

Pixel to Binary Embedding Towards Robustness for CNNs

Ikki Kishida and Hideki Nakayama

The University of Tokyo

7-3-1 Hongo Bunkyo-ku, Tokyo, Japan

kishida@nlab.ci.i.u-tokyo.ac.jp and nakayama@ci.i.u-tokyo.ac.jp

Abstract—There are several problems with the robustness of Convolutional Neural Networks (CNNs). For example, the prediction of CNNs can be changed by adding a small magnitude of noise to an input, and the performances of CNNs are degraded when the distribution of input is shifted by a transformation never seen during training (e.g., the blur effect). There are approaches to replace pixel values with binary embeddings to tackle the problem of adversarial perturbations, which successfully improve robustness. In this work, we propose Pixel to Binary Embedding (P2BE) to improve the robustness of CNNs. P2BE is a learnable binary embedding method as opposed to previous hand-coded binary embedding methods. P2BE outperforms other binary embedding methods in robustness against adversarial perturbations and visual corruptions that are not shown during training.

I. INTRODUCTION

Convolutional Neural Networks (CNNs) have several issues with robustness. One of the problems is adversarial perturbations: they can maliciously modify CNN's prediction by adding a small magnitude of noise to an input [38]. Since the finding of adversarial perturbations, many types of attacking methods [11], [28] and defensive methods [43], [41] have been proposed. We also know that CNNs do not generalize on input distributions other than the one they are trained on [9]. For example, CNNs trained with regular images fail to generalize on images with the blur effect [40]. CIFAR-C and ImageNet-C [14] are proposed to investigate the generalization ability of trained models on such visually corrupted images. Since then, some training strategies [15], [33] and ensemble techniques [22] have been proposed to improve the robustness against visual corruptions which do not appear during training. Robustness matters for applying the computer vision system to real-world applications since some malicious exploitations may occur using the above flaws.

For the robustness against adversarial perturbations, there are approaches to replace pixel values with binary embeddings (e.g., one-hot and thermometer encoding [3]). They empirically show that binarized input successfully improves the robustness against adversarial perturbations. These binary encodings are based on hand-coded simple rules even though vision tasks are diverse and complex. We consider that such a simple binary encoding would not be optimal for all problems. It is a promising direction to learn the rule of binary encoding for each problem by using data.

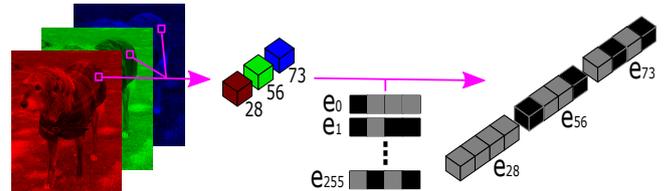


Fig. 1: Overview of P2BE (Pixel to Binary Embedding). Each value of RGB image have 256 types of magnitude (i.e., 0 to 255) and P2BE replaces each RGB value with learned binary embedding $e_k \in \{0, 1\}^M$ where $k \in [0, \dots, 255]$. $M \in \mathbb{N}$ is the dimension of the embedding and hyperparameter controlling the expression ability of e_k . The black and grey colors indicate 0 and 1, respectively. The figure illustrates the case of $M = 4$.

In this work, we propose Pixel to Binary Embedding (P2BE), which is a learnable binary embedding method as opposed to previous hand-coded binary embeddings [3]. In addition, we propose embedding smoothness loss to introduce the effect of quantization which effectively works with adversarial perturbations.

To measure the robustness against visual corruptions, we use two benchmark datasets for image classification (i.e., CIFAR-100-C and ImageNet-C datasets). We show that P2BE outperforms other binary encoding methods for robustness against visual corruptions across various CNNs. In addition, we show that P2BE achieves the best robustness performances against adversarial perturbations. In our analyses, we show that the performance of P2BE is not sensitive to the dimension size M , the proposed embedding smoothness loss is essential to improve robustness against adversarial perturbations, and ImageNet-1k pretrained P2BE has the transferability to the other task.

Our contributions are summarized as follows:

- We propose P2BE, which is a learnable binary embedding unlike other hand-coded binary embedding methods. P2BE shows the best robustness performances among RGB and other binary embedding methods on various datasets.
- The embedding smoothness loss is proposed to realize the effect of quantization in P2BE. Our analysis shows that the embedding smoothness loss improves robustness against adversarial perturbations in P2BE.

- Binary embedding methods have been evaluated only from the view of robustness against adversarial perturbations. In this work, we additionally assess binary embedding methods on robustness against visual corruptions. The results reveal that the approach of binary embedding effectively improves robustness against visual corruptions.

II. RELATED WORK

A. Robustness of CNNs

Adversarial Perturbations. A small amount of adversarial noise can intentionally change the prediction of trained CNNs. This phenomenon is called adversarial perturbations, and it was initially reported in [38]. Since the finding of adversarial perturbations, many types of attacking methods [11], [28], [31], [4], [30], [45], [37] and defensive methods [43], [3], [27], [12], [8], [41], [36], [35], [1], [24], [5], [48] have been proposed. However, it has been reported that the most robust defensive method is still defective against some attacks [21].

