# VISUAL PROMPTING FOR ADVERSARIAL ROBUSTNESS

Aochuan Chen\*,1

Peter Loren $z^{\star,2}$  Yu

<sup>2</sup> Yuguang Yao<sup>1</sup>

Sijia Liu<sup>1</sup>

<sup>1</sup>Michigan State University, USA

<sup>2</sup>Fraunhofer ITWM and Fraunhofer Center of Machine Learning, Germany

<sup>3</sup>IBM Research, USA

## ABSTRACT

In this work, we leverage visual prompting (VP) to improve adversarial robustness of a fixed, pre-trained model at test time. Compared to conventional adversarial defenses, VP allows us to design universal (*i.e.*, data-agnostic) input prompting templates, which have plug-and-play capabilities at test time to achieve desired model performance without introducing much computation overhead. Although VP has been successfully applied to improving model generalization, it remains elusive whether and how it can be used to defend against adversarial attacks. We investigate this problem and show that the vanilla VP approach is *not* effective in adversarial defense since a universal input prompt lacks the capacity for robust learning against sample-specific adversarial perturbations. To circumvent it, we propose a new VP method, termed Class-wise Adversarial Visual Prompting (C-AVP), to generate class-wise visual prompts so as to not only leverage the strengths of ensemble prompts but also optimize their interrelations to improve model robustness. Our experiments show that C-AVP outperforms the conventional VP method, with 2.1× standard accuracy gain and 2× robust accuracy gain. Compared to classical test-time defenses, C-AVP also yields a  $42 \times$  inference time speedup. Code is available at https://github.com/Phoveran/vp-for-adversarial-robustness.

*Index Terms*— visual prompting, adversarial defense, adversarial robustness

### 1. INTRODUCTION

Machine learning (ML) models can easily be manipulated (by an adversary) to output drastically different classifications. Thereby, model robustification against adversarial attacks is now a major focus of research. Yet, a large volume of existing works focused on training recipes and/or model architectures to gain robustness. Adversarial training (AT) [1], one of the most effective defense, adopted min-max optimization to minimize the worst-case training loss induced by adversarial attacks. Extended from AT, various defense methods were proposed, ranging from supervised learning to semi-supervised learning, and further to unsupervised learning [2–11]. Although the design for robust training has made tremendous success in improving model robustness [12, 13], it typically takes an intensive computation cost with poor defense scalability to a fixed, pre-trained ML model. Towards circumventing this difficulty, the problem of test-time defense arises; see the seminal work in Croce *et. al.* [14]. Test-time defense alters either a test-time input example or a small portion of the pre-trained model. Examples include input (anti-adversarial) purification [15–17] and model refinement by augmenting the pre-trained model with auxiliary components [18–20]. However, these defense techniques inevitably raise the inference time and hamper the test-time efficiency [14]. Inspired by that, our work will advance the test-time defense technology by leveraging the idea of *visual prompting* (**VP**) [21], also known as model reprogramming [22–25].

Pin-Yu Chen<sup>3</sup>

Generally speaking, VP [21] creates a universal (i.e., dataagnostic) input prompting template (in terms of input perturbations) in order to improve the generalization ability of a pre-trained model when incorporating such a visual prompt into test-time examples. It enjoys the same idea as model reprogramming [22–25] or unadversarial example [26], which optimizes a universal perturbation pattern to maneuver (i.e., reprogram) the functionality of a pre-trained model towards the desired criterion, *e.g.*, cross-domain transfer learning [24], out-of-distribution generalization [26], and fairness [25]. However, it remains elusive whether or not VP could be designed as an effective solution to adversarial defense. We will investigate this problem, which we call adversarial visual prompting (AVP) in this work. Compared to conventional test-time defense methods, AVP significantly reduces the inference time overhead since visual prompts can be designed offline over training data and have the plug-and-play capability applied to any testing data. We summarize our contributions as below. • We formulate and investigate the problem of AVP for the first time and empirically show the conventional data-agnostic VP design is incapable of gaining adversarial robustness.

