

TOWARDS ADVERSARIALLY ROBUST CONTINUAL LEARNING

Tao Bai^{*1} Chen Chen² Lingjuan Lyu^{2†} Jun Zhao³ Bihan Wen¹

¹ School of Electrical and Electronic Engineering, Nanyang Technological University

² Sony AI

³ School of Computer Science and Engineering, Nanyang Technological University

ABSTRACT

Recent studies show that models trained by continual learning can achieve the comparable performances as the standard supervised learning and the learning flexibility of continual learning models enables their wide applications in the real world. Deep learning models, however, are shown to be vulnerable to adversarial attacks. Though there are many studies on the model robustness in the context of standard supervised learning, protecting continual learning from adversarial attacks has not yet been investigated. To fill in this research gap, we are the first to study adversarial robustness in continual learning and propose a novel method called **Task-Aware Boundary Augmentation (TABA)** to boost the robustness of continual learning models. With extensive experiments on CIFAR-10 and CIFAR-100, we show the efficacy of adversarial training and TABA in defending adversarial attacks.

Index Terms— Adversarial training, continual learning, data augmentation

1. INTRODUCTION

Continual learning studies the problem of learning from an infinite stream of data, with the goal of gradually extending acquired knowledge and using it for future learning [1]. The major challenge is to learn without catastrophic forgetting: performance on a previously learned task or domain should not significantly degrade over time when new tasks or domains are added. To this end, researchers have proposed various methods [2, 3, 4, 5] to reduce the computational costs while maintaining the performances for learned tasks. As such, continual learning has made a wide range of real-world applications into reality recently [6, 7].

Though the training process of continual learning is quite different from regular supervised learning, the model trained with continual learning is exactly the same as the regular supervised learning during inference. Recent studies [8, 9] on adversarial examples reveal the vulnerabilities of well-trained deep learning models, which are easy to break through. Thus,

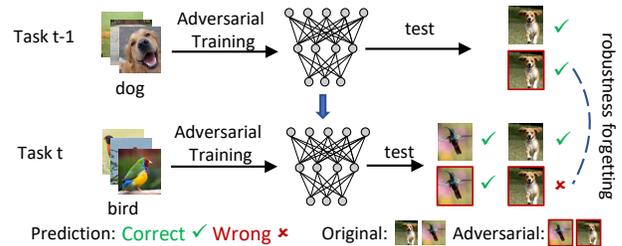


Fig. 1. Robust continual learning and the issue of robustness forgetting. (e.g. “adversarial dog” is predicted wrong after training on Task t).

it’s natural to assume that models trained with continual learning suffer from adversarial examples as well. Considering the real-world applications of continual learning models, it is essential to protect continual learning models against adversarial attacks. There have been a number of studies exploring how to secure the deep learning models against adversarial examples [10, 11], but surprisingly, protecting continual learning from adversarial attacks has not been fully studied.

To bridge the gap between continual learning and adversarial robustness, we focus on the replay-based continual learning methods and take the first step to develop robust continual learning methods. As we stated above, nevertheless, data from previously learned tasks in continual learning are partially accessible, causing the imbalance between previous tasks and new tasks. In this case, models trained in current stage usually tend to overfit the new class data. As such, the catastrophic forgetting of adversarial robustness is inevitable in robust continual learning, *i.e.* when leveraging adversarial training [10, 11] in continual learning for adversarial robustness (see Fig. 1) Preventing forgetting, or in other words, preserving learned knowledge refers to maintaining the previously learned decision boundary among classes [12]. We thus propose a novel approach called **Task-Aware Boundary Augmentation (TABA)** for maintaining the decision boundaries for adversarial training in continual learning settings.

Our contributions are summarized as follows:

1. To the best of our knowledge, we are the first to investigate the security issues in continual learning and

^{*}Work done during internship at Sony AI.

[†]Corresponding Author.

improving the adversarial robustness by leveraging adversarial training.

2. We further identify the catastrophic forgetting of adversarial robustness and propose a novel approach called **Task-Aware Boundary Augmentation (TABA)** for enhancing adversarial training and continual learning.
3. With experiments on popular datasets like CIFAR10 and CIFAR100, we show the efficacy of TABA in different continual learning scenarios.