Some defensive methods improve robustness against adversarial perturbations by transforming the input $x \in [0, 1, \dots, 255]$ into an M -dimensional binary embedding. One-hot encoding $D_{\text{one-hot}}(x) \in \{0, 1\}^M$ is a simple binary discretization method [3] as follows:

$$D_{\text{one-hot}}(x)_i = \begin{cases} 1, & \text{if } \frac{i-1}{M} \leq \frac{x}{255} < \frac{i}{M} \text{ or } x = 255 \\ 0, & \text{otherwise,} \end{cases} \quad (1)$$

where $i \in \{1, \dots, M\}$, and $M \in \mathbb{N}^+$ is a hyperparameter controlling the size of the dimension of the binary embedding. $D_{\text{one-hot}}$ improves the robustness of CNNs, however, it degrades the performance on clean images by losing the information on relative distance (e.g., $D_{\text{one-hot}}(0.03)$ is equally far from $D_{\text{one-hot}}(0.48)$ and $D_{\text{one-hot}}(0.92)$). To overcome such flaw, thermometer encoding $D_{\text{thermo}}(x) \in \{0, 1\}^M$ [3] is proposed as follows:

$$D_{\text{thermo}}(x)_i = \begin{cases} 1, & \text{if } \frac{x}{255} < \frac{i}{M} \text{ or } x = 255 \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

The examples of binary encoding methods are summarized in Table I.

Common Visual Corruptions. CNNs fail to generalize on the images with visual corruptions, which are not shown during training [9], [40]. It is essential to measure the robustness against such common visual corruptions (e.g., blur, brightness, contrast, and so on) for the reliability of computer vision systems. For benchmarking the robustness of image classification, ImageNet-C and CIFAR-C [14] are proposed, and there are 15 types of common visual corruptions for evaluation. PASCAL-C, COCO-C, and Cityscapes-C are proposed for evaluating the robustness of object detection [29] by using the same types of common visual corruptions.

Robust Training. Some training methods are proposed to improve the robustness of a vision system against adversarial perturbations or visual corruptions. [4] proposed Adversarial

TABLE I: Examples of binary embedding methods. M is the dimensional size of embeddings. The examples in the table represent the case of $M = 10$. P2BE is a learnable binary embedding method. Thus the binary translation rules of P2BE depend on tasks and training strategies.

	One-hot [3]	Thermometer [3]	P2BE (ours)
0.03	[1000000000]	[1111111111]	[0111101011]
0.48	[0000100000]	[0000111111]	[1111101001]
0.92	[0000000001]	[0000000001]	[1011011110]

Training (AT) and they train neural networks on only adversarial perturbed images. AT improves the performance on adversarial perturbed images in exchange for dropping the performance on clean images [47], [39].

Several approaches aim to improve robustness against visual corruptions which never seen during training. One of the approaches is `augxmix` training method [15]. `augxmix` introduces the regularization loss L_{aug} to enforce the model to do consistent predictions between the original and visually transformed images. L_{aug} is computed based on Jensen-Shannon divergence (JSD) between original image (i.e., x) and two visually transformed images (i.e., x_{aug1} and x_{aug2}) as follows:

$$L_{\text{aug}}(p(x); p(x_{\text{aug1}}); p(x_{\text{aug2}})) = \frac{1}{3}(\text{KL}[p(x)||V] + \text{KL}[p(x_{\text{aug1}})||V] + \text{KL}[p(x_{\text{aug2}})||V]), \quad (3)$$

where V is $\frac{1}{3}(p(x) + p(x_{\text{aug1}}) + p(x_{\text{aug2}}))$, KL is Kullback–Leibler divergence and p is CNN’s prediction from the softmax layer. Another approach is an adversarial training method to use an adversarial noise generator [33]. They show that being robust against noise improves the robustness against common visual corruptions.

B. Binary Neural Networks

Deep Neural Networks (DNNs) require heavy matrix computations. Therefore, it is hard to deploy DNNs on devices that have limited computational ability. To overcome such issues, Binary Neural Networks (BNNs) have been proposed [17], [32], [6], [49], [46], [25], [10]. In BNNs, the matrix multiplications are replaced with the combinations of bitwise XNOR and bit count operations, which are lightweight calculations. BNNs accelerate inference time, save up storage, and improve energy efficiency. However, BNNs suffer from performance degradation compared to non-binary DNNs. The motivation of BNNs and binary embedding methods are different. On the one hand, BNNs aim to make DNNs resource-efficient. On the other hand, binary embedding methods translate only input into binary values for improving robustness.

III. METHOD

A. Pixel to Binary Embedding

We show the overview of P2BE in Fig 1. Our method transforms an RGB image $x \in \{0, \dots, 255\}^{3 \times H \times W}$ to the learnable binary embedding $b \in \{0, 1\}^{3M \times H \times W}$ where $M \in \mathbb{N}^+$ is the

dimension size of the binary embedding. There are two steps in P2BE: Binarization and Embedding Smoothness Loss.