**②** We propose a new VP method, termed class-wise AVP (**C**-**AVP**), which produces multiple, class-wise visual prompts with explicit optimization on their couplings to gain better adversarial robustness.

• We provide insightful experiments to demonstrate the pros and cons of VP in adversarial defense.

<sup>\*</sup> Equal contribution.

#### 2. RELATED WORK

**Visual prompting.** Originated from the idea of in-context learning or prompting in natural language processing (NLP) [27–30], VP was first proposed in Bahng *et. al.* [21] for vision models. Before formalizing VP in Bahng *et. al.* [21], the underlying prompting technique has also been devised in computer vision (CV) with different naming. For example, VP is closely related to *adversarial reprogramming* or *model reprogramming* [22–24, 31–33], which focused on altering the functionality of a fixed, pre-trained model across domains by augmenting test-time examples with an additional (universal) input perturbation pattern. *Unadversarial learning* also enjoys the similar idea to VP. In [26], unadversarial examples that perturb original ones using 'prompting' templates were introduced to improve out-of-distribution generalization. Yet, the problem of VP for adversarial defense is under-explored.

Adversarial defense. The lack of adversarial robustness is a weakness of ML models. Adversarial defense, such as adversarial detection [19, 34–38] and robust training [2, 6, 9, 10, 18, 39], is a current research focus. In particular, adversarial training (AT) [1] is the most widely-used defense strategy and has inspired many recent advances in adversarial defense [12, 13, 20, 40-42]. However, these AT-type defenses (with the goal of robustness-enhanced model training) are computationally intensive due to min-max optimization over model parameters. To reduce the computation overhead of robust training, the problem of test-time defense arises [14], which aims to robustify a given model via lightweight unadversarial input perturbations (a.k.a input purification) [15,43] or minor modifications to the fixed model [44,45]. In different kinds of test-time defenses, the most relevant work to ours is anti-adversarial perturbation [17].

## 3. PROBLEM STATEMENT

**Visual prompting.** We describe the problem setup of VP following Bahng *et. al.* [21, 23–25]. Specifically, let  $\mathcal{D}_{tr}$  denote a training set for supervised learning, where  $(\mathbf{x}, y) \in \mathcal{D}_{tr}$  signifies a training sample with feature  $\mathbf{x}$  and label y. And let  $\boldsymbol{\delta}$  be a visual prompt to be designed. The prompted input is then given by  $\mathbf{x}+\boldsymbol{\delta}$  with respect to (w.r.t.)  $\mathbf{x}$ . Different from the problem of adversarial attack generation that optimizes  $\boldsymbol{\delta}$  for erroneous prediction, VP drives  $\boldsymbol{\delta}$  to minimize the performance loss  $\ell$  of a pre-trained model  $\boldsymbol{\theta}$ . This leads to

$$\begin{array}{l} \underset{\boldsymbol{\delta}}{\text{minimize}} \quad \mathbb{E}_{(\mathbf{x},y)\in\mathcal{D}_{\text{tr}}}[\ell(\mathbf{x}+\boldsymbol{\delta};y,\boldsymbol{\theta})] \\ \text{subject to} \quad \boldsymbol{\delta}\in\mathcal{C}, \end{array}$$
(1)

where  $\ell$  denotes prediction error given the training data  $(\mathbf{x}, y)$ and base model  $\theta$ , and C is a perturbation constraint. Following Bahng *et. al.* [21, 23, 24], C restricts  $\delta$  to let  $\mathbf{x} + \delta \in [0, 1]$ for any  $\mathbf{x}$ . Projected gradient descent (PGD) [1, 26] can then be applied to solving problem (1). In the evaluation,  $\delta$  is integrated into test data to improve the prediction ability of  $\theta$ . **Adversarial visual prompting.** Inspired by the usefulness of VP to improve model generalization [21, 24], we ask: (AVP problem) Can VP (1) be extended to robustify  $\theta$  against adversarial attacks?

