\mathcal{W}}(w) \left( \sin(\omega_0 w) - \gamma \, \text{sign}\left( \sin(\omega_0 w) \right) \right)^2 \, dw, \quad (5)$$

where $f_{\mathcal{W}}$ is the density distribution function of the latent weights. Using the fact that $|x| = x\text{Sign}(x)$ along with the properties of the absolute value, we can rewrite Eq. (5) as

$$\text{QE} = \int_0^{+\infty} \frac{1}{b} \exp\left( \frac{-w}{b} \right) \left( |\sin(\omega_0 w)| - \gamma \right)^2 \, dw. \quad (6)$$

The solution to this integral is given by

$$\text{QE} = \frac{2(\omega_0 b)^2}{4(\omega_0 b)^2 + 1} - \frac{2\gamma\omega_0 b \left( e^{\pi/\omega_0 b} + 1 \right)}{(\omega_0 b)^2 + 1) \left( e^{\pi/\omega_0 b} - 1 \right)} + \gamma^2. \quad (7)$$

On the other hand, the optimal solution of the scaling factor $\gamma$ in (5) can be computed as

$$\gamma = \mathbb{E}\{|\sin(\omega_0 w)|\} = \frac{\omega_0 b \left( e^{\pi/\omega_0 b} + 1 \right)}{(\omega_0 b)^2 + 1) \left( e^{\pi/\omega_0 b} - 1 \right)}. \quad (8)$$

Replacing $\gamma$ from (8) into Eq. (7), we can rewrite the QE as a function of the frequency $\omega_0$ and the parameter $b$. Figure 3 illustrates the QE as a function of the frequency $\omega_0$ for different values of $b$. It can be seen that the maximum QE is 0.102835, which occurs when the product $b\omega_0 \approx 0.954882$ and can be reduced by varying the frequency.
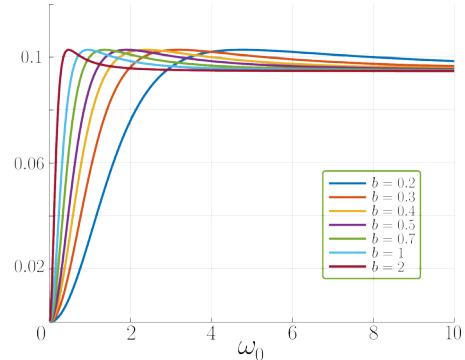


*Figure 3.* Quantization error as a function of the frequency $\omega_0$ for different values of $b$. The proposed BiPer approach is able to control QE with the frequency of the periodic function.

It is worth noting that in contrast to current approaches that progressively reduce the QE to zero, BiPer does not meet this QE value. Nonetheless, further explorations can adapt state-of-the-art surrogate estimators to smoothly converge from the sine function to the square wave

## 3. Experiments

We evaluated BiPer for image classification with a widely used neural network architecture, i.e., ResNet, trained on Imagenet. In the following, we first describe the experiments setup. Then, we compare BiPer with state-of-the-art BNN approaches in terms of performance and complexity.

3

## 3.1. Experiments Setup

**Dataset:** ImageNet [19] is a challenging data set because of its larger size and more diverse image categories. Among its multiple versions, we adopted the widely used ILSVRC12 version, divided into 1,000 categories, from which 1.2 million are training images and, 50,000 test images. ImageNet is the most widely used data set to report results on binary networks and, it allows us to show for the first time that binary networks can perform competitively on a large-scale data set.

**Network Architectures** We chose to binarize ResNet-18/34. We adopted the double skip connections as in [20] to provide fair comparisons. Following [21], the downsampling layers are not quantized, and the double skip connections [20] were included. Following standard procedures of the comparison methods, we binarized all layers but the first and last.

**Training Details and Procedures** All experiments used SGD optimization with $0.9$ momentum. We followed the data augmentation strategies in [22], which include random crops and horizontal flips.