2. RELATED WORKS

Continual learning is widely studied in the last few years, which assumes data comes in a sequential way [2, 3, 13]. There are, however, only a few works studying the security issues in continual learning [14, 15]. It is empirically shown the importance of robust features in continual learning [16]. Authors of [17] proposed to incorporate adversarial training with continual learning to enhance the robustness, as adversarial training has been validated in other deep learning tasks [18, 19, 11]. In this paper, we study how to leverage adversarial training in continual learning and alleviate the catastrophic forgetting of adversarial robustness.

3. APPROACH

3.1. Problem Definition

In this work, we focus on the robust multi-class classification problem, which involves the sequential learning of \mathcal{T} stages/tasks consisting of disjoint class sets. Formally, at learning stage $t \in \{2, \dots, \mathcal{T}\}$, given a model trained on an old dataset \mathcal{X}_o^{t-1} from stage $\{1, \dots, t-1\}$, our goal is to learn a unified classifier for both old classes \mathcal{C}_o and new classes \mathcal{C}_n . The training data at stage t is denoted as $\mathcal{X}^t = \mathcal{X}_n^t \cup \tilde{\mathcal{X}}_o^{t-1}$, where $\tilde{\mathcal{X}}_o^{t-1}$ is a tiny subset of \mathcal{X}_o^{t-1} . Thus, the challenge in continual learning is retraining the original model with the severely imbalanced \mathcal{X}^t to boost the robustness on all seen classes while avoiding catastrophic forgetting.

3.2. Revisiting Distillation for Catastrophic Forgetting

Knowledge distillation [20] is firstly introduced to continual learning by Learning without forgetting (LwF) [21] and adapted by iCaRL [4] for the *multi-class* continual learning problem. Typically, the loss function of such distillation-based methods consists of two terms for each training sample x : the classification loss \mathcal{L}_{ce} and the distillation loss \mathcal{L}_{dis} . Specifically, the classification loss \mathcal{L}_{ce} is expressed as

$$\mathcal{L}_{ce}(x) = - \sum_{i=1}^{|\mathcal{C}|} y_i \log(p_i), \quad (1)$$

where $\mathcal{C} = \mathcal{C}_o \cup \mathcal{C}_n$, y_i is the i_{th} value of the one-hot ground truth y , and p_i is the i_{th} value of predicted class probability p . The goal of \mathcal{L}_{dis} is to preserve knowledge obtained from previous data, which is expressed as

$$\mathcal{L}_{dis}(x) = - \sum_{i=1}^{|\mathcal{C}_o|} (p^*) \log(p), \quad (2)$$

where p^* is the soft label of x generated by the old model. It, however, is observed in [4] that there is tendency of classifying test samples to new classes by LwF. Thus, iCaRL utilized *herd selection* to better approximate the class mean vector of old classes, where samples that are close to the center of old classes are selected.

Recall that our goal is to obtain a robust model trained in the continual learning manner. To gain robustness, adversarial training is inevitable, which requires augmenting datasets with adversarial examples in every training iteration. Following the definition of continual learning, we can derive the loss function of **Robust Continual Learning (RCL)**. With adversarial training, we should replace the input x in Equation (1) and (2) with its adversarial counterpart x_{adv} , which is solved by

$$x_{adv} = \underset{\|x_{adv}-x\|_p \leq \epsilon}{\operatorname{argmax}} (\mathcal{L}_{ce}(x_{adv})), \quad (3)$$

where ϵ is the allowed magnitude of perturbations in p -norm. Thus, the loss function of robust continual learning would be

$$\mathcal{L}_{RCL} = \mathcal{L}_{ce}(x_{adv}) + \mathcal{L}_{dis}(x_{adv}) \quad (4)$$

Nevertheless, simply combining adversarial training with continual learning is not enough. From the perspective of adversarial training, centered exemplars are not helpful for the forgetting of adversarial robustness. Recent studies [22, 23] pointed out that not all data points contribute equally during adversarial training and samples that are close to the decision boundaries should be emphasised. Therefore, how to deal with the exemplar set during adversarial training is essential for robust continual learning. In addition, adversarial training is more data-hungry than standard training. The significant imbalance between old classes and new classes can be more severe. In this work, we aim to tackle these problems by incorporating data augmentation with adversarial training.