At the first glance, the AVP problem seems trivial if we specify the performance loss  $\ell$  as the adversarial training loss [1,2]:

$$\ell_{\mathrm{adv}}(\mathbf{x} + \boldsymbol{\delta}; y, \boldsymbol{\theta}) = \underset{\mathbf{x}': \|\mathbf{x}' - \mathbf{x}\|_{\infty} \le \epsilon}{\operatorname{maximize}} \ell(\mathbf{x}' + \boldsymbol{\delta}; y, \boldsymbol{\theta}), \quad (2)$$

where  $\mathbf{x}'$  denotes the adversarial input that lies in the  $\ell_{\infty}$ -norm ball centered at  $\mathbf{x}$  with radius  $\epsilon > 0$ .

Recall from (1) that the conventional VP requests  $\delta$  to be universal across training data. Thus, we term *universal AVP* (**U-AVP**) the following problem by integrating (1) with (2):

$$\begin{array}{l} \underset{\delta:\delta\in\mathcal{C}}{\text{minimize}} \quad \lambda \mathbb{E}_{(\mathbf{x},y)\in\mathcal{D}_{\text{tr}}}[\ell(\mathbf{x}+\delta;y,\boldsymbol{\theta})] + \\ \mathbb{E}_{(\mathbf{x},y)\in\mathcal{D}_{\text{tr}}}[\ell_{\text{adv}}(\mathbf{x}+\delta;y,\boldsymbol{\theta})] \end{array}$$
(U-AVP)

where  $\lambda > 0$  is a regularization parameter to strike a balance between generalization and adversarial robustness [2].



**Fig. 1**: Example of designing U-AVP for adversarial defense on (CIFAR-10, ResNet18), measured by robust accuracy against PGD attacks [1] of different steps. The robust accuracy of 0 steps is the standard accuracy.

The problem (U-AVP) can be effectively solved using a standard minmax optimization method, which involves two alternating optimization routines: inner maximization and outer minimization. The former generates adversarial examples as AT, and the latter produces the visual prompt  $\delta$  like (1). At test time, the effectiveness of  $\delta$  is measured from two aspects: (1) standard accuracy, *i.e.*, the accuracy

of  $\delta$ -integrated benign examples, and (2) robust accuracy, *i.e.*, the accuracy of  $\delta$ -integrated adversarial examples (against the victim model  $\theta$ ). Despite the succinctness of (U-AVP), Fig. 1 shows its *ineffectiveness* to defend against adversarial attacks. Compared to the vanilla VP (1), it also suffers a significant standard accuracy drop (over 50% in Fig. 1 corresponding to 0 PGD attack steps) and robust accuracy is only enhanced by a small margin (around 18% against PGD attacks). The negative results in Fig. 1 are not quite surprising since a data-agnostic input prompt  $\delta$  has limited learning capacity to enable adversarial defense. Thus, it is non-trivial to tackle the problem of AVP.

## 4. CLASS-WISE ADVERSARIAL VISUAL PROMPT

No free lunch for class-wise visual prompts. A direct extension of (U-AVP) is to introduce multiple adversarial visual prompts, each of which corresponds to one class in the training set  $\mathcal{D}_{tr}$ . If we split  $\mathcal{D}_{tr}$  into class-wise training sets  $\{\mathcal{D}_{tr}^{(i)}\}_{i=1}^{N}$ 

(for N classes) and introduce class-wise visual prompts  $\{\delta^{(i)}\}\)$ , then the direct C-AVP extension from (U-AVP) becomes

$$\underset{\{\boldsymbol{\delta}^{(i)} \in \mathcal{C}\}_{i \in [N]}}{\operatorname{minimize}} \quad \frac{1}{N} \sum_{i=1}^{N} \left\{ \lambda \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_{\mathrm{tr}}^{(i)}} [\ell(\mathbf{x} + \boldsymbol{\delta}^{(i)}; y, \boldsymbol{\theta})] + \\ \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_{\mathrm{tr}}^{(i)}} [\ell_{\mathrm{adv}}(\mathbf{x} + \boldsymbol{\delta}^{(i)}; y, \boldsymbol{\theta})] \right\}$$
(C-AVP-v0)

where [N] denotes the set of class labels  $\{1, 2, ..., N\}$ . It is worth noting that C-AVP-v0 is *decomposed* over class labels. Although the class-wise separability facilitates numerical optimization, it introduces challenges (C1)-(C2) when applying class-wise visual prompts for adversarial defense.

