# Dual-Rate Alpha Binary Networks: An Energy-Efficient Optimization Approach for Waste Classification

1st Duc-Anh Pham
*Department of Electrical Engineering and*
*Computer Science, Graduate School of Engineering*
*Tokyo University of Agriculture and Technology*
Tokyo, Japan
s241318t@st.go.tuat.ac.jp

2nd Hironori Nakajo
*Division of Advanced Information Technology and*
*Computer Science, Institute of Engineering*
*Tokyo University of Agriculture and Technology*
Tokyo, Japan
nakajo@cc.tuat.ac.jp

*Abstract*—The current global waste management issue continues to grow, with waste segregation becoming an increasingly urgent challenge. In response, AI-based solutions have shown promising potential for waste classification. However, these solutions often encounter performance bottlenecks due to computational constraints. To address these challenges, we propose a novel approach for waste classification, namely the Dual-Rate Alpha Binary Networks: An Energy-Efficient Optimization Approach for Waste Classification. This method introduces a Binary Neural Network with a new trainable parameter using Dual Variable Learning Rates and employs a flexible weight/activation quantization strategy with Parametric Rectified Linear Unit (PReLU). By combining these techniques with early stopping, our method demonstrates improved performance. The proposed approach is evaluated on the CIFAR10 dataset and shows favorable results compared to existing models. Further experiments on Convolution Neural Network (CNN) and conventional Binary Neural Network (BNN) models, utilizing a publicly available Garbage Classification dataset with over 15,000 images, demonstrate high accuracy and low FLOPS. This research opens up new possibilities for applying binary neural networks to resource-constrained devices.

*Index Terms*—binary neural network, dual variable learning rate, waste classification, energy-efficient computing

## I. INTRODUCTION

### A. Waste Management Problem

With societal development, waste disposal has become a significant challenge for environmental sustainability, particularly in densely populated and developing countries.

Over the years, numerous efforts have been made to address pollution caused by waste, including recycling, disposal, and treatment. However, owing to the manual classification of various types of waste, individuals often encounter challenges in sorting them efficiently, and the significant consumption of time, effort, and financial resources has become an urgent issue that requires immediate attention.

The question that must be addressed is how modern technologies can be applied to automate the waste classification process, thereby minimizing manual labor, enhancing environ-

mental protection, and promoting global economic development.

### B. Problem Statement

Researchers are actively developing methods for automating waste classification using computer vision and deep learning technologies. Convolutional Neural Networks (CNN) and their advanced variants, such as AlexNet [1], VGG [2], Inception [3], and ResNet [4], have shown promising results for sorting tasks. For example, Kang et al. [5] applied a ResNet34-based algorithm for waste classification, incorporating features such as multifeature fusion and residual unit reuse. Sha Meng et al. [8] proposed the X-DenseNet model, which integrates Xception with DenseNet, achieving an accuracy of 94.1 % on the TrashNet dataset.

However, CNNs typically use full-precision (32-bit) weights and activation functions, which often result in a large number of parameters, complex architectures, and lengthy processing times, making them impractical for resource-limited environments such as mobile devices and wearables. For instance, YOLOv3 [7] achieved only 60% accuracy in garbage detection while processing at 20 fps, which is suitable for real-time applications but not ideal for classification tasks on low-power devices. To address these challenges, one promising solution is the use of quantization techniques, which significantly reduce the model size and computational complexity while maintaining acceptable performance in resource-constrained environments.

### C. Proposed Solution

Quantization is a technique used to reduce the computational complexity and memory requirements of deep neural networks by reducing the bit-widths of the weights and activations. This process compresses the model, making it more suitable for resource-constrained environments. Quantization techniques can be applied to different parts of a network, including weights, activations, or both, using various optimization approaches to minimize the accuracy degradation.

Among these techniques, Binary Neural Networks (BNNs) represent an extreme case of quantization, where parameters are restricted to binary values (+1 or -1). Although BNNs reduce computational and memory requirements, they often lead to reduced accuracy owing to excessive quantization. Consequently, many studies have focused on modifying the model, refining binaryization methods for weights and activations, and applying optimization techniques such as regularization, learning rate adjustment, and flexible activation functions to mitigate the negative impact of extreme quantization. Notable approaches include BinaryConnect [10], XNOR-Net [11], and differentiable surrogates such as the straight-through estimator (STE) [12] or Soft Binary Activation [13].

