

---

# Visual Prompting for Adversarial Robustness

---

**Aochuan Chen\***  
Michigan State University  
chenaoch@msu.edu

**Peter Lorenz\***  
Fraunhofer ITWM<sup>†</sup>  
peter.lorenz@itwm.fhg.de

**Yuguang Yao**  
Michigan State University  
yaoyugua@msu.edu

**Pin-Yu Chen**  
IBM Research  
pin-yu.chen@ibm.com

**Sijia Liu**  
Michigan State University  
liusiji5@msu.edu

## Abstract

In this work, we leverage visual prompting (VP) to improve adversarial robustness of a fixed, pre-trained model at testing 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 testing 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  $2.1\times$  standard accuracy gain and  $2\times$  robust accuracy gain. Compared to classical test-time defenses, C-AVP also yields a  $42\times$  inference time speedup. Code is available at github.

## 1 Introduction

Current machine learning (ML) models, *e.g.*, vision models in particular, can easily be manipulated (by an adversary) to output drastically different classifications and can be done so in a controlled and directed way. This process is known as *adversarial attack* and is considered as one of the major hurdles in using ML models in high-stakes applications [Goodfellow et al., 2014, Carlini and Wagner, 2017]. Thereby, model robustification against adversarial attacks is now a major focus of research. Yet, a large volume of existing works focused on advancing training recipes and/or model architectures to gain robustness. For example, adversarial training (AT) [Madry et al., 2017], one of the most effective defense methods, adopted min-max optimization to minimize the worst-case (maximum) training loss induced by adversarial attacks. Extended from AT, various empirical and certified defense methods were proposed in various learning paradigms, ranging from supervised learning to semi-supervised learning, and further to unsupervised learning [Zhang et al., 2019b, Shafahi et al., 2019, Zhang et al., 2019a, Carmon et al., 2019, Wong and Kolter, 2017, Raghunathan et al., 2018, Xie et al., 2019, Chen et al., 2020, Fan et al., 2021].

Although the design for robust training has made tremendous success in improving model robustness [Athalye et al., 2018, Croce and Hein, 2020], it typically takes an intensive computation cost with

---

\*Equal Contribution. <sup>†</sup> and Fraunhofer Center of Machine Learning

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., 2022]. Test-time defense alters either a test-time input example or a small portion of the pre-trained model for adversarial defense. Examples include input (anti-adversarial) purification [Yoon et al., 2021, Mao et al., 2021, Alfarrá et al., 2022] and model refinement by augmenting the pre-trained model with auxiliary components [Salman et al., 2020, Gong et al., 2022, Kang et al., 2021]. However, these defense techniques inevitably raise the inference time and hamper the test-time efficiency [Croce et al., 2022]. Inspired by that, our work will advance the test-time defense technology by leveraging the idea of *visual prompting* (VP) [Bahng et al., 2022], also known as model reprogramming [Chen, 2022, Elsayed et al., 2018, Tsai et al., 2020, Zhang et al., 2022].

Generally speaking, VP [Bahng et al., 2022] creates a *universal* (i.e., *data-agnostic*) 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 [Chen, 2022, Elsayed et al., 2018, Tsai et al., 2020, Zhang et al., 2022] or unadversarial example [Salman et al., 2021], which optimizes the 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 [Tsai et al., 2020], out-of-distribution generalization [Salman et al., 2021], and fairness [Zhang et al., 2022]. 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 will significantly reduce 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** below.

- ❶ We formulate and investigate the problem of AVP for the first time. We empirically show that the conventional data-agnostic VP design is incapable of gaining adversarial robustness.
- ❷ 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 adversarial robustness.
- ❸ We provide insightful experiments to demonstrate the pros and cons of VP in adversarial defense.