• (C1) *Test-time prompt selection*: After acquiring the visual prompts  $\{\delta^{(i)}\}$  from (C-AVP-v0), it remains unclear how a class-wise prompt should be selected for application to a test-time example  $\mathbf{x}_{\text{test}}$ . An intuitive way is to use the inference pipeline of  $\boldsymbol{\theta}$  by aligning its top-1 prediction with the prompt selection. That is, the selected prompt  $\boldsymbol{\delta}$  and the predicted class  $i^*$  are determined by

$$\boldsymbol{\delta} = \boldsymbol{\delta}^*, \ i^* = \operatorname*{arg\,max}_{i \in [N]} f_i(\mathbf{x}_{test} + \boldsymbol{\delta}^{(i)}; \boldsymbol{\theta}), \tag{3}$$

where  $f_i(\mathbf{x}; \boldsymbol{\theta})$  denotes the *i*th-class prediction confidence. However, the seemingly correct rule (3) leads to a large prompt selection error (thus poor prediction accuracy) due to (**C2**).

• (C2) Backdoor effect of class mis-matched prompts: Given  $\delta^{(i)}$  from (C-AVP-v0), if the test-time example  $\mathbf{x}_{\text{test}}$  is drawn from class *i*, the visual prompt  $\delta^{(i)}$  then helps prediction. However, if  $\mathbf{x}_{\text{test}}$  is *not* originated from class *i*, then  $\delta^{(i)}$ could serve as a backdoor attack trigger [46] with the targeted backdoor label *i* for the 'prompted input'  $\mathbf{x}_{\text{test}} + \delta^{(i)}$ . Since the backdoor attack is also input-agnostic, the class-discriminative ability of  $\mathbf{x}_{\text{test}} + \delta^{(i)}$  enabled by  $\delta^{(i)}$  could result in incorrect prediction towards the target class *i* for  $\mathbf{x}_{\text{test}}$ .

Joint prompts optimization for C-AVP. The failure of C-AVP-v0 inspires us to rethink the value of class-wise separability. As illustrated in challenges (C1)-(C2), the compatibility with the test-time prompt selection rule and the interrelationship between class-wise visual prompts should be taken into account. To this end, we develop a series of new AVP principles below. Fig. 2 provides a schematic overview of C-AVP and its comparison with U-AVP and the predictor without VP.



Fig. 2: Overview of C-AVP over two classes (red and green) vs. U-AVP and the prompt-free learning pipeline.

First, to bake the prompt selection rule (3) into C-AVP, we enforce the correct prompt selection, *i.e.*, under the condition that  $f_y(\mathbf{x} + \boldsymbol{\delta}^{(y)}; \boldsymbol{\theta}) > \max_{k:k \neq y} f_k(\mathbf{x} + \boldsymbol{\delta}^{(k)}; \boldsymbol{\theta})$  for  $(\mathbf{x}, y) \in \mathcal{D}^{(y)}$ . The above can be cast as a CW-type loss [47]:

$$\ell_{\text{C-AVP},1}(\{\boldsymbol{\delta}^{(i)}\}; \mathcal{D}_{\text{tr}}, \boldsymbol{\theta}) = \\ \mathbb{E}_{(\mathbf{x},y)\in\mathcal{D}_{\text{tr}}} \max\{\max_{k\neq y} f_k(\mathbf{x} + \boldsymbol{\delta}^{(k)}; \boldsymbol{\theta}) - f_y(\mathbf{x} + \boldsymbol{\delta}^{(y)}; \boldsymbol{\theta}), -\tau\},$$
<sup>(4)</sup>

where  $\tau > 0$  is a confidence threshold. The rationale behind (4) is that given a data sample  $(\mathbf{x}, y)$ , the minimum value of  $\ell_{C-AVP,1}$  is achieved at  $-\tau$ , indicating the desired condition with the confidence level  $\tau$ . Compared with (C-AVP-v0), another key characteristic of  $\ell_{C-AVP,1}$  is its non-splitting over class-wise prompts { $\delta^{(i)}$ }, which benefits the joint optimization of these prompts.