Despite significant progress in the aforementioned research, gradient loss remains a major challenge. Our method addresses this by using learnable parameters ($\alpha_{\text{DVLR}}$) for each parameter, bypassing the traditional approach and effectively optimizing gradient loss, resulting in:

- Proposing a weight and activation quantization model without using STE.
- Applying an efficient activation function, minimizing computational cost (FLOPS) without sacrificing accuracy.
- The effectiveness of the method was demonstrated through experiments on CIFAR-10 and a 12-class public garbage classification dataset.

The remainder of this paper is organized as follows. Related works are reviewed in Section II to provide the research context. The methodology is detailed in Section III. Section IV describes the experimental setup and implementation. Results and analysis are presented in Section V. Finally, conclusions and future work are discussed in Section VI.

## II. RELATED WORK

### A. Binary Neural Networks

A BNN is a specialized type of neural network in which some or all the weights and activations are constrained to 1-bit values, except for the input and output layers as shown in Fig. 1. The process of quantizing from 32-bit precision to 1-bit, known as binarization, is expected to reduce complexity and computation costs. Rastegari et al. [14] demonstrated that XNOR-Net reduces memory usage by approximately 32 times compared to conventional CNNs while maintaining effective performance. Moreover, the use of BNNs introduces a novel approach to deep learning computations by leveraging bitwise operations such as XNOR and bit-counting, which significantly enhance computational efficiency. These operations allow BNNs to replace conventional matrix multiplications with lightweight binary operations, making them highly suitable for deployment in resource-constrained environments such as edge devices and mobile applications.

In Binary Neural Networks (BNNs), the binarization function is typically the sign function, which converts the output into either +1 or -1, as shown in (1).

$$\text{sign}(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases} \tag{1}$$
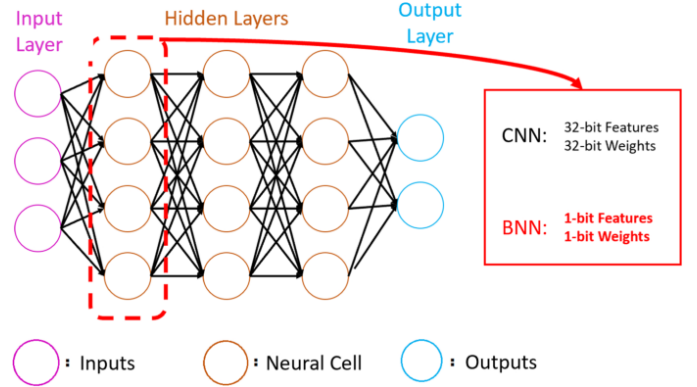


Fig. 1. Binary Neural Networks (Source: Yuan and Agaian [15]).

While this binarization simplifies the computation, it introduces a significant challenge during the backpropagation phase. Specifically, the sign function is non-differentiable, which means that the gradient of the activation with respect to the weights is zero almost everywhere, making it impossible to update the weights using traditional gradient descent. STE is commonly used to approximate the gradient of the sign function by treating it as an identity function during the forward pass and allowing the gradient to pass through as if the activation was continuous during the backward pass. This enables the use of gradient-based optimization techniques despite the non-differentiability of the sign function. An illustration of the process using both the sign function and STE is shown in Fig. 3. This approach has been widely adopted in BNNs to enable efficient training while maintaining binary activation.

### B. ResNet Model

Residual Networks (ResNets) were introduced by He et al. [4], revolutionized deep learning by addressing the vanishing gradient problem using residual connections. These connections allow the gradients to propagate more effectively, thereby enabling the training of very deep networks. In particular, ResNet-18 [16], which consists of 18 layers, provides a simpler and more efficient architecture. ResNet-18 exhibited significant improvements in image classification tasks, outperforming previous architectures on benchmarks such as ImageNet. For BNN, ResNet-18 was adapted for the BNN, allowing for effective training with binary weights and activations.