## 1.1 Related work

**Visual prompting.** Originated from the idea of in-context learning or prompting in natural language processing (NLP) [Brown et al., 2020, Li and Liang, 2021, Radford et al., 2021], VP was first proposed in [Bahng et al., 2022] for vision models. Before formalizing VP in [Bahng et al., 2022], 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* [Elsayed et al., 2018, Chen, 2022, Tsai et al., 2020, Neekhara et al., 2022, Yang et al., 2021, Zheng et al., 2021], 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 [Salman et al., 2021], 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 [Grosse et al., 2017, Yang et al., 2019, Metzen et al., 2017, Meng and Chen, 2017, Wójcik et al., 2020, Gong et al., 2022] and robust training [Wong and Kolter, 2017, Zhang et al., 2019b, Salman et al., 2020, Chen et al., 2020, Boopathy et al., 2020, Fan et al., 2021], is a current research focus. In particular, adversarial training (AT) [Madry et al., 2017] is the most widely-used defense strategy and has inspired many recent advances in adversarial defense [Athalye et al., 2018, Ye et al., 2019, Croce and Hein, 2020, Mohapatra et al., 2020, Kang et al., 2021, Wang et al., 2021]. 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 [Croce et al., 2022], which aims to robustify a given model via lightweight unadversarial input perturbations (a.k.a input purification) [Shi et al., 2021, Yoon et al., 2021] or minor modifications to the fixed model [Chen et al., 2021]. In different kinds of test-time defenses, the most relevant work to ours is anti-adversarial perturbation [Alfarrá et al., 2022].

## 2 Problem Statement

In this section, we will begin by providing a brief background on VP, and then introduce the problem of our interest—*adversarial visual prompting (AVP)*—which aims at generating visual prompts to improve adversarial robustness of a pre-trained, fixed model. Through a warm-up example, we will empirically show that the conventional design of VP is difficult to apply to the paradigm of AVP.

**Visual prompting.** We describe the problem setup of VP following [Bahng et al., 2022, Elsayed et al., 2018, Tsai et al., 2020, Zhang et al., 2022]. Specifically, let  $\mathcal{D}_{\text{tr}}$  denote a training set for supervised learning, where  $(\mathbf{x}, y) \in \mathcal{D}_{\text{tr}}$  signifies a training sample with feature  $\mathbf{x}$  and label  $y$ . And let  $\delta$  be a visual prompt to be designed. The prompted input is then given by  $\mathbf{x} + \delta$  with respect to (w.r.t.)  $\mathbf{x}$ . Different from the problem of adversarial attack generation that optimizes  $\delta$  for erroneous prediction, VP drives  $\delta$  to minimize the performance loss  $\ell$  of a pre-trained model  $\theta$ . This leads to

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

where  $\ell$  denotes a certain performance loss (e.g., prediction error [Bahng et al., 2022]) given the prior knowledge of training data  $(\mathbf{x}, y)$  and base model  $\theta$ , and  $\mathcal{C}$  is a perturbation constraint. Following [Elsayed et al., 2018, Tsai et al., 2020, Bahng et al., 2022],  $\mathcal{C}$  restricts  $\delta$  to be located in an image’s boundary region and requests the perturbation magnitude within a normalized input space, i.e.,  $\mathbf{x} + \delta \in [0, 1]$  for any  $\mathbf{x}$ . Projected gradient descent (PGD) [Madry et al., 2017, Salman et al., 2021] can then be applied to solving problem (1). At inference time, the designed  $\delta$  will be integrated into test-time examples to improve the prediction ability of  $\theta$ .

**Adversarial visual prompting.** Inspired by the usefulness of VP to improve model generalization [Tsai et al., 2020, Bahng et al., 2022], we ask:

**(AVP problem)** Can VP (1) be extended to robustify  $\theta$  against adversarial attacks?

At the first glance, the AVP problem seems trivial only if we specify the performance loss  $\ell$  as the adversarial training (AT) loss [Madry et al., 2017, Zhang et al., 2019b]:

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

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

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

where  $\lambda > 0$  is a regularization parameter to strike a balance between generalization and adversarial robustness [Zhang et al., 2019b].

The problem (U-AVP) can be effectively solved using a standard min-max 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 testing 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.

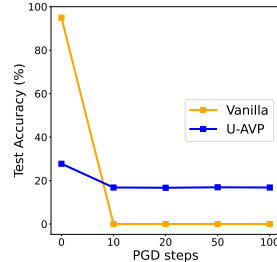


Fig. 1: Example of designing U-AVP for adversarial defense on (CIFAR-10, ResNet18), measured by robust accuracy against PGD attacks [Madry et al., 2017] of different steps. The robust accuracy of 0 steps is the standard accuracy.

### 3 Class-wise Adversarial Visual Prompting

In this section, we will develop a new VP approach, termed Class-wise AVP (**C-AVP**), which improves (U-AVP) in adversarial robustness. Different from U-AVP, C-AVP expands the designing space of VP by associating each image class with an adversarial visual prompt and taking the couplings of these class-wise visual prompts into account for robustness enhancement.