Binarization. In P2BE, the learnable embeddings $w_k \in \mathbb{R}^M$ are translated into the binary embedding $e_k \in \{0, 1\}^M$ where $k \in [0, \dots, 255]$ corresponds to the magnitude of each RGB value and $M \in \mathbb{N}^+$ is the hyperparameter to controlling the dimension size of e_k . The binarization is based on the sign function as follows:

$$\text{sign}(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ -1, & \text{otherwise.} \end{cases} \quad (4)$$

In P2BE, we calculate the binary embedding e_k as follows:

$$e_k = 0.5 \times \text{sign}(w_k) + 0.5. \quad (5)$$

Since the sign function is non-differentiable, straight-through estimator (STE) has been proposed to make it differentiable [2]. STE approximates $\frac{\partial \text{sign}(x)}{\partial x}$ by the derivative of the identity function. However, using $\frac{\partial \text{identity}(x)}{\partial x}$ as backward function leads unstabilities in learning since the identity and sign functions are greatly different. Since then, differentiable functions closer to the sign function have been proposed for better approximation of $\frac{\partial \text{sign}(x)}{\partial x}$ [17], [25], [10], [7], [23], [18].

In P2BE, we use a function called approximate sign (i.e., $\text{sign}_{\text{approx}}$) [25] and its derivatives are as follows:

$$\frac{\partial \text{sign}_{\text{approx}}(x)}{\partial x} = \begin{cases} 2 + 2x, & \text{if } -1 \leq x < 0 \\ 2 - 2x, & \text{if } 0 \leq x < 1 \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

$\text{sign}_{\text{approx}}$ function approximates sign function with quadratic equation. When we calculate the gradients of Eq 5 to w_k , we use $\frac{\partial \text{sign}_{\text{approx}}(w_k)}{\partial w_k}$ instead of $\frac{\partial \text{sign}(w_k)}{\partial w_k}$. P2BE transforms the each RGB value $x_{c,h,w} \in [0, \dots, 255]$ to binary values as follows:

$$D_{\text{p2be}}(x_{c,h,w}) = e_{x_{c,h,w}}, \quad (7)$$

where c , h and w are the channel, height and width of an input, respectively. The pseudocode of P2BE is shown in Algorithm 1.

Embedding Smoothness Loss. As [11] have hypothesized, the adversarial perturbations are caused by the linearity of trained neural networks with respect to an input. Let us consider the case of the single linear layer with sigmoid function σ .

$$y(\hat{x}) = \sigma(w(x + \epsilon)) = \sigma(w \cdot x + w \cdot \epsilon), \quad (8)$$

where $w \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^m$ and $\epsilon \in \mathbb{R}^m$ are the weight, input and adversarial noise, respectively. ϵ satisfies $\|\epsilon\|_\infty \leq C$ where $C \in \mathbb{R}^+$ is small enough. If the dimension m is large enough, the output of the sigmoid function can be changed significantly by $w \cdot \epsilon$.

As [3] have mentioned, the quantization can be a reasonable approach against adversarial perturbations since quantized $D_{\text{one-hot}}(x)$ may be equivalent to quantized $D_{\text{one-hot}}(\hat{x})$ and the term of $w \cdot \epsilon$ disappears. As can be seen in Eq 1 and

Algorithm 1: P2BE Pseudocode

Input : Image $x \in \{0, 1, \dots, 255\}^{3 \times K \times N}$,
Learnable Embedding $w \in \mathbb{R}^{256 \times M}$,
Loss L

Initialization: $w \sim \mathcal{N}(0, 1)$

Forwarding (x):

```
//  $e \in \{0, 1\}^{256 \times M}$ 
 $e = \frac{\text{sign}(w)+1}{2}$ 
for  $c = 0, \dots, 2$  do
  for  $k = 0, \dots, K - 1$  do
    for  $n = 0, \dots, N - 1$  do
      for  $m = 0, \dots, M - 1$  do
        //  $b_{Mc+m,k,n} \in \{0, 1\}$ 
         $b_{Mc+m,k,n} = e_{x_{c,k,n},m}$ 
      end
    end
  end
end
```

end
// $b \in \{0, 1\}^{MC \times K \times N}$
return b

Backwarding ($\frac{\delta L}{\delta b}$):

```
//  $\frac{\delta L}{\delta b} \in \mathbb{R}^{3M \times K \times N}$ 
for  $c = 0, \dots, 2$  do
  for  $k = 0, \dots, K - 1$  do
    for  $n = 0, \dots, N - 1$  do
      for  $m = 0, \dots, M - 1$  do
        //  $\times$  is scalar multiplication
        //  $\frac{\delta L}{\delta w_{x_{c,k,n},m}} \in \mathbb{R}$ 
         $\frac{\delta L}{\delta w_{x_{c,k,n},m}} +=$ 
         $\frac{1}{2} \left( \frac{\delta L}{\delta b_{Mc+m,k,n}} \times \frac{\partial \text{sign}_{\text{approx}}(w_{x_{c,k,n},m})}{\delta w_{x_{c,k,n},m}} \right)$ 
      end
    end
  end
end
```

end
// $\frac{\delta L}{\delta w} \in \mathbb{R}^{256 \times M}$
return $\frac{\delta L}{\delta w}$