3.3. Task-Aware Boundary Augmentation

Preventing catastrophic forgetting of adversarial robustness in continual learning is equivalent to maintaining the decision boundary learned by adversarial training. One direct way to do so is to introduce some samples close to the decision boundaries to the exemplar set (named Boundary Exemplar in Section 4 and Table 1). However, this makes the exemplar selection process more sophisticated because the ratio of centered samples and boundary samples is hard to decide. In addition, such mixed exemplar set may have negative influence on the approximation of old classes, which may downgrade the model performance. Another potential solution is

Mixup [24], where the dataset is augmented by interpolating different samples linearly. Mixup, however, is not specially designed for adversarial training or continual learning. It breaks the local invariance of adversarially trained models by linear interpolation and worsens the imbalance between old tasks and new tasks.

Inspired by Mixup, we propose **Task-Aware Boundary Augmentation (TABA)** to augment the training data \mathcal{X} by synthesizing more boundary data, which can be plugged in RCL easily. Compared to Mixup, TABA is specially designed for adversarial training and continual learning. The differences are summarized as below. *First*, TABA doesn't select samples in the whole dataset but from the boundary data. The reason is that boundary data is easier to attack and contributes more to adversarial robustness [22]. We can obtain the boundary data for free because adversarial training requires generating adversarial examples. Misclassified samples in the previous iteration are marked as the boundary data, which is denoted by \mathcal{B} . *Second*, to deal with the data imbalance issue in continual learning, TABA selects samples from two sets: one is boundary data from $\tilde{\mathcal{X}}_o^t$ and the other is boundary data from \mathcal{X}_n^t , denoted as \mathcal{B}_o and \mathcal{B}_n , respectively. In this way, the augmented data can help maintain the learned decision boundaries in the previous stage. *Third*, we restrict the interpolation weight λ to a interval of $[0.45, 0.55]$ rather than $[0, 1]$ in Mixup to avoid the linearity, which is decided empirically. The augmented samples can also be closer to the decision boundaries, compared to samples provided by Mixup.

The augmented sample (\bar{x}, \bar{y}) by our TABA can be defined as follows:

$$\begin{aligned}\bar{x} &= \lambda x_o + (1 - \lambda)x_n \\ \bar{y} &= \lambda y_o + (1 - \lambda)y_n,\end{aligned}\tag{5}$$

where λ is the interpolation weight, $(x_o, y_o) \in \mathcal{B}_o$ and $(x_n, y_n) \in \mathcal{B}_n$.

Accordingly, the final loss function of RCL with TABA (RCL-TABA) would be

$$\begin{aligned}\mathcal{L}_{final} &= \mathcal{L}_{TABA} + \mathcal{L}_{RCL} \\ \mathcal{L}_{TABA} &= \mathcal{L}_{ce}(\bar{x}_{adv}) + \mathcal{L}_{dis}(\bar{x}_{adv}).\end{aligned}\tag{6}$$

The training details of RCL-TABA are in Algorithm 1.

4. EXPERIMENTS

4.1. Settings

Datasets. We conduct our experiments on two popular datasets: CIFAR-10 and CIFAR-100 [25]. A common setting is to train the model on data with equal classes in each stage (**Setting I**). Based on this, we set five stages for both CIFAR-10 and CIFAR-100, i.e., 2/20 classes in each stage. In addition, we further take the unequal-class scenario for different stages (**Setting II**), which is more realistic in practice. The classes for each stage is randomly sampled and we make