**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  $\mathcal{D}_{\text{tr}}$ . If we split  $\mathcal{D}_{\text{tr}}$  into class-wise training sets  $\{\mathcal{D}_{\text{tr}}^{(i)}\}_{i=1}^N$  (for  $N$  classes) and introduce class-wise visual prompts  $\{\delta^{(i)}\}$ , then the direct C-AVP extension from (U-AVP) becomes

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

where  $[N]$  denotes the set of class labels  $\{1, 2, \dots, N\}$ . It is worth noting that C-AVP-v0 is *decomposed* over class labels. That is, solving the above problem is equivalent to solving a sequence of sub-problems: For each class  $i$ ,

$$\underset{\delta^{(i)} \in \mathcal{C}}{\text{minimize}} \quad \lambda \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_{\text{tr}}^{(i)}} [\ell(\mathbf{x} + \delta^{(i)}; y, \theta)] + \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_{\text{tr}}^{(i)}} [\ell_{\text{adv}}(\mathbf{x} + \delta^{(i)}; y, \theta)] \quad (3)$$

Although the class-wise separability facilitates numerical optimization, it introduces two challenges **(C1)**-**(C2)** when applying class-wise visual prompts to defend adversarial attacks.

- **(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  $\theta$  by aligning its top-1 prediction with the prompt selection. That is, the selected prompt  $\delta$  and the predicted class  $i^*$  are determined by

$$\delta = \delta^*, \quad i^* = \arg \max_{i \in [N]} f_i(\mathbf{x}_{\text{test}} + \delta^{(i)}; \theta), \quad (4)$$

where  $f_i(\mathbf{x}; \theta)$  denotes the  $i$ th-class prediction confidence of using  $\theta$  at  $\mathbf{x}$ . However, the seemingly correct rule (4) 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 (3), 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 [Gu et al., 2017] 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  $\mathbf{x}_{\text{test}}$ . Our empirical experiments justified the above: Nearly all testing samples will be (mis)classified as the prompt’s class regardless of their true labels.

#### Joint prompts optimization for C-AVP.

The failure of C-AVP-v0 inspires us to re-think the value of class-wise separability in (3). 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 original predictor without VP.

First, to bake the prompt selection rule (4) into C-AVP, we enforce the prompt design along the correct selection path, *i.e.*, under the condition that  $f_y(\mathbf{x} + \delta^{(y)}; \theta) > \max_{k: k \neq y} f_k(\mathbf{x} + \delta^{(k)}; \theta)$  for  $(\mathbf{x}, y) \in \mathcal{D}^{(y)}$ . The above can be cast as a CW-type loss [Carlini and Wagner, 2017]:

$$\ell_{\text{C-AVP},1}(\{\delta^{(i)}\}; \mathcal{D}_{\text{tr}}, \theta) = \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_{\text{tr}}} \max \left\{ \max_{k \neq y} f_k(\mathbf{x} + \delta^{(k)}; \theta) - f_y(\mathbf{x} + \delta^{(y)}; \theta), -\tau \right\}, \quad (5)$$

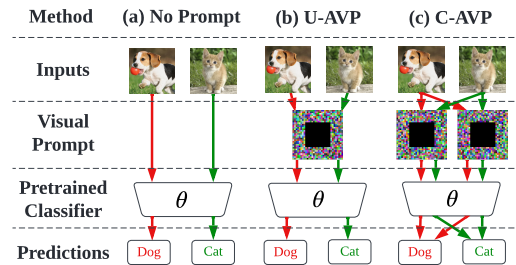


Fig. 2: Overview of C-AVP over two classes (red and green) vs. U-AVP and the prompt-free learning pipeline.

where  $\tau > 0$  is a confidence threshold. The rationale behind (5) 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 class 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 mis-matched training data. For the prompt  $\delta^{(i)}$ , this leads to

$$\ell_{C-AVP,2}(\{\delta^{(i)}\}; \mathcal{D}_{tr}, \theta) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_{tr}^{(-i)}} \max\{f_i(\mathbf{x} + \delta^{(i)}; \theta) - f_y(\mathbf{x} + \delta^{(i)}; \theta), -\tau\}, \quad (6)$$

where  $\mathcal{D}_{tr}^{(-i)}$  denotes the training data set by excluding  $\mathcal{D}_{tr}^{(i)}$ . The rationale behind (6) is that 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 of this matched case should surpass that of a mis-matched case. This leads to

$$\ell_{C-AVP,3}(\{\delta^{(i)}\}; \mathcal{D}_{tr}, \theta) = \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_{tr}} \max\{\max_{k \neq y} f_y(\mathbf{x} + \delta^{(k)}; \theta) - f_y(\mathbf{x} + \delta^{(y)}; \theta), -\tau\}. \quad (7)$$