### C. Dual Variable Learning Rates

Dual Variable Learning Rates (DVLR) [17] is a method designed to optimize the training process of neural networks by adjusting the learning rate for different variables based on their importance. This method helps achieve better convergence by applying distinct learning rates to different parameter groups, ensuring that critical parameters are updated faster, whereas less important parameters are updated more slowly. This approach has been shown to improve the training efficiency and model performance across a variety of tasks,
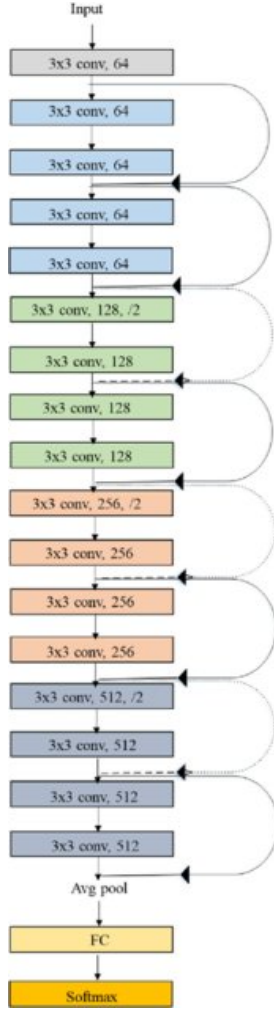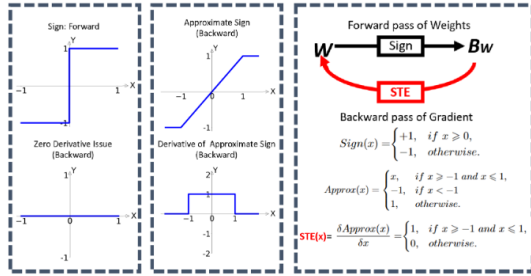
Fig. 2. ResNet18 architecture



Fig. 3. The Sign function and STE used in Binary Neural Networks (BNNs) (Source: Yuan and Agaian [15]).
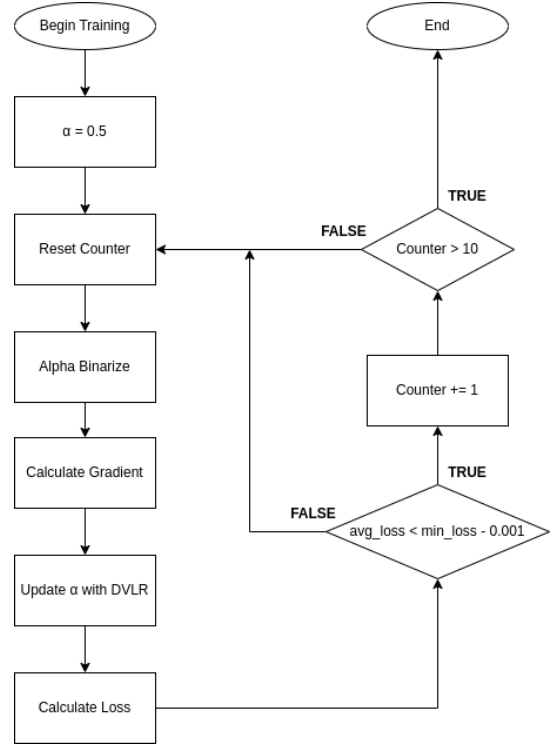


Fig. 4. Illustration of the proposed DRA-BNN framework, highlighting key components and data flow.

$W$ in a network with corresponding gradients $\nabla W$ under DVLR can be expressed as:

$$W^{(t+1)} = W^{(t)} - \eta_i \nabla W^{(t)}, \quad i \in \{1, 2\} \tag{2}$$

where $\eta_1$ and $\eta_2$ are the learning rates applied to different parameter groups. This differential approach to learning rate application ensures that more important parameters are updated at a higher rate, whereas less significant parameters are adjusted more conservatively.

By incorporating DVLR, it becomes possible to stabilize training and improving the performance of the BNNs. This method can be used to apply different learning rates to binary weights, enhance convergence and prevent the model from becoming stuck in suboptimal solutions. Thus, the DVLR is a promising approach for future research on BNNs, potentially leading to more efficient and accurate binary models.

## III. METHODOLOGY

In this section, we address the limitations of traditional BNNs during the training process, such as information loss, reduced accuracy, and limited flexibility, particularly when the input data vary significantly. To address these issues, we propose the Dual-Rate Alpha Binary Networks (DRA-BNN). This approach enhances quantization while maintaining a meaningful gradient flow during backpropagation, effectively preventing gradient vanishing.The flowchart in Fig 4 illustrates the flow of the proposed method.

including image classification and other common machine learning benchmarks.