2, the quantization of $D_{\text{one-hot}}$ and D_{thermo} [3] are pre-defined (e.g., location and the step size of quantization). We propose the embedding smoothness loss to introduce the effect of quantization in P2BE. The embedding smoothness loss L_{smooth} is computed by cosine similarity $\cos(a \angle b): \mathbb{R}^d \times \mathbb{R}^d \rightarrow [1, -1]$ as follows:

$$\begin{aligned} L_{\text{smooth}}(w) &= \sum_{k \in [0, \dots, 254]} 1 - \cos(w_k \angle w_{k+1}) \\ &= \sum_{k \in [0, \dots, 254]} 1 - \frac{\langle w_k, w_{k+1} \rangle}{\|w_k\| \|w_{k+1}\|}, \end{aligned} \quad (9)$$

where $\langle \cdot, \cdot \rangle$ denotes the dot product and $\|\cdot\|$ represents the l2 norm. The smaller L_{smooth} , the angle of neighbored embeddings w_k and w_{k+1} is closer to 0.

IV. EXPERIMENTS

A. Preparations

Dataset. We use CIFAR-100 [19] and ImageNet-1k [34] for our experiments. CIFAR-100 are image classification datasets with 100 classes. They contain 50000 training images and 10000 validation images. We use three types of models: Wide ResNet 40-2 [44], DenseNet-BC ($k=12, d=100$) [16] and ResNeXt-29 (32×4) [42].

ImageNet-1k is the large-scale dataset for image classification with 1.28M training images and 50k validation images of 1000 classes. In this work, we use ResNet50 [13] as the baseline model for ImageNet-1k experiments.

Evaluation. We evaluate the robustness of models in two aspects: common visual corruptions and adversarial perturbations. As the benchmarks of the robustness against common visual corruptions, we use CIFAR-100-C and ImageNet-C datasets [14]. Fifteen types of visual corruptions c transform the images with five different severities s (e.g., blurring, contrasting). On CIFAR-100-C, we calculate the average of the test error $E_{s,c}$ across all corruptions c and severities s . On ImageNet-C, we calculate the mean Corruption Error (mCE) for the measurement of the robustness as proposed in [14]. mCE is the average of Corruption Error (CE_c) across all corruptions c . CE_c is normalized test error as follows: $CE_c = \sum_{s=1}^5 E_{s,c} / \sum_{s=1}^5 E_{s,c}^{\text{alexnet}}$ where $E_{s,c}^{\text{alexnet}}$ is the test error of alexnet [20].

To evaluate robustness against adversarial perturbations, we measure the test error on adversarially perturbed test images. Since binary embedding methods are not differentiable, ordinal attacking methods of adversarial perturbations are not applicable. Thus, we generate adversarial noise for testing by using the LS-PGA attacking method [3], which is specifically designed for the network with binary embeddings.

B. Common Visual Corruptions

Implementation Details. As the optimizer for the classification models, we use Momentum SGD with momentum of 0.9. On the CIFAR-100 dataset, we train the models for 200 epochs with 128 batches. We set the coefficient of weight decay to 5.0×10^{-4} , and the learning rate is scheduled by cosine annealing strategy, which starts from 0.1 and ends at 1.0×10^{-5} . On the ImageNet-1k dataset, we train the models for 180 epochs with 64 batches. The coefficient of weight decay is 1.0×10^{-4} and the initial learning rate is set to 0.1 and divided by 10 at 60 and 120 epochs.

As the optimizer for embedding parameters of P2BE, we use AdamW [26]. The learning rate, β_1 and β_2 are set to 1.0×10^{-4} , 0.999 and 0.999, respectively. We do not apply the scheduling of the learning rate, and the coefficient of weight decay is set to 1.0×10^{-4} . The dimension of binary embedding M is 64.

The total loss is defined as follows:

$$L_{\text{total}} = L_{\text{ce}} + \alpha L_{\text{aug}} + \lambda L_{\text{smooth}}, \quad (10)$$

where L_{ce} is cross-entropy loss for the classification and $\alpha \in \mathbb{R}^+$ is the hyperparameter for `augmix` regularization

TABLE II: The test error of CIFAR-100-C. The all models are trained with Eqn. 10. λ is set to 0 in the case of RGB, One-hot and Thermometer. The values on the table represent mean \pm std across 5 runs. The values in parenthesis are the test error of CIFAR-100 (i.e., clean images).