**Two stage training:** Recent works have shown that an appropriate initialization is often required to improve network performance. Two-stage training strategies are generally employed to alleviate feature quantization adverse effects, [23], [24]. Particularly, in the first stage the network is trained with real weights and binary features. Then, in the second stage, a warm weight initialization is employed based on the binary representation of the output weights from the first stage, and the model is fully trained to binarize the weights. In BiPer, we propose a two-stage training where the first stage uses real-valued weights $\hat{w}$, and the second stage uses the weight binarization from Eq. (1). By testing different frequency values, we experimentally found that the hyperparameter frequency $\omega_0 = 20$ balances the QE and precision of the full binary model. The frequency is the same and fixed for both stages. The learning rate was set to $0.1$ in the first stage and, $0.01$ in the second stage. In both stages the learning rate was adjusted by the cosine scheduler.

## 3.2. Comparison with SOTA methods

We evaluate the proposed BiPer approach using ResNet-18 and ResNet-34, and training on the large-scale ImageNet dataset. Table 1 shows a number of SOTA quantization methods over ResNet-18 and ResNet-34, including XNOR-Net [10], Bi-Real Net [20], PCNN [25], IR-Net [26], BONN [27], LCR-BNN [28], HWGQ [29], RBNN [30], FDA [31], ReSTE [32], ReCU [11], and DIR-Net [33]. We can observe that the proposed BiPer approach in the 1W/1A setting achieves the best Top-1 and top-5 accuracy for both network architectures. Specifically, for ResNet-18, we attained a top-1 validation accuracy of $61.4\%$, outperforming the second-

*Table 1.* BiPer performance comparison with state-of-the-art BNN on ImageNet. W/A: bit length of weights and activations. FP: full precision model.

| Network | Method | W/A | Top-1 | Top-5 |
|---------|--------|-----|-------|-------|
| ResNet-18 | FP | 32/32 | 69.6% | 89.2% |
| | XNOR-Net | 1/1 | 51.2% | 73.2% |
| | Bi-Real Net | 1/1 | 56.4% | 79.5% |
| | PCNN | 1/1 | 57.3% | 80.0% |
| | IR-Net | 1/1 | 58.1% | 80.0% |
| | BONN | 1/1 | 59.3% | 81.6% |
| | LCR-BNN | 1/1 | 59.6% | 81.6% |
| | HWGQ | 1/1 | 59.6% | 82.2% |
| | RBNN | 1/1 | 59.9% | 81.9% |
| | FDA | 1/1 | 60.2% | 82.3% |
| | ReSTE | 1/1 | 60.88% | 82.59% |
| | ReCU | 1/1 | 61.0% | 82.6% |
| | DIR-Net | 1/1 | 60.4% | 81.9% |
| | **BiPer (Ours)** | 1/1 | **61.4**% | **83.14**% |
| ResNet-34 | FP | 32/32 | 73.3% | 91.3% |
| | Bi-Real Net | 1/1 | 62.2% | 83.9% |
| | IR-Net | 1/1 | 62.9% | 84.1% |
| | RBNN | 1/1 | 63.1% | 84.4% |
| | ReSTE | 1/1 | 65.05% | 85.78% |
| | ReCU | 1/1 | 65.1% | 85.8% |
| | DIR-Net | 1/1 | 64.1% | 85.3% |
| | **BiPer (Ours)** | 1/1 | **65.73**% | **86.39**% |

best result of $61.0\%$ achieved by ReCu. Furthermore, our top-5 performance reached $83.14\%$, surpassing the second-best result of $82.6\%$, also achieved by ReCU. Likewise, for ResNet-34, we achieved the highest top-1 and top-5 accuracies, namely $65.73\%$ and $86.39\%$, respectively. These results improve the second-best method (ReCU) by $0.63\%$ and $0.59\%$ in top-1 and Top-5 accuracies, respectively.