Second, to mitigate the backdoor effect of mis-matched prompts, we propose additional two losses, noted by  $\ell_{C-AVP,2}$  and  $\ell_{C-AVP,3}$ , to penalize the data-prompt mismatches. Specifically,  $\ell_{C-AVP,2}$  penalizes the backdoor-alike targeted prediction accuracy of a class-wise visual prompt when applied to mismatched training data. For the prompt  $\delta^{(i)}$ , this leads to

$$\ell_{\text{C-AVP},2}(\{\boldsymbol{\delta}^{(i)}\}; \mathcal{D}_{\text{tr}}, \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{E}_{(\mathbf{x},y)\in\mathcal{D}_{\text{tr}}^{(-i)}} \max\{f_i(\mathbf{x} + \boldsymbol{\delta}^{(i)}; \boldsymbol{\theta}) - f_y(\mathbf{x} + \boldsymbol{\delta}^{(i)}; \boldsymbol{\theta}), -\tau\},$$
(5)

where  $\mathcal{D}_{tr}^{(-i)}$  denotes the training data set by excluding  $\mathcal{D}_{tr}^{(i)}$ . The class *i*-associated prompt  $\delta^{(i)}$  should *not* behave as a backdoor trigger to non-*i* classes' data. Likewise, if the prompt is applied to the correct data class, then the prediction confidence should surpass that of a mis-matched case. This leads to

$$\ell_{\text{C-AVP},3}(\{\boldsymbol{\delta}^{(i)}\}; \mathcal{D}_{\text{tr}}, \boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x},y)\in\mathcal{D}_{\text{tr}}} \max\{\max_{k\neq y} f_y(\mathbf{x} + \boldsymbol{\delta}^{(k)}; \boldsymbol{\theta}) - f_y(\mathbf{x} + \boldsymbol{\delta}^{(y)}; \boldsymbol{\theta}), -\tau\}.$$
<sup>(6)</sup>

Let  $\ell_{\text{C-AVP},0}(\{\delta^{(i)}\}; \mathcal{D}_{\text{tr}}, \theta)$  denote the objective function of (C-AVP-v0). Integrated with  $\ell_{\text{C-AVP},q}(\{\delta^{(i)}\}; \mathcal{D}_{\text{tr}}, \theta)$ for  $q \in \{1, 2, 3\}$ , the desired class-wise AVP design is cast as

where  $\gamma > 0$  is a parameter for class-wise prompting penalties.

#### 5. EXPERIMENTS

**Experiment setup.** We conduct experiments on CIFAR-10 with a pretrained ResNet18 of testing accuracy of 94.92% on standard test dataset. We use PGD-10 (*i.e.*, PGD attack with 10 steps [1]) to generate adversarial examples with  $\epsilon = 8/255$  during visual prompts training, and with a cosine learning rate scheduler starting at 0.1. Throughout experiments, we choose  $\lambda = 1$  in (U-AVP), and  $\tau = 0.1$  and  $\gamma = 3$  in (C-AVP). The width of visual prompt is set to 8 (see Fig. 3 for the visualization).

**C-AVP outperforms conventional VP.** Tab. 1 demonstrates the effectiveness of proposed C-AVP approach vs. U-AVP (the



Fig. 3: C-AVP visualization. One image is chosen from each CIFAR-10 class with the corresponding C-AVP.

**Table 1:** VP performance comparison in terms of standard (std) accuracy (acc) and robust accuracy against PGD attacks with  $\epsilon = 8/255$  and multiple PGD steps on (CIFAR-10, ResNet18).