Encoding	Model	Test error (%)
RGB	WideResNet	35.2 \pm 0.3 (22.2 \pm 0.3)
	DenseNet	37.3 \pm 0.1 (22.6 \pm 0.3)
	ResNeXt	33.9 \pm 0.3 (20.5 \pm 0.4)
One-hot	WideResNet	34.4 \pm 0.4 (23.9 \pm 0.2)
	DenseNet	36.6 \pm 0.2 (23.9 \pm 0.2)
	ResNeXt	33.2 \pm 0.3 (21.8 \pm 0.2)
Thermometer	WideResNet	35.1 \pm 0.2 (22.8 \pm 0.2)
	DenseNet	36.9 \pm 0.2 (23.0 \pm 0.1)
	ResNeXt	33.9 \pm 0.3 (21.3 \pm 0.3)
P2BE (ours)	WideResNet	34.2\pm0.2 (22.8 \pm 0.3)
	DenseNet	36.3 \pm0.2 (23.3 \pm 0.3)
	ResNeXt	32.6\pm0.2 (20.7 \pm 0.4)

loss. $\lambda \in \mathbb{R}^+$ is the hyperparameter controlling the degree of quantization. When the λ is larger, the neighbored embeddings (i.e., w_k and w_{k+1}) tend to have similar directions. For the training of P2NE, the coefficients λ are set to 1.0, and 10.0 on CIFAR-100-C and ImageNet-C, respectively. α is set to 12 which is the same value used in [15].

Results. The result of CIFAR-100-C is shown in Table II. It shows that the approaches of binary embedding generally improve the robustness against visual corruptions with small performance drops on clean images compared to RGB input space. This finding is interesting since the binary embeddings have only been evaluated from the aspect of robustness against adversarial perturbations. It implies that designing a sophisticated input space may be a promising way to improve robustness against never-seen visual corruptions.

As can be seen in Table II, one-hot encoding has the bigger performance drops on clean images and bigger improvements in the robustness against never-seen visual corruptions among the three binary embedding methods. Thermometer encoding has the opposite tendencies: the smaller performance drops on clean images and the smaller robustness improvement against never-seen visual corruptions. P2BE has the good properties of both methods that smaller performance drops on clean images and a bigger improvement of robustness against never-seen visual corruptions.

The result of ImageNet-C is shown in Table III. The results are similar to those in CIFAR-100-C datasets and P2BE shows the best robustness against never-seen visual corruptions. As can be seen in Table III, the binary embedding methods tend to improve robustness against corruptions of the noise and digital categories. However, it has almost no effect on the weather category. This result indicates the limitation of the approach of binary embeddings.

C. Adversarial Perturbations

Implementation Details. The same hyperparameters for the training of CIFAR-100-C in Sec IV-B are used except for some

TABLE III: The clean error of ImageNet-1k and Corruption Error (CE_c) of ImageNet-C. CE_c is the normalized test error by the test error of alexnet. mCE is the averaged CE_c across 15 different corruptions c . The detailed definition of mCE is denoted in Sec. IV-A. The all models are trained with Eqn. 10. λ is set to 0 in the case of RGB, One-hot and Thermometer. The lower the values on the table, the better performances.

	Clean	Noise				Blur				Weather				Digital				mCE
		Gauss	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	Jpeg		
RGB	22.6	66	66	65	69	82	65	66	73	73	62	58	63	80	67	70	68.2	
One-hot	23.5	63	63	59	68	78	64	67	73	72	65	59	64	74	56	68	66.1	
Thermometer	22.9	62	62	61	69	79	65	67	71	71	63	57	63	78	61	68	66.2	
P2BE (ours)	23.3	62	61	61	67	77	63	64	71	71	63	58	63	74	61	67	65.6	

TABLE IV: The test error with LS-PGA white-box attack on CIFAR-10 and CIFAR-100. The values in parenthesis is the test error on clean images.

		One-hot	Thermometer	P2BE (ours)
CIFAR-10	ConTrain	66.1 (11.6)	48.0 (16.0)	45.3 (15.7)
	AdvTrain	52.1 (16.3)	51.6 (15.0)	N/A
CIFAR-100	ConTrain	91.4 (38.1)	73.0 (42.7)	71.5 (42.2)
	AdvTrain	75.2 (41.7)	75.3 (41.0)	N/A

specific parts in adversarial training. Since binary embedding methods are not differentiable, ordinal attacking methods of adversarial perturbations are not applicable. In this work, we use the LS-PGA attacking method [3] which is specially developed to attack binary embeddings. For LS-PGA attacking, we use the seven steps for iterative attack with the annealing rate $\delta = 1.2$ as [3] use. However, we notice that the step size of $\xi = 0.031$ in [3] is too small for the convergence on CIFAR-datasets. In our experiments, we set ξ to 1.0 for the convergence, and the results of one-hot and thermometer encoding are much worse than the scores reported in [3].

To be robust against adversarial perturbations, we use two types of adversarial training methods: AdvTrain and ConTrain. The standard adversarial training is to train the model on only adversarially perturbed images as proposed in [38] and we call it AdvTrain in this work. ConTrain is a variant of the `augmix` training method. We use L_{con} instead of L_{aug} as follows:

$$L_{\text{con}}(p(x); p(x_{\text{adv}})) = \frac{1}{2}(\text{KL}[p(x)||V_{\text{con}}] + \text{KL}[p(x_{\text{adv}})||V_{\text{con}}]), \quad (11)$$

where x and x_{adv} are the original and the adversarially perturbed images, respectively. V_{con} is $\frac{1}{2}(p(x) + p(x_{\text{adv}}))$ and p is the CNN’s prediction from the softmax layer.