# 4. Conclusions

An approach for neural network binarization using a binary periodic function or square wave, dubbed BiPer, has been proposed. To improve gradient stability we employed a sinusoidal function with the same period of the square wave as a differentiable surrogate during the backward pass. This simple, yet powerful modification tackles the problem of standard gradient mismatch between forward and backward steps during network training, providing a suitable alternative that can be incorporated within back-propagation algorithms based on stochastic gradient descent. Mathematical analysis of BiPer quantization error demonstrated that it can be controlled by the frequency of the periodic function. Comparisons with respect to state-of-the-art BNN approaches showed that BiPer outperforms prior works by up to $0.63\%$ on Imagenet. Although this work tested the BiPer approach for classification, it can be easily extended to other high-level tasks without increasing the complexity.

# References

[1] C. Wang, A. Bochkovskiy, and H. Liao, "Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Los Alamitos, CA, USA: IEEE Computer Society, Jun. 2023, pp. 7464–7475. DOI: 10.1109/CVPR52729.2023.00721. [Online]. Available: https://doi.ieeecomputersociety.org/10.1109/CVPR52729.2023.00721.

[2] A. G. Howard, M. Zhu, B. Chen, *et al.*, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *CoRR*, vol. abs/1704.04861, 2017. arXiv: 1704.04861. [Online]. Available: http://arxiv.org/abs/1704.04861.

[3] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the 36th International Conference on Machine Learning*, K. Chaudhuri and R. Salakhutdinov, Eds., ser. Proceedings of Machine Learning Research, vol. 97, PMLR, Sep. 2019, pp. 6105–6114. [Online]. Available: https://proceedings.mlr.press/v97/tan19a.html.

[4] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and ¡1mb model size," *ArXiv*, vol. abs/1602.07360, 2016. [Online]. Available: https://api.semanticscholar.org/CorpusID:14136028.

[5] Y. He, J. Lin, Z. Liu, H. Wang, L.-J. Li, and S. Han, "Amc: Automl for model compression and acceleration on mobile devices," in *Proceedings of the European Conference on Computer Vision (ECCV)*, Sep. 2018.

[6] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, vol. 129, pp. 1789–1819, 2021. DOI: https://doi.org/10.1007/s11263-021-01453-z.

[7] E. Eban, Y. Movshovitz-Attias, H. Wu, *et al.*, "Structured multi-hashing for model compression," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 900–11 909. DOI: 10.1109/CVPR42600.2020.01192.

[8] C. Yuan and S. S. Agaian, "A comprehensive review of binary neural network," *Artificial Intelligence Review*, vol. 56, pp. 12 949–13 013, 2023. DOI: https://doi.org/10.1007/s10462-023-10464-w.

[9] R. Sayed, H. Azmi, H. Shawkey, A. H. Khalil, and M. Refky, "A systematic literature review on binary neural networks," *IEEE Access*, vol. 11, pp. 27 546–27 578, 2023. DOI: 10.1109/ACCESS.2023.3258360.

[10] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., Springer International Publishing, 2016, pp. 525–542, ISBN: 978-3-319-46493-0.

[11] Z. Xu, M. Lin, J. Liu, *et al.*, "Recu: Reviving the dead weights in binary neural networks," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 5178–5188. DOI: 10.1109/ICCV48922.2021.00515.

[12] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks," in *NeurIPS*, 2016.

[13] R. Gong, X. Liu, S. Jiang, *et al.*, "Differentiable soft quantization: Bridging full-precision and low-bit neural networks," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Los Alamitos, CA, USA: IEEE Computer Society, Nov. 2019, pp. 4851–4860. DOI: 10.1109/ICCV.2019.00495. [Online]. Available: https://doi.ieeecomputersociety.org/10.1109/ICCV.2019.00495.