The effectiveness of DVLR was validated through extensive experiments on standard datasets, demonstrating its ability to accelerate training and achieve higher accuracy. Specifically, it has been shown to outperform traditional methods with a single learning rate, highlighting its potential for optimizing neural network training.

In terms of implementation, the update rule for parameters

## A. Alpha Quantization

In DRA-BNN, the coefficient $\alpha_{DVLR}$ serves as a learnable parameter used to scale the binarized weights after applying the sign($\cdot$) function. The training of $\alpha_{DVLR}$ is performed via gradient descent by using the DVLR method.

In general, the binarized values are computed as follows:

$$\text{binary\_values} = \text{sign(value)} \cdot \alpha_{DVLR} \tag{3}$$

Where:

- sign(weight): The sign function, returning $+1$ if weight $> 0$, $-1$ if weight $< 0$, and 0 if weight $= 0$.
- $\alpha_{DVLR}$: A learnable scaling factor, initialized to 0.5.

By adjusting the $\alpha_{DVLR}$ coefficient based on the data distribution, the binarization coefficients are better characterized and distributed according to the task.

## B. Gradient of $\alpha_{DVLR}$ in the Backward Pass

The gradient of $\alpha_{DVLR}$ is calculated based on the gradient of the loss with respect to the output (grad_output) and the binarized weight values:

$$\text{grad}_{\alpha_{DVLR}} = \sum(\text{grad\_output} \cdot \text{sign(weight)}) \tag{4}$$

Where:

- grad_output: The gradient propagated back from the subsequent layer.
- $\sum$: summation of all elements in the weight tensor.

## C. Updating $\alpha$ in DVLR

The primary goal is to adaptively adjust $\alpha$ with DVLR during training to improve the binarization process while ensuring that the gradient flow remains meaningful.

First, $\alpha$ is constrained within the range $[0.1, 2.0]$ for stability:

$$\alpha_{t+1} = \text{clamp}(\alpha_{t+1}, 0.1, 2.0) \tag{5}$$

To update $\alpha$, we apply a gradient descent rule:

$$\alpha_{t+1} = \alpha_t - \eta_\alpha \cdot \nabla_\alpha L \tag{6}$$

Where $\alpha_t$ represents the $\alpha$ coefficient at time step $t$, $\eta_\alpha$ is the learning rate for $\alpha$, and $\nabla_\alpha L$ is the gradient of the loss function with respect to $\alpha$.

The gradient $\nabla_\alpha L$ ensures that the scaling factor $\alpha$ contributes to preserving the critical information during binarization. This allows the model to adaptively adjust the scaling factor based on the characteristics of the data, there by improving the overall performance without relying on the traditional STE approach.

By updating $\alpha$ in this manner, DVLR enables the model to maintain the computational efficiency while improving the expressiveness and accuracy of the binarized model. The independent and adaptive learning of both weights and $\alpha$ results in a more robust network that can manage imbalanced distributions and vanishing gradients effectively.

## D. PReLU Activation Function

In this study, we explored the use of a more effective activation function to enhance the training performance of the model. Specifically, we consider the use of a Parametric ReLU (PReLU) [20] activation function as an alternative to the traditional ReLU. The PReLU function is defined as (7):

$$f(x) = \begin{cases} x & \text{if } x > 0, \\ \alpha_{\text{PReLU}} x & \text{if } x \le 0, \end{cases} \tag{7}$$

where $\alpha_{\text{PReLU}}$ is a learnable parameter during training.

Unlike ReLU, which outputs zero for negative inputs, PReLU allows the slope of the negative part of the activation function to be learned via the parameter $\alpha_{\text{PReLU}}$. This addresses the "dying ReLU" problem, in which certain neurons may become inactive if their inputs are consistently negative, causing them to stop updating and resulting in poor learning. By learning the parameter $\alpha_{\text{PReLU}}$, PReLU ensures that neurons can still learn from negative inputs, improving model performance compared to ReLU.