Algorithm 1 Robust continual learning with task-aware boundary augmentation (RCL-TABA)

```

1: Randomly initialize model  $f^0$ , old task data  $\tilde{\mathcal{X}}_o^0 = \emptyset$ 
2: for  $t = \{1, \dots, \mathcal{T}\}$  do
3:   Input: model  $f^{t-1}$ , new task data  $\mathcal{X}_n^t$ , training epochs  $E$ , number of batches  $M$ , original batch size  $m$ , interpolation batch size  $m'$ 
4:   Output: model  $f^t$ 
5:    $f^t \leftarrow f^{t-1}$ ,  $\mathcal{X}^t = \mathcal{X}_n^t \cup \tilde{\mathcal{X}}_o^{t-1}$ ,  $\mathcal{B}_0 = \mathcal{X}^t$ 
6:   for  $e = \{1, \dots, E\}$  do
7:      $\mathcal{B}_e = \emptyset$ 
8:     Compute augmentation set  $\tilde{\mathcal{X}}^t$  from  $\mathcal{B}_{e-1}$  by Eq. (5)

9:     for  $mini - batch = \{1, \dots, M\}$  do
10:      Randomly sample  $\{(x_i, y_i)\}_{i=1}^m$  from  $\mathcal{X}$ 
11:      for  $i = \{1, \dots, m\}$  do
12:        Generate adversarial data  $x_i^{adv}$  by Eq. (3)
13:        if  $f(x_i^{adv}) \neq y_i$  then
14:           $\mathcal{B}_e \leftarrow \mathcal{B}_e \cup (x_i, y_i)$ 
15:        end if
16:      end for
17:      Randomly sample  $\{(\bar{x}_i, \bar{y}_i)\}_{i=m+1}^{m+m'}$  from  $\tilde{\mathcal{X}}$ 
18:      for  $i = \{m+1, \dots, m+m'\}$  do
19:        Generate adversarial data  $\bar{x}_i^{adv}$  by Eq. (3)
20:      end for
21:      optimize  $f_t$  on  $\{(\bar{x}_i, \bar{y}_i)\}_{i=1}^{m+m'}$  by Eq. (6)
22:    end for
23:  end for
24:  update  $\tilde{\mathcal{X}}_o^t$  by class using herd selection [4]
25: end for

```

sure there is no overlap between different stages. Note that Setting II is only for CIFAR-100, where the variance of class numbers is large enough for observation.

Implementation Details. All the models are implemented with PyTorch and trained on NVIDIA Tesla V100. We use ResNet18 [26] as our backbone model for experiments. For adversarial training on both datasets, we set the maximal magnitude of perturbations ϵ to 8/255 and utilize the 7-step Projected Gradient Descent (PGD) to generate adversarial examples, where the step size is 2/255. For evaluation, we not only test the standard accuracy on clean samples but also the robust accuracy with adversarial attacks. We denote the standard accuracy as **SA**, robust accuracy under PGD attacks as **RA(PGD)**, and robust accuracy under AutoAttack as **RA(AA)**, respectively. The ϵ and parameters of PGD attacks for evaluation is set to be the same as for training.

During training, the class order for datasets is fixed for fair comparisons. For reserving samples in previous stages, we use *herd selection strategy* in [4] and set the memory capacity to be 2000 samples for both CIFAR-10 and CIFAR-100. The capacity is independent of the number classes and the number of exemplar for each class is $\frac{2000}{\# \text{of seen classes}}$.

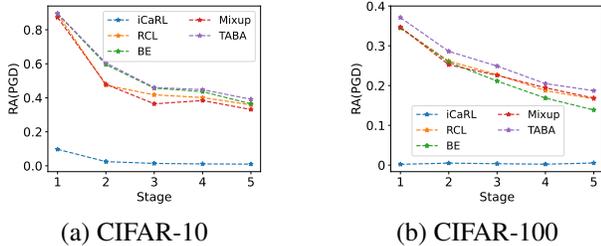


Fig. 2. Robustness evaluation on all seen classes at different stages. BE is short for Boundary Exemplar to save space.

Baselines. As we stated, the adversarial robustness of continual learning is firstly studied in this paper and there is no previous work on this topic. Thus, we choose *iCaRL*, the representative method for continual learning as the baseline. To obtain adversarial robustness, we adopt adversarial training in continual learning and build upon *iCaRL*, named *RCL*, as another baseline. In addition, we introduce *Boundary Exemplar* to verify the influence of boundary data for *RCL* and *Mixup*, which is closely related to *TABA*. *Boundary Exemplar*, *Mixup* and *TABA* are augmentation methods for improving *RCL*.