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

$$\underset{\{\delta^{(i)}\}_{i \in [N]}}{\text{minimize}} \quad \ell_{C-AVP,0}(\{\delta^{(i)}\}; \mathcal{D}_{tr}, \theta) + \gamma \sum_{q=1}^3 \ell_{C-AVP,q}(\{\delta^{(i)}\}; \mathcal{D}_{tr}, \theta), \quad (\text{C-AVP})$$

where  $\gamma > 0$  is a regularization parameter to control our emphasis on the class-wise prompting penalties. It is worth mentioning that since  $\ell_{C-AVP,q}$  (for  $q > 0$ ) is a hinge loss with hard threshold  $\tau$ , its optimization could automatically stop if a prompting regulation is satisfied. To solve (C-AVP), we will use the min-max optimizer similar to the approach used for solving (C-AVP-v0).

## 4 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 [Madry et al., 2017]) 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). To evaluate test-time adversarial robustness, we generate PGD attacks of different steps under  $\epsilon = 8/255$ .

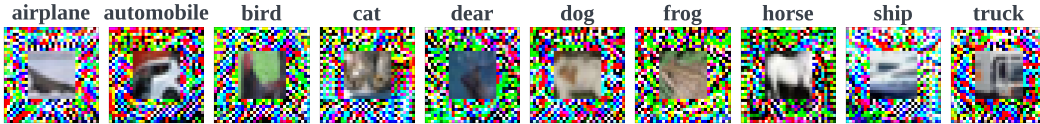


Fig. 3: C-AVP visualization. One image is chosen from each CIFAR-10 class with the corresponding C-AVP.

**C-AVP outperforms conventional VP.** Tab. 1 demonstrates the effectiveness of proposed C-AVP approach vs. U-AVP (the direction extension of VP to adversarial defense) and the C-AVP-v0 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 (evaluated using PGD attacks with different step sizes). We also observe that compared to the pre-trained model and the vanilla VP, the robustness-induced VP variants bring in an evident standard accuracy drop as the cost of robustness enhancement. This leaves a future research direction to optimize the accuracy-robustness trade-off of visual prompts.

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	Robust acc vs PGD w/ step #			
		10	20	50	100
Pre-trained	<b>94.92</b>	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
<b>C-AVP (ours)</b>	<b>57.57</b>	<b>34.75</b>	<b>34.62</b>	<b>34.51</b>	<b>33.63</b>

**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 improving 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, as evidenced by the comparable performance of S3 vs. S5. This is because such a regularization is accompanied with the prompt selection rule (4). If training cost is taken into consideration, 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).

Table 2: Sensitivity analysis of prompting regularizations in C-AVP on (CIFAR-10, ResNet18).

Setting	$\ell_{C-AVP,1}$	$\ell_{C-AVP,2}$	$\ell_{C-AVP,3}$	Std Acc (%)	PGD-10 Acc (%)
S1	✗	✗	✗	19.69	13.91
S2	✓	✗	✗	22.72	13.01
S3	✗	✓	✗	40.01	25.40
S4	✗	✗	✗	17.44	11.78
S5	✓	✓	✗	57.03	32.39
S6	✓	✗	✗	26.02	15.80
S7	✓	✓	✓	<b>57.57</b>	<b>34.75</b>

**Class-wise prediction error analysis.**

Fig. 4 shows a comparison of the classification confusion matrix over benign test dataset. Here the row index signifies the test data per class, and the column index refers to the selected prompt for prediction when using C-AVP-v0 or C-AVP.

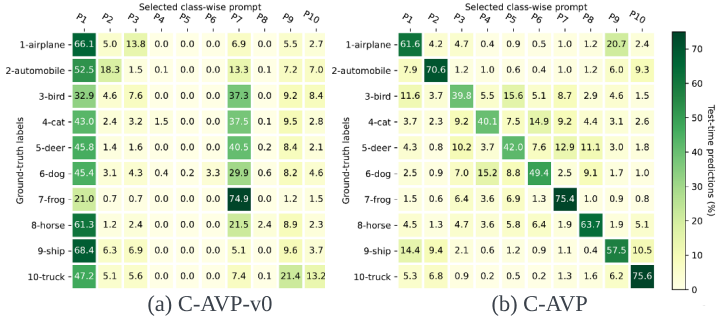


Fig. 4: The prediction error analysis of C-AVP vs. C-AVP-v0 on (CIFAR10, ResNet18). Each row corresponds to testing samples from one class, and each column corresponds to the prompt (‘P’) selection across 10 image classes.