## IV. EXPERIMENTS

### A. Dataset

In this experiment, we first validated our model using the CIFAR-10 dataset [18], a widely used benchmark for image classification tasks. Subsequently, we applied our model to a more specific dataset, the Garbage Classification Dataset [19], which consists of 15,150 images categorized into 12 classes of household garbage: paper, cardboard, biological, metal, plastic, green-glass, brown-glass, white-glass, clothes, shoes, batteries, and trash.

This dataset was created to improve the recycling process by classifying waste into more specific categories, allowing for better sorting and higher recycling efficiency. Unlike many other datasets that focus on fewer classes (categories 2-6), the inclusion of 12 categories can significantly enhance recycling efforts by allowing a more precise identification of waste types.

### B. Model Setup and Training

The experiments were conducted on a system with an NVIDIA GeForce GTX 1660 GPU and an Intel Core i5-10400 CPU. The software environment includes Python 3.9.21 and PyTorch 2.5.1.

The model used was BinaryResNet18, which was created by binarizing all layers except the input and output layers to reduce the computational cost. The initial value of the $\alpha_{DVLR}$ parameter was set to 0.5, and the learning rate was initialized at 0.001. The Adam optimizer was used with the CrossEntropyLoss function for multi-class classification, and the learning rate was adjusted using the ReduceLROnPlateau scheduler, which decreases the rate by 50% if the loss did not improve after five epochs, with a minimum LR of 0.0001.

Training was performed for up to 400 epochs, with early stopping applied if there was no significant improvement in the loss after 10 epochs (delta = 0.001). The dataset was split into 80% for training and 20% for validation using data

| Model | Accuracy (%) | FLOPS (M) |
|---|---|---|
| CNN-ResNet18 | 93.3 | 1820 |
| XNOR-Net | 89.83 | 410 |
| BinaryNet | 88.7 | 700 |
| BinaryConnect | 85 | 550 |
| LAB-Net | 87.7 | 700-1000 |
| DRA-BNN (ReLU) | 85.47 | 237 |
| DRA-BNN (PReLU) | **89.80** | **237** |

| Model | FLOPS (M) | Accuracy |
|---|---|---|
| CNN-ResNet18 | 1820 | 91.62% (epoch 167) |
| BNN-STE | 237.63 | 42.89% (epoch 61) |
| DRA-BNN | 237.63 | **82.66%** (epoch 97) |

augmentation techniques such as rotation, flipping, and color adjustment applied to the training set. The validation set was then resized and normalized.

## V. RESULT

The performance of the model was evaluated based on its accuracy and its computational efficiency was measured by calculating FLOPS.

### A. Model Performance Evaluation on CIFAR-10

In Table I, our approach, using PReLU activation, outperforms ReLU in terms of accuracy, with DRA-BNN (PReLU) achieving 89.80% accuracy compared with 85.47% for DRA-BNN (ReLU). This demonstrates that the PReLU allows the model to adapt better to negative inputs, leading to an improved overall performance.

Moreover, while CNN-ResNet18 use higher bit-widths for both weights and activations, our DRA-BNN model still achieves competitive results with significantly reduced computational costs (237M FLOPS). This highlights the feasibility of our approach, providing a strong trade-off between performance and computational efficiency, making it suitable for resource-constrained environments.

### B. Dataset-Specific Result

Based on the promising results on CIFAR-10, our method was evaluated for waste classification. As shown in Table II, the results demonstrate that DRA-BNN achieved an accuracy of 82.66% after 97 epochs, outperforming BNN-STE, which only achieves 42.89% after 61 epochs. Although CNN-ResNet18 achieves a higher accuracy (91.62%), the DRA-BNN method, with much lower FLOPS (237.63M compared to 1820M for CNN-ResNet18), demonstrates the feasibility of the model in computationally constrained environments while maintaining good performance compared to other methods.