4.2. Experimental Results

First, we conduct experiments on CIFAR-10 and CIFAR-100 in Setting I. Robustness changes over stages are visualized in Fig. 2 and the experimental results are summarized in Table. 1. We can observe that models trained by *iCaRL* are not robust under all adversarial attacks, showing nearly 0 robust accuracy against PGD attack, and 0 robust accuracy against AutoAttack. With adversarial training, the adversarial robustness for continual learning models is greatly improved, though there is a drop of standard accuracy. Compared to all other methods, our *TABA* clearly shows strong performances: On both CIFAR-10 and CIFAR-100, *TABA* shows the best or second best robustness under PGD attacks and AutoAttacks while maintaining the standard accuracy. Though *Mixup* achieves the highest robustness under AutoAttack on CIFAR10, it brings a large drop of 20% for standard accuracy.

Second, we conduct experiments in Setting II on CIFAR-100. In this setting, the class numbers for each stage are randomly selected and the sum of classes in all stages is guaranteed to be 100. We run the experiments for 3 times and the class numbers for different stages varies from 5 to 45. The average results are reported in Table. 2 (the variance are close to zero and not reported here). We can see that *TABA* achieves the best overall performances. Compared to *Mixup*, *TABA* has comparable RA(AA) and much higher SA. The large drop of SA in *Mixup* should be avoided.

Table 1. Robustness evaluation on CIFAR-10 and CIFAR-100 in Setting I. The best results (the higher, the better) in each column are in **bold text**.

		SA	RA(PGD)	RA(AA)
CIFAR10	iCaRL	67.17%	1.00%	0.00%
	RCL	60.36%	36.83%	16.71%
	Boundary Exemplar	66.52%	36.91%	10.88%
	Mixup	46.96%	33.11%	20.36%
	TABA	65.97%	38.41%	19.74%
CIFAR100	iCaRL	58.31%	0.53%	0.00%
	RCL	46.67%	16.67%	9.99%
	Boundary Exemplar	38.08%	14.15%	6.51%
	Mixup	46.58%	16.86%	10.03%
	TABA	45.16%	18.71%	11.21%

Table 2. Robustness evaluation on CIFAR-100 in Setting II. The best results (the higher, the better) in each column are in **bold text**.

Method	SA	RA(PGD)	RA(AA)
iCaRL	49.68%	0.04%	0.01%
RCL	44.55%	17.49%	9.71%
Mixup	28.53%	16.08%	11.77%
TABA	42.79%	18.72%	11.43%

Table 3. Effects of three modifications in *TABA*.

Boundary	Task-aware	λ	SA	RA(PGD)	RA(AA)
✗	✗	✗	46.96%	33.11%	20.36%
✓	✗	✗	54.84%	31.09%	15.45%
✓	✓	✗	59.61%	32.18%	15.87%
✓	✓	✓	65.97%	38.41%	19.74%

✓: w/ ✗: w/o

4.3. Ablation Study

Inspired by *Mixup*, we propose *TABA* for relieving the forgetting of adversarial robustness in continual learning. Compared to *Mixup*, *TABA* is different in three ways: *boundary data*, *task-aware sample selection* and *the range of λ* . Here we investigate the effects of these modifications and results are summarized in Table. 3. We can see the improvements when we make modifications sequentially on *Mixup*.

5. CONCLUSION

In this paper, we study the continual learning problem in the adversarial settings. It is verified that models trained in continual learning ways are also vulnerable to adversarial examples. We thus propose *RCL-TABA*, which consists of adversarial training and a novel data augmentation method *TABA*, to secure continual learning. As this is the very first step to studying the intersection of adversarial training and continual learning, we hope our findings provide useful insights and motivate researchers to explore deeper.