Evaluation metrics (%)	Std acc	Robus 10	st acc vs 20	PGD w/ 50	step # 100
Pre-trained	94.92	0	0	0	0
Vanilla VP	94.48	0	0	0	0
U-AVP	27.75	16.9	16.81	16.81	16.7
C-AVP-v0	19.69	13.91	13.63	13.6	13.58
C-AVP (ours)	57.57	34.75	34.62	34.51	33.63

direct extension of VP to adversarial defense) and the C-AVPv0 method in the task of robustify a normally-trained ResNet18 on CIFAR-10. For comparison, we also report the standard accuracy of the pre-trained model and the vanilla VP solution given by (1). As we can see, C-AVP outperforms U-AVP and C-AVP-v0 in both standard accuracy and robust accuracy. We also observe that compared to the pretrained model and the vanilla VP, the robustness-induced VP variants bring in an evident standard accuracy drop as the cost of robustness.

**Prompting regularization effect in (C-AVP).** Tab. 2 shows different settings of prompting regularizations used in C-AVP, where 'S*i*' represents a certain loss configuration. As we can see, the use of  $\ell_{C-AVP,2}$  contributes most to the performance of learned visual prompts (see S3). This is not surprising, since we design  $\ell_{C-AVP,2}$  for mitigating the backdoor effect of class-wise prompts, which is the main source of prompting selection error. We also note that  $\ell_{C-AVP,1}$  is the second most important regularization. This is because such a regularization is accompanied with the prompt selection rule (3). Tab. 2 also indicates that the combination of  $\ell_{C-AVP,1}$  and  $\ell_{C-AVP,2}$  is a possible computationally lighter alternative to (C-AVP).

**Class-wise prediction error analysis.** Fig. 4 shows a comparison of the classification confusion matrix. Each row corresponds to testing samples from one class, and each column corresponds to the prompt ('P') selection across 10 image classes. As we can see, our proposal outperforms C-AVP-v0 since the former's higher main diagonal entries indicate less prompt selection error than the latter.

**Comparisons with other test-time defenses.** In Tab. 3, we compare our proposed C-AVP with three test-time defense methods selected from Croce *et. al.* [14]. Note that all methods are applied to robustifying a fixed, standardly pre-trained ResNet18. Following Croce *et. al.* [14], we divide the considered defenses into different categories, relying on their defense principles (*i.e.*, IP or MA) and needed test-time operations (*i.e.*, IA, AN, and R). As we can see, our method C-AVP falls into the IP category but requires no involved test-time operations. This leads to the least inference overhead. Although there



**Fig. 4**: The predictions of C-AVP-v0 vs. C-AVP on (CIFAR10, ResNet18). **Table 2**: Sensitivity analysis of prompting regularization in C-AVP on (CIFAR-10, ResNet18).

Setting	$\ell_{\rm C-AVP,1}$	$\ell_{\rm C-AVP,2}$	$\ell_{\rm C-AVP,3}$	Std Acc (%)	PGD-10 Acc (%)
<b>S</b> 1	×	×	×	19.69	13.91
<b>S</b> 2	<ul> <li>✓</li> </ul>	×	×	22.72	13.01
<b>S</b> 3	×	<ul> <li>✓</li> </ul>	×	40.01	25.40
<b>S</b> 4	×	×	<ul> <li></li> </ul>	17.44	11.78
<b>S</b> 5	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>	×	57.03	32.39
<b>S</b> 6	<ul> <li>✓</li> </ul>	×	<ul> <li></li> </ul>	26.02	15.80
S7	<ul> <li>✓</li> </ul>	<ul> <li>Image: A second s</li></ul>	<ul> <li>✓</li> </ul>	57.57	34.75

exists a performance gap with the test-time defense baselines, we hope that our work could pave a way to study the pros and cons of visual prompting in adversarial robustness.

**Table 3:** Comparison of C-AVP with other SOTA test-time defenses. Per the benchmark in Croce *et. al.* [14], the involved test-time operations in these defenses include: IP (input purification), MA (model adaption), IA (iterative algorithm), AN (auxiliary network), and R (randomness). And inference time (IT), standard accuracy (SA), and robust accuracy (RA) against PGD-10 are used as performance metrics.