[14] H. Le, R. K. Høier, C.-T. Lin, and C. Zach, "Adaste: An adaptive straight-through estimator to train binary neural networks," in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Los Alamitos, CA, USA: IEEE Computer Society, Jun. 2022, pp. 460–469. DOI: 10.1109/CVPR52688.2022.00055. [Online]. Available: https://doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.00055.

[15] C. Baskin, N. Liss, E. Schwartz, *et al.*, "Uniq: Uniform noise injection for non-uniform quantization of neural networks," *ACM Transactions on Computer Systems (TOCS)*, vol. 37, no. 1-4, pp. 1–15, 2021.

[16] Z. Zhang, W. Shao, J. Gu, X. Wang, and P. Luo, "Differentiable dynamic quantization with mixed precision and adaptive resolution," in *International Conference on Machine Learning*, PMLR, 2021, pp. 12 546–12 556.

[17] R. Banner, Y. Nahshan, and D. Soudry, "Post training 4-bit quantization of convolutional networks for rapid-deployment," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[18] H. Pishro-Nik, *Introduction to probability, statistics, and random processes*. Kappa Research, LLC Blue Bell, PA, USA, 2014.

[19] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, ser. NIPS'12, Lake Tahoe, Nevada: Curran Associates Inc., 2012, pp. 1097–1105.

[20] Z. Liu, W. Luo, B. Wu, X. Yang, W. Liu, and K.-T. Cheng, "Bi-real net: Binarizing deep network towards real-network performance," *International Journal of Computer Vision*, vol. 128, pp. 202–219, 2020.

[21] J. Bethge, H. Yang, and C. Meinel, "Training accurate binary neural networks from scratch," in *2019 IEEE International Conference on Image Processing (ICIP)*, IEEE, 2019, pp. 899–903.

[22] 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*, 2016, pp. 770–778.

[23] B. Martinez, J. Yang, A. Bulat, and G. Tzimiropoulos, "Training binary neural networks with real-to-binary convolutions," in *International Conference on Learning Representations*, 2019.

[24] A. Bulat, G. Tzimiropoulos, J. Kossaifi, and M. Pantic, "Improved training of binary networks for human pose estimation and image recognition," *arXiv preprint arXiv:1904.05868*, 2019.

[25] J. Gu, C. Li, B. Zhang, *et al.*, "Projection convolutional neural networks for 1-bit cnns via discrete back propagation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, 2019, pp. 8344–8351.

[26] H. Qin, R. Gong, X. Liu, *et al.*, "Forward and backward information retention for accurate binary neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2020.

[27] J. Gu, J. Zhao, X. Jiang, *et al.*, "Bayesian optimized 1-bit cnns," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 4909–4917.

[28] Y. Shang, D. Xu, B. Duan, Z. Zong, L. Nie, and Y. Yan, "Lipschitz continuity retained binary neural network," in *European conference on computer vision*, Springer, 2022, pp. 603–619.

[29] Z. Cai, X. He, J. Sun, and N. Vasconcelos, "Deep learning with low precision by half-wave gaussian quantization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5918–5926.

[30] M. Lin, R. Ji, Z. Xu, *et al.*, "Rotated binary neural network," *Advances in neural information processing systems*, vol. 33, pp. 7474–7485, 2020.

[31] Y. Xu, K. Han, C. Xu, Y. Tang, C. Xu, and Y. Wang, "Learning frequency domain approximation for binary neural networks," *Advances in Neural Information Processing Systems*, vol. 34, pp. 25 553–25 565, 2021.

[32] X.-M. Wu, D. Zheng, Z. Liu, and W.-S. Zheng, "Estimator meets equilibrium perspective: A rectified straight through estimator for binary neural networks training," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct. 2023, pp. 17 055–17 064.

[33] H. Qin, X. Zhang, R. Gong, Y. Ding, Y. Xu, and X. Liu, "Distribution-sensitive information retention for accurate binary neural network," *International Journal of Computer Vision*, vol. 131, no. 1, pp. 26–47, 2023.