**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., 2022]. Note that all methods are applied to robustifying a fixed, normally pre-trained ResNet18. Following [Croce et al., 2022], we divide the considered defenses into different categories, relying on their defense principles (*i.e.*, IP or MA) as well as their 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 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.

**5 Conclusion**

In this work, we develop a novel VP method, *i.e.*, C-AVP, to improve adversarial robustness of a fixed model at testing time. Compared to existing VP methods, this is the first work to peer into how VP could be in adversarial defense. We show that the direct integration of VP into robust learning does *not* offer an effective adversarial defense at testing time for the 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.

Table 3: Comparison of C-AVP with other SOTA test-time defenses. Per the benchmark in [Croce et al., 2022], 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 (%)
Shi et al. [2021]	✓	✗	✓	✗	✗	518 ×	85.9%	0.4%
Yoon et al. [2021]	✓	✗	✓	✓	✓	176 ×	91.1%	40.3%
Chen et al. [2021]	✗	✗	✓	✓	✓	59 ×	56.1%	50.6%
C-AVP	✓	✗	✗	✗	✗	<b>1.4 ×</b>	<b>57.6%</b>	<b>34.3%</b>

## References

- Motasem Alfarra, Juan C Pérez, Ali Thabet, Adel Bibi, Philip HS Torr, and Bernard Ghanem. Combating adversaries with anti-adversaries. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5992–6000, 2022.
- Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:1802.00420*, 2018.
- Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Visual prompting: Modifying pixel space to adapt pre-trained models. *arXiv preprint arXiv:2203.17274*, 2022.
- Akhilan Boopathy, Sijia Liu, Gaoyuan Zhang, Cynthia Liu, Pin-Yu Chen, Shiyu Chang, and Luca Daniel. Proper network interpretability helps adversarial robustness in classification. In *ICML*, 2020.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on S&P*, 2017.
- Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, John C Duchi, and Percy S Liang. Unlabeled data improves adversarial robustness. *Advances in Neural Information Processing Systems*, 32, 2019.
- Pin-Yu Chen. Model reprogramming: Resource-efficient cross-domain machine learning. *arXiv preprint arXiv:2202.10629*, 2022.
- Tianlong Chen, Sijia Liu, Shiyu Chang, Yu Cheng, Lisa Amini, and Zhangyang Wang. Adversarial robustness: From self-supervised pre-training to fine-tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 699–708, 2020.
- Zhuotong Chen, Qianxiao Li, and Zheng Zhang. Towards robust neural networks via close-loop control. *arXiv preprint arXiv:2102.01862*, 2021.
- Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International Conference on Machine Learning*, pages 2206–2216. PMLR, 2020.
- Francesco Croce, Sven Gowal, Thomas Brunner, Evan Shelhamer, Matthias Hein, and Taylan Cemgil. Evaluating the adversarial robustness of adaptive test-time defenses. *arXiv preprint arXiv:2202.13711*, 2022.
- Gamaleldin F Elsayed, Ian Goodfellow, and Jascha Sohl-Dickstein. Adversarial reprogramming of neural networks. *arXiv preprint arXiv:1806.11146*, 2018.
- Lijie Fan, Sijia Liu, Pin-Yu Chen, Gaoyuan Zhang, and Chuang Gan. When does contrastive learning preserve adversarial robustness from pretraining to finetuning? *Advances in Neural Information Processing Systems*, 34:21480–21492, 2021.
- Yifan Gong, Yuguang Yao, Yize Li, Yimeng Zhang, Xiaoming Liu, Xue Lin, and Sijia Liu. Reverse engineering of imperceptible adversarial image perturbations. *arXiv preprint arXiv:2203.14145*, 2022.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick McDaniel. On the (statistical) detection of adversarial examples. *arXiv preprint arXiv:1702.06280*, 2017.
- Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.