Method	IP	MA	IA	AN	R	IT	SA (%)	RA (%)
[43]	~	×	~	×	×	518 ×	85.9%	0.4%
[15]	V .	×	~	<b>~</b>	1	176 ×	91.1%	40.3%
[44]	×	~	~	× .	×	59 ×	56.1%	50.6%
C-AVP	<b>v</b>	×	×	×	×	1.4 ×	57.6%	34.8%

### 6. CONCLUSION

In this work, we develop a novel VP method, *i.e.*, C-AVP, to improve adversarial robustness of a fixed model at test time. Compared to existing VP methods, this is the first work to peer into how VP could be in adversarial defense. We show the direct integration of VP into robust learning is *not* an effective adversarial defense at test time for a fixed model. To address this problem, we propose C-AVP to create ensemble visual prompts and jointly optimize their interrelations for robustness enhancement. We empirically show that our proposal significantly reduces the inference overhead compared to classical adversarial defenses which typically call for computationally-intensive test-time defense operations.

#### 7. REFERENCES

- [1] Aleksander Madry et al., "Towards deep learning models resistant to adversarial attacks," *arXiv:1706.06083*, 2017.
- [2] Hongyang Zhang, Yaodong Yu, et al., "Theoretically principled trade-off between robustness and accuracy," *ICML*, 2019.
- [3] Ali Shafahi, Mahyar Najibi, et al., "Adversarial training for free!," in *NeurIPS*, 2019.
- [4] Dinghuai Zhang, Tianyuan Zhang, et al., "You only propagate once: Accelerating adversarial training via maximal principle," arXiv:1905.00877, 2019.
- [5] Yair Carmon, Aditi Raghunathan, et al., "Unlabeled data improves adversarial robustness," *NeurIPS*, 2019.
- [6] Eric Wong and J Zico Kolter, "Provable defenses against adversarial examples via the convex outer adversarial polytope," *arXiv*:1711.00851, 2017.
- [7] Aditi Raghunathan, Jacob Steinhardt, et al., "Certified defenses against adversarial examples," *arXiv:1801.09344*, 2018.
- [8] Cihang Xie, Yuxin Wu, et al., "Feature denoising for improving adversarial robustness," in *CVPR*, 2019.
- [9] Tianlong Chen, Sijia Liu, et al., "Adversarial robustness: From self-supervised pre-training to fine-tuning," in CVPR, 2020.
- [10] Lijie Fan, Sijia Liu, et al., "When does contrastive learning preserve adversarial robustness from pretraining to finetuning?," *NeurIPS*, 2021.
- [11] Jinghan Jia et al., "Clawsat: Towards both robust and accurate code models," *arXiv:2211.11711*, 2022.
- [12] Anish Athalye, Nicholas Carlini, et al., "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," *arXiv*:1802.00420, 2018.
- [13] Francesco Croce and Matthias Hein, "Reliable evaluation of adversarial robustness with an ensemble of diverse parameterfree attacks," in *ICML*. PMLR, 2020.
- [14] Francesco Croce et al., "Evaluating the adversarial robustness of adaptive test-time defenses," *arXiv:2202.13711*, 2022.
- [15] Jongmin Yoon, Sung Ju Hwang, et al., "Adversarial purification with score-based generative models," in *ICML*. PMLR, 2021.
- [16] Chengzhi Mao, Mia Chiquier, et al., "Adversarial attacks are reversible with natural supervision," in *ICCV*, 2021.
- [17] Motasem Alfarra, Juan C Pérez, et al., "Combating adversaries with anti-adversaries," in *AAAI*, 2022.
- [18] Hadi Salman, Mingjie Sun, et al., "Denoised smoothing: A provable defense for pretrained classifiers," *NeurIPS*, 2020.
- [19] Yifan Gong et al., "Reverse engineering of imperceptible adversarial image perturbations," arXiv:2203.14145, 2022.
- [20] Qiyu Kang et al., "Stable neural ode with lyapunov-stable equilibrium points for defending against adversarial attacks," *NeurIPS*, 2021.
- [21] Hyojin Bahng et al., "Visual prompting: Modifying pixel space to adapt pre-trained models," *arXiv:2203.17274*, 2022.
- [22] Pin-Yu Chen, "Model reprogramming: Resource-efficient crossdomain machine learning," arXiv:2202.10629, 2022.