To further illustrate the effectiveness of the proposed model, qualitative testing results are shown in Figure 5, where several example predictions from the test set are visualized. Furthermore, the loss and accuracy plots are shown in Fig.(s) 6 and 7,

respectively, and the confusion matrix, displayed in Fig 8 provides a clearer perspective on the classification results.
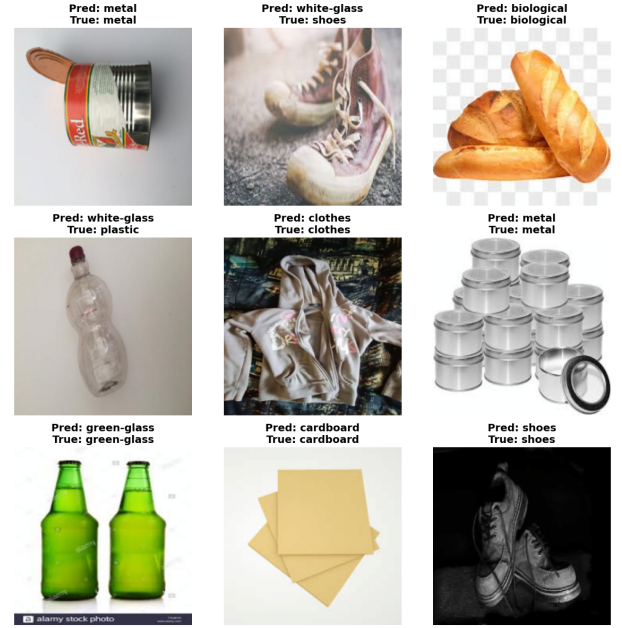


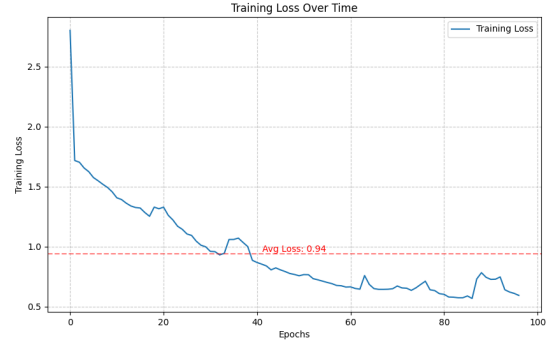Fig. 5. Garbage Classification Testing Result.



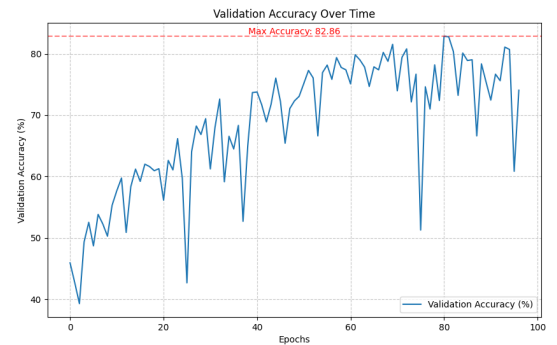Fig. 6. Loss curves on Public Garbage Dataset training phase.



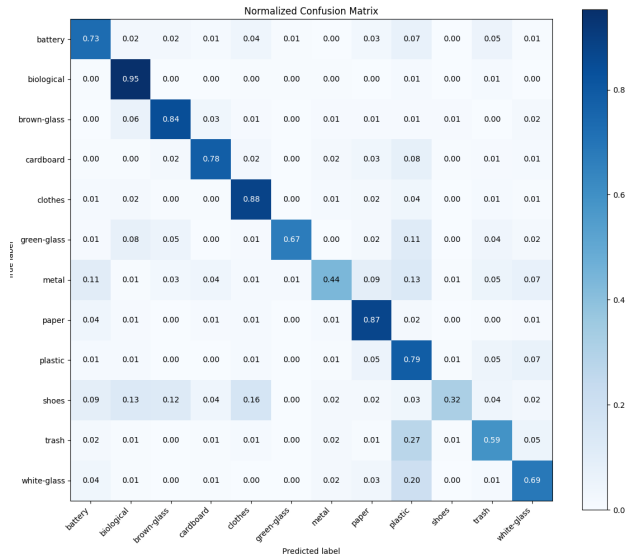Fig. 7. Accuracy on Public Garbage Dataset training phase.

Fig. 8. Confusion matrix result on Public Garbage Dataset.

## VI. Conclusion

In this paper, we propose Dual-Rate Alpha Binary Networks (DRA-BNN) for weight/activation quantization. In the DRA-BNN, the trainable DVLR alpha parameter improves the performance more effectively than in traditional BNN models. By binarizing weights relative to the mean value and computing gradients without relying on STE, this approach provides a compelling trade-off between model size, computational efficiency, and accuracy.