- Qiyu Kang, Yang Song, Qinxu Ding, and Wee Peng Tay. Stable neural ode with Lyapunov-stable equilibrium points for defending against adversarial attacks. *Advances in Neural Information Processing Systems*, 34:14925–14937, 2021.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Chengzhi Mao, Mia Chiquier, Hao Wang, Junfeng Yang, and Carl Vondrick. Adversarial attacks are reversible with natural supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 661–671, 2021.
- Dongyu Meng and Hao Chen. Magnet: a two-pronged defense against adversarial examples. *arXiv preprint arXiv:1705.09064*, 2017.
- Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. *arXiv preprint arXiv:1702.04267*, 2017.
- Jeet Mohapatra, Ching-Yun Ko, Sijia Liu, Pin-Yu Chen, Luca Daniel, et al. Rethinking randomized smoothing for adversarial robustness. *arXiv preprint arXiv:2003.01249*, 49, 2020.
- Paarth Neekhara, Shehzeen Hussain, Jinglong Du, Shlomo Dubnov, Farinaz Koushanfar, and Julian McAuley. Cross-modal adversarial reprogramming. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2427–2435, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. *arXiv preprint arXiv:1801.09344*, 2018.
- Hadi Salman, Mingjie Sun, Greg Yang, Ashish Kapoor, and J Zico Kolter. Denoised smoothing: A provable defense for pretrained classifiers. *NeurIPS*, 2020.
- Hadi Salman, Andrew Ilyas, Logan Engstrom, Sai Vemprala, Aleksander Madry, and Ashish Kapoor. Unadversarial examples: Designing objects for robust vision. *Advances in Neural Information Processing Systems*, 34:15270–15284, 2021.
- Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *Advances in Neural Information Processing Systems*, pages 3353–3364, 2019.
- Changhao Shi, Chester Holtz, and Gal Mishne. Online adversarial purification based on self-supervision. *arXiv preprint arXiv:2101.09387*, 2021.
- Yun-Yun Tsai, Pin-Yu Chen, and Tsung-Yi Ho. Transfer learning without knowing: Reprogramming black-box machine learning models with scarce data and limited resources. *arXiv preprint arXiv:2007.08714*, 2020.
- Ren Wang, Kaidi Xu, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Chuang Gan, and Meng Wang. On fast adversarial robustness adaptation in model-agnostic meta-learning. In *International Conference on Learning Representations*, 2021.
- Bartosz Wójcik, Paweł Morawiecki, Marek Śmieja, Tomasz Krzyżek, Przemysław Spurek, and Jacek Tabor. Adversarial examples detection and analysis with layer-wise autoencoders. *arXiv preprint arXiv:2006.10013*, 2020.
- Eric Wong and J Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. *arXiv preprint arXiv:1711.00851*, 2017.



- Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 501–509, 2019.
- Chao-Han Huck Yang, Yun-Yun Tsai, and Pin-Yu Chen. Voice2series: Reprogramming acoustic models for time series classification. In *International Conference on Machine Learning*, pages 11808–11819. PMLR, 2021.
- Puyudi Yang, Jianbo Chen, Cho-Jui Hsieh, Jane-Ling Wang, and Michael I Jordan. MI-loo: Detecting adversarial examples with feature attribution. *arXiv preprint arXiv:1906.03499*, 2019.
- Shaokai Ye, Kaidi Xu, Sijia Liu, Jan-Henrik Lambrechts, Huan Zhang, Aojun Zhou, Kaisheng Ma, Yanzhi Wang, and Xue Lin. Adversarial robustness vs model compression, or both? *arXiv e-prints*, pages arXiv–1903, 2019.
- Jongmin Yoon, Sung Ju Hwang, and Juho Lee. Adversarial purification with score-based generative models. In *International Conference on Machine Learning*, pages 12062–12072. PMLR, 2021.
- Dinghui Zhang, Tianyuan Zhang, Yiping Lu, Zhanxing Zhu, and Bin Dong. You only propagate once: Accelerating adversarial training via maximal principle. *arXiv preprint arXiv:1905.00877*, 2019a.
- Guanhua Zhang, Yihua Zhang, Yang Zhang, Wenqi Fan, Qing Li, Sijia Liu, and Shiyu Chang. Fairness reprogramming. *arXiv preprint arXiv:2209.10222*, 2022.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. *International Conference on Machine Learning*, 2019b.
- Yang Zheng, Xiaoyi Feng, Zhaoqiang Xia, Xiaoyue Jiang, Ambra Demontis, Maura Pintor, Battista Biggio, and Fabio Roli. Why adversarial reprogramming works, when it fails, and how to tell the difference. *arXiv preprint arXiv:2108.11673*, 2021.