- [23] Gamaleldin F Elsayed, Ian Goodfellow, et al., "Adversarial reprogramming of neural networks," arXiv:1806.11146, 2018.
- [24] Yun-Yun Tsai et al., "Transfer learning without knowing: Reprogramming black-box machine learning models with scarce data and limited resources," *arXiv:2007.08714*, 2020.
- [25] Guanhua Zhang, Yihua Zhang, et al., "Fairness reprogramming," arXiv:2209.10222, 2022.
- [26] Hadi Salman, Andrew Ilyas, et al., "Unadversarial examples: Designing objects for robust vision," *NeurIPS*, 2021.
- [27] Tom Brown, Benjamin Mann, et al., "Language models are few-shot learners," *NeurIPS*, 2020.
- [28] Xiang Lisa Li and Percy Liang, "Prefix-tuning: Optimizing continuous prompts for generation," *arXiv:2101.00190*, 2021.
- [29] Alec Radford et al., "Learning transferable visual models from natural language supervision," in *ICML*. PMLR, 2021.
- [30] Yimeng Zhang et al., "Text-visual prompting for efficient 2d temporal video grounding," *arXiv:2303.04995*, 2023.
- [31] Paarth Neekhara, Shehzeen Hussain, et al., "Cross-modal adversarial reprogramming," in *WACV*, 2022.
- [32] Chao-Han Huck Yang, Yun-Yun Tsai, et al., "Voice2series: Reprogramming acoustic models for time series classification," in *ICML*. PMLR, 2021.
- [33] Yang Zheng, Xiaoyi Feng, et al., "Why adversarial reprogramming works, when it fails, and how to tell the difference," *arXiv:2108.11673*, 2021.
- [34] Kathrin Grosse, Praveen Manoharan, et al., "On the (statistical) detection of adversarial examples," arXiv:1702.06280, 2017.
- [35] Puyudi Yang et al., "Ml-loo: Detecting adversarial examples with feature attribution," arXiv:1906.03499, 2019.
- [36] Jan Hendrik Metzen, Tim Genewein, et al., "On detecting adversarial perturbations," arXiv:1702.04267, 2017.
- [37] Dongyu Meng and Hao Chen, "Magnet: a two-pronged defense against adversarial examples," arXiv:1705.09064, 2017.
- [38] Bartosz Wójcik, Paweł Morawiecki, et al., "Adversarial examples detection and analysis with layer-wise autoencoders," *arXiv*:2006.10013, 2020.
- [39] Akhilan Boopathy et al., "Proper network interpretability helps adversarial robustness in classification," in *ICML*, 2020.
- [40] Shaokai Ye, Kaidi Xu, et al., "Adversarial robustness vs model compression, or both?," arXiv e-prints, 2019.
- [41] Jeet Mohapatra, Ching-Yun Ko, et al., "Rethinking randomized smoothing for adversarial robustness," arXiv:2003.01249, 2020.
- [42] Ren Wang, Kaidi Xu, et al., "On fast adversarial robustness adaptation in model-agnostic meta-learning," in *ICLR*, 2021.
- [43] Changhao Shi, Chester Holtz, et al., "Online adversarial purification based on self-supervision," arXiv:2101.09387, 2021.
- [44] Zhuotong Chen, Qianxiao Li, et al., "Towards robust neural networks via close-loop control," arXiv:2102.01862, 2021.
- [45] Yimeng Zhang, Yuguang Yao, et al., "How to robustify black-box ml models? a zeroth-order optimization perspective," arXiv:2203.14195, 2022.

- [46] Tianyu Gu, Brendan Dolan-Gavitt, et al., "Badnets: Identifying vulnerabilities in the machine learning model supply chain," *arXiv:1708.06733*, 2017.
- [47] Nicholas Carlini et al., "Towards evaluating the robustness of neural networks," in *IEEE Symposium on S&P*, 2017.