The results presented on general datasets such as CIFAR-10, and specifically on the Garbage Dataset show that the DRA-BNN model, utilizing PReLU, achieves stable performance without complicating the model. Therefore, we conclude that this solution is feasible and opens up new avenue for utilizing high-accuracy BNN models in embedded systems or specialized FPGAs for real-time applications.

In this research, current experiments assume clean, well-annotated input data and do not account for domain shifts or deployment-specific noise. In future work, we plan to investigate the robustness of DRA-BNN under real-world variations and further optimize its computational primitives for efficient FPGA deployment, such as exploiting parallel binary operations and reducing memory access latency in hardware pipelines.

## Acknowledgment

## References

[1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Proc. Neural Information Processing Systems (NIPS)*, 2012, pp. 1097–1105.

[2] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," in *Proc. International Conference on Learning Representations (ICLR)*, 2015.

[3] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, "Rethinking the Inception Architecture for Computer Vision," *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 2818–2826, doi: 10.1109/CVPR.2016.308.

[4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–778, doi: 10.1109/CVPR.2016.90.

[5] Z. Kang, J. Yang, G. Li, Z. Zhang, "An Automatic Garbage Classification System Based on Deep Learning," IEEE Access, Vol. 8, pp. 140019-140029, 2020, doi: 10.1109/ACCESS.2020.3010496.

[6] B. D. Carolis, F. Ladogana and N. Macchiarulo, "YOLO TrashNet: Garbage Detection in Video Streams," 2020 IEEE Conference on Evolving and Adaptive Intelligent Systems (EAIS), Bari, Italy, 2020, pp. 1-7, doi: 10.1109/EAIS48028.2020.9122693.

[7] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," 2018, arXiv:1804.02767, doi: 10.1109/AAAI.2018.01928

[8] Sha Meng, Ning Zhang, and Yunwen Ren, "X-DenseNet: Deep Learning for Garbage Classification Based on Visual Images," *Journal of Physics: Conference Series*, vol. 1575, 5th Annual International Conference on Information System and Artificial Intelligence [ISAI2020], 22-23 May 2020, Zhejiang, China, doi: 10.1088/1742-6596/1575/1/012139.

[9] M. Courbariaux and Y. Bengio, "BinaryNet: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1," CoRR, vol. abs/1602.02830, 2016. Available at: http://arxiv.org/abs/1602.02830.

[10] 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 (NeurIPS)*, 2015, pp. 3123–3131. Available: https://arxiv.org/abs/1511.00363

[11] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2425–2433.

[12] Y. Bengio, P. Simard, and P. Frasconi, "Learning Deep Architectures with a Softmax Activation Function," in *Journal of Machine Learning Research*, vol. 3, pp. 157–191, 2013.

[13] M. Courbariaux, Y. Bengio, and J. David, "BinaryConnect: Training Deep Neural Networks with Binary Weights During Propagation," in *Proc. of Neural Information Processing Systems (NeurIPS)*, 2015.

[14] 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, Cham, 2016, pp. 525–542. doi: 10.1007/978-3-319-46493-0_32.

[15] C. Yuan and S. S. Agaian, "A comprehensive review of Binary Neural Network," *Artificial Intelligence Review*, vol. 56, pp. 12949–13013, 2023. [Online]. Available: https://doi.org/10.1007/s10462-023-10464-w

[16] F. Ramzan, M. U. Khan, A. Rehmat, S. Iqbal, T. Saba, A. Rehman, and Z. Mehmood, "A deep learning approach for automated diagnosis and multi-class classification of Alzheimer's disease stages using resting-state fMRI and residual neural networks," *Journal of Medical Systems*, vol. 44, Dec. 2019, doi: 10.1007/s10916-019-1475-2.

[17] E. Liner and R. Miikkulainen, "Improving neural network learning through dual variable learning rates," *arXiv preprint arXiv:2002.03428*, 2020. doi: 10.48550/arXiv.2002.03428.

[18] A. Krizhevsky, "CIFAR-10 dataset," 2009. Available at: https://www.cs.toronto.edu/ kriz/cifar.html.

[19] M. Abla, "Garbage Classification Dataset," Kaggle, 2023. Available at:https://www.kaggle.com/datasets/mostafaabla/garbage-classification/data.

[20] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification," arXiv preprint arXiv:1502.01852, 2015.