The total loss of ConTrain is defined as follows:

$$L_{\text{total}} = L_{\text{ce}} + \alpha L_{\text{con}} + \lambda L_{\text{smooth}}, \quad (12)$$

where the coefficients λ are set to 1.0 and 0.1 on CIFAR-10 and CIFAR-100 datasets, respectively.

Results. We show the results of adversarial perturbations in Table IV. As can be seen in Table IV, P2BE with ConTrain achieves the best robustness against adversarial perturbations on both CIFAR-10 and CIFAR-100. On one-hot and thermometer encoding, AdvTrain and ConTrain are suitable

TABLE V: The test error against LS-PGA white-box attack on CIFAR-10 and CIFAR-100 with various λ . The value in parenthesis is the test error on clean images. λ is the coefficient for embedding smoothness loss in Eqn. 12.

		Test Error (%)
CIFAR-10	ConTrain ($\lambda = 0.0$)	59.3 (12.0)
	ConTrain ($\lambda = 1.0$)	45.3 (15.7)
CIFAR-100	ConTrain ($\lambda = 0.0$)	72.5 (43.4)
	ConTrain ($\lambda = 0.1$)	71.5 (42.2)

to improve the robustness against adversarial perturbations, respectively. Interestingly, P2BE fails to learn with AdvTrain. It seems P2BE requires non-perturbed images for the stability of learning.

V. ANALYSIS

Learned Binary Embedding in P2BE. We show the cosine similarities of binary embedding in Fig 2. The distance space of ImageNet-1k trained P2BE is shown in Fig 2-(d). It shows that the distances of P2BE embeddings are periodically larger and smaller. Such properties of distance space do not exist in RGB, one-hot, and thermometer encoding.

Effect of Dimension Size of Embedding M . We conduct an experiment to investigate the relationship between the robustness performance and dimension size M of P2BE. The result is shown in Fig 3. The worst performance is obtained when M is 128. As can be seen in Table II, it is still better than the results of other baselines. We claim that P2BE is not sensitive to the size of M since the biggest performance gap is 0.5% within various M .

Effect of Embedding Smoothness Loss. We investigate the effectiveness of embedding smoothness loss. Table V shows the performances against adversarial perturbations on CIFAR-10 and CIFAR-100. As can be seen in Table V, P2BE with $\lambda = 0.1$ outperforms P2BE with $\lambda = 0.0$ which is the case without the embedding smoothness loss. This result implies that embedding smoothness loss is an effective regularization of P2BE to improve robustness against adversarial perturbations.

Does P2BE Require Longer Training? We investigate the relationship between the training length and the performances on CIFAR-100-C. The results are shown in Table VI. As can be seen in Table VI, it tends to show better performance with the longer training regardless of input space. Since binary embeddings are learnable in P2BE, the input space is being

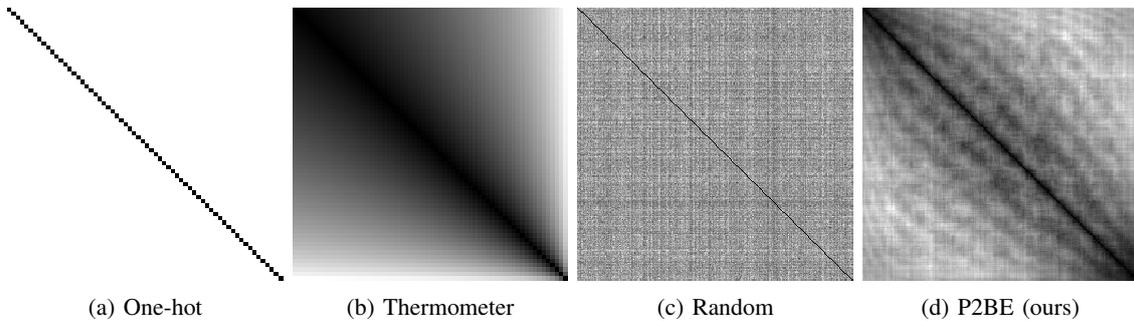


Fig. 2: The cosine similarity of binary embeddings. The vertical and horizontal axis represents that the indices $i, j \in [0, \dots, 255]$ corresponding to the magnitude for each RGB value. In table, the cell at the coordinate (i, j) represents the cosine similarity between binary embeddings e_i and e_j . The black and white colors indicate that the cosine similarities are 1.0 and 0.0, respectively. Figure (c) is the binary embedding generated by the standard normal distribution. Figure (d) is calculated by ImageNet-1k trained P2BE with ResNet50.

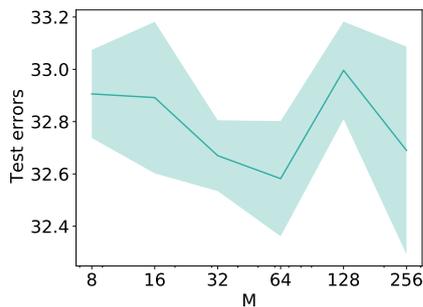


Fig. 3: The CIFAR-100-C results of P2BE across various M . The horizontal axis is M and the vertical axis is the test error of CIFAR-100-C.