6. REFERENCES

- [1] Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars, “A continual learning survey: Defying forgetting in classification tasks,” *IEEE TPAMI*, p. 1–1, 2021.
- [2] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al., “Overcoming catastrophic forgetting in neural networks,” *PNAS*, 2017.
- [3] Sang-Woo Lee, Jin-Hwa Kim, Jaehyun Jun, Jung-Woo Ha, and Byoung-Tak Zhang, “Overcoming catastrophic forgetting by incremental moment matching,” *NIPS*, vol. 30, 2017.
- [4] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert, “icarl: Incremental classifier and representation learning,” in *CVPR*, 2017, pp. 2001–2010.
- [5] David Lopez-Paz and Marc’Aurelio Ranzato, “Gradient episodic memory for continual learning,” *NIPS*, vol. 30, 2017.
- [6] Cecilia S Lee and Aaron Y Lee, “Clinical applications of continual learning machine learning,” *The Lancet Digital Health*, vol. 2, no. 6, pp. e279–e281, 2020.
- [7] Pankaj Gupta, Yatin Chaudhary, Thomas Runkler, and Hinrich Schuetze, “Neural topic modeling with continual lifelong learning,” in *ICML*, 2020.
- [8] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus, “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
- [9] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [10] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu, “Towards Deep Learning Models Resistant to Adversarial Attacks,” in *ICLR*. 2018, OpenReview.net.
- [11] Tao Bai, Jinqi Luo, Jun Zhao, Bihan Wen, and Qian Wang, “Recent Advances in Adversarial Training for Adversarial Robustness,” in *IJCAI-21*, 2021.
- [12] Fei Zhu, Zhen Cheng, Xu-yao Zhang, and Cheng-lin Liu, “Class-Incremental Learning via Dual Augmentation,” in *NeurIPS*, 2021, vol. 34, pp. 14306–14318.
- [13] Jiahua Dong, Lixu Wang, Zhen Fang, Gan Sun, Shichao Xu, Xiao Wang, and Qi Zhu, “Federated class-incremental learning,” in *CVPR*, 2022.
- [14] Hikmat Khan, Pir Masoom Shah, Syed Farhan Alam Zaidi, et al., “Susceptibility of continual learning against adversarial attacks,” *arXiv preprint arXiv:2207.05225*, 2022.
- [15] Yunhui Guo, Mingrui Liu, Yandong Li, Liqiang Wang, Tianbao Yang, and Tajana Rosing, “Attacking lifelong learning models with gradient reversion,” 2020.
- [16] Hikmat Khan, Nidhal Carla Bouaynaya, and Ghulam Rasool, “Adversarially robust continual learning,” in *IJCNN*, 2022, pp. 1–8.
- [17] Ting-Chun Chou, Jhih-Yuan Huang, and Wei-Po Lee, “Continual learning with adversarial training to enhance robustness of image recognition models,” in *2022 International Conference on Cyberworlds (CW)*, 2022, pp. 236–242.
- [18] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu, “Towards deep learning models resistant to adversarial attacks,” in *ICLR*, 2018.
- [19] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I. Jordan, “Theoretically principled trade-off between robustness and accuracy,” in *ICML*, 2019.
- [20] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al., “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, vol. 2, no. 7, 2015.
- [21] Zhizhong Li and Derek Hoiem, “Learning without forgetting,” *IEEE TPAMI*, 2017.
- [22] Jingfeng Zhang, Jianing Zhu, Gang Niu, Bo Han, Masashi Sugiyama, and Mohan Kankanhalli, “Geometry-aware instance-reweighted adversarial training,” *arXiv preprint arXiv:2010.01736*, 2020.
- [23] Chen Chen, Jingfeng Zhang, Xilie Xu, Lingjuan Lyu, Chaochao Chen, Tianlei Hu, and Gang Chen, “Decision Boundary-aware Data Augmentation for Adversarial Training,” *IEEE IDSC*, pp. 1–1, 2022.
- [24] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz, “mixup: Beyond empirical risk minimization,” in *ICLR*, 2018.
- [25] Alex Krizhevsky, Geoffrey Hinton, et al., “Learning multiple layers of features from tiny images,” 2009.
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.