TABLE VI: The test error of CIFAR-100-C with different epochs for training. The hyperparameters are all same as in Sec IV-B except for epochs.

Epochs	Network	RGB	One-hot	Thermometer	P2BE (ours)
50	WideResNet	36.2	36.3	36.3	35.9
	DenseNet	39.6	40.1	39.7	40.3
	ResNext	34.6	34.1	33.8	33.9
100	WideResNet	34.9	34.5	35.1	34.4
	DenseNet	37.6	37.5	37.9	38.1
	ResNext	33.5	33.1	34.3	33.6
200	WideResNet	35.4	34.4	35.1	34.2
	DenseNet	37.5	36.6	36.9	36.3
	ResNext	34.2	33.2	33.9	32.5

changed during training, and we have expected that P2BE requires the model to train longer for convergence. However, P2BE shows comparable or better performances with short training periods on the CIFAR-100-C dataset.

Transferability of Learned Binary Embedding in P2BE. It is known that ImageNet-1k pre-trained classification models have good features that are transferable to other tasks. Then, it is reasonable to consider whether the ImageNet-1k trained P2BE is transferable to other tasks or not.

In this work, we verify the transferability of ImageNet-1k trained P2BE by using CIFAR-10, CIFAR-10-C, CIFAR-100,

TABLE VII: The test errors of CIFAR and CIFAR-C datasets with P2BE and fixed ImageNet-1k trained P2BE on WideResNet. Fixed P2BE represents the case that the embedding of P2BE is fixed during the training of the classification model with Eqn. 10.

	Embeddings	Test Error (%)
CIFAR-10	P2BE	4.8
	Fixed P2BE (ImageNet-1k)	4.7
CIFAR-10-C	P2BE	10.3
	Fixed P2BE (ImageNet-1k)	10.1
CIFAR-100	P2BE	23.2
	Fixed P2BE (ImageNet-1k)	22.8
CIFAR-100-C	P2BE	34.2
	Fixed P2BE (ImageNet-1k)	33.9

and CIFAR-100-C. We show the results in Table VII and the fixed ImageNet-1k trained P2BE outperforms the result with P2BE. This result indicates that we may be able to get better binary embeddings by using more complex and large-scale datasets.

VI. CONCLUSION

We propose Pixel to Binary Embedding (P2BE) for improving the robustness of CNNs. P2BE is a learnable binary embedding method as opposed to hand-coded binary embedding methods (e.g., one-hot and thermometer encoding). We show that P2BE achieves the best robustness performances against adversarial perturbations and common visual corruptions among other hand-coded binary embedding methods.

ACKNOWLEDGEMENT

These research results were obtained from the commissioned research (No. 225) by National Institute of Information and Communications Technology (NICT), Japan. This work was also supported by Institute of AI and Beyond of the University of Tokyo and JSPS KAKENHI Grant Number JP19H04166.

REFERENCES

- [1] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," *ICML*, 2018.
- [2] Y. Bengio, N. Léonard, and A. C. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *CoRR*, 2013.
- [3] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow, "Thermometer encoding: One hot way to resist adversarial examples," *ICLR*, 2018.
- [4] N. Carlini and D. A. Wagner, "Towards evaluating the robustness of neural networks," *IEEE Symposium on Security and Privacy*, 2017.
- [5] J. Chen, X. Wu, V. Rastogi, Y. Liang, and S. Jha, "Towards understanding limitations of pixel discretization against adversarial attacks," *2019 IEEE European Symposium on Security and Privacy (EuroS&P)*, 2019.
- [6] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," *NIPS*, 2015.
- [7] S. Darabi, M. Belbahri, M. Courbariaux, and V. P. Nia, "BNN+: improved binary network training," *CoRR*, 2018.
- [8] G. S. Dhillon, K. Azizzadenesheli, J. D. Bernstein, J. Kossaifi, A. Khanna, Z. C. Lipton, and A. Anandkumar, "Stochastic activation pruning for robust adversarial defense," *ICLR*, 2018.
- [9] R. Geirhos, C. R. M. Temme, J. Rauber, H. H. Schütt, M. Bethge, and F. A. Wichmann, "Generalisation in humans and deep neural networks," *NeurIPS*, 2018.
- [10] R. Gong, X. Liu, S. Jiang, T. Li, P. Hu, J. Lin, F. Yu, and J. Yan, "Differentiable soft quantization: Bridging full-precision and low-bit neural networks," *ICCV*, 2019.
- [11] I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *ICLR*, 2015.
- [12] C. Guo, M. Rana, M. Cisse, and L. van der Maaten, "Countering adversarial images using input transformations," *ICLR*, 2018.
- [13] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CVPR*, 2016.
- [14] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," *ICLR*, 2019.
- [15] D. Hendrycks, N. Mu, E. D. Cubuk, B. Zoph, J. Gilmer, and B. Lakshminarayanan, "Augmix: A simple data processing method to improve robustness and uncertainty," *ICLR*, 2020.
- [16] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," *CVPR*, 2017.
- [17] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks," *NIPS*, 2016.
- [18] H. Kim, K. Kim, J. Kim, and J.-J. Kim, "Binaryduo: Reducing gradient mismatch in binary activation network by coupling binary activations," *ICLR*, 2020.
- [19] A. Krizhevsky, "Learning multiple layers of features from tiny images," Tech. Rep., 2009.
- [20] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *NIPS*, 2012.
- [21] A. Kurakin, I. J. Goodfellow, S. Bengio, Y. Dong, F. Liao, M. Liang, T. Pang, J. Zhu, X. Hu, C. Xie, J. Wang, Z. Zhang, Z. Ren, A. L. Yuille, S. Huang, Y. Zhao, Y. Zhao, Z. Han, J. Long, Y. Berdibekov, T. Akiba, S. Tokui, and M. Abe, "Adversarial attacks and defences competition," *The NIPS '17 Competition: Building Intelligent Systems. The Springer Series on Challenges in Machine Learning*, 2017.
- [22] J. Lee, T. Won, T. K. Lee, H. Lee, G. Gu, and K. Hong, "Compounding the performance improvements of assembled techniques in a convolutional neural network," *CoRR*, 2020.
- [23] C. Liu, W. Ding, X. Xia, B. Zhang, J. Gu, J. Liu, R. Ji, and D. David, "Circulant binary convolutional networks: Enhancing the performance of 1-bit dcnn with circulant back propagation," *CVPR*, 2019.
- [24] X. Liu, M. Cheng, H. Zhang, and C. Hsieh, "Towards robust neural networks via random self-ensemble," *ECCV*, 2018.
- [25] Z. Liu, B. Wu, W. Luo, X. Yang, W. Liu, and K. Cheng, "Bi-real net: Enhancing the performance of 1-bit dcnn with improved representational capability and advanced training algorithm," *ECCV*, 2018.
- [26] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *ICLR*, 2019.
- [27] X. Ma, B. Li, Y. Wang, S. M. Erfani, S. Wijewickrema, G. Schoenebeck, M. E. Houle, D. Song, and J. Bailey, "Characterizing adversarial subspaces using local intrinsic dimensionality," *ICLR*, 2018.
- [28] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *ICLR*, 2018.
- [29] C. Michaelis, B. Mitzkus, R. Geirhos, E. Rusak, O. Bringmann, A. S. Ecker, M. Bethge, and W. Brendel, "Benchmarking robustness in object detection: Autonomous driving when winter is coming," *CoRR*, 2019.
- [30] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," *CVPR*, 2017.
- [31] S. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: A simple and accurate method to fool deep neural networks," *CVPR*, 2016.
- [32] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," *ECCV*, 2016.
- [33] E. Rusak, L. Schott, R. S. Zimmermann, J. Bitterwolf, O. Bringmann, M. Bethge, and W. Brendel, "A simple way to make neural networks robust against diverse image corruptions," *ECCV*, 2020.
- [34] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *IJCV*, 2015.
- [35] C. Song, K. He, J. Lin, L. Wang, and J. E. Hopcroft, "Robust local features for improving the generalization of adversarial training," *ICLR*, 2020.
- [36] Y. Song, T. Kim, S. Nowozin, S. Ermon, and N. Kushman, "Pixeldefend: Leveraging generative models to understand and defend against adversarial examples," *ICLR*, 2018.
- [37] J. Su, D. V. Vargas, and K. Sakurai, "One pixel attack for fooling deep neural networks," *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 5, pp. 828–841, 2019.
- [38] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *ICLR*, 2014.
- [39] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, "Robustness may be at odds with accuracy," *ICLR*, 2019.
- [40] I. Vasiljevic, A. Chakrabarti, and G. Shakhnarovich, "Examining the impact of blur on recognition by convolutional networks," *CoRR*, 2016.
- [41] C. Xie, J. Wang, Z. Zhang, Z. Ren, and A. Yuille, "Mitigating adversarial effects through randomization," *ICLR*, 2018.
- [42] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," *CVPR*, 2017.
- [43] W. Xu, D. Evans, and Y. Qi, "Feature squeezing: Detecting adversarial examples in deep neural networks," *25th Annual Network and Distributed System Security Symposium*, 2018.
- [44] S. Zagoruyko and N. Komodakis, "Wide residual networks," *BMVC*, 2016.
- [45] D. Zhang, T. Zhang, Y. Lu, Z. Zhu, and B. Dong, "You only propagate once: Accelerating adversarial training via maximal principle," *NeurIPS*, 2019.
- [46] D. Zhang, J. Yang, D. Ye, and G. Hua, "Lq-nets: Learned quantization for highly accurate and compact deep neural networks," *ECCV*, 2018.
- [47] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. E. Ghaoui, and M. Jordan, "Theoretically principled trade-off between robustness and accuracy," *ICML*, 2019.
- [48] Y. Zhang and P. Liang, "Defending against whitebox adversarial attacks via randomized discretization," *PMLR*, 2019.
- [49] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou, "Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients," *CoRR*, 2016.