# VPA: Fully Test-Time Visual Prompt Adaptation

Jiachen Sun*
University of Michigan
Ann Arbor, MI, USA
jiachens@umich.edu

Mark Ibrahim
Meta AI
New York, NY, USA
marksibrahim@meta.com

Melissa Hall
Meta AI
New York, NY, USA
melissahall@meta.com

Ivan Evtimov
Meta AI
Seattle, WA, USA
ivanevtimov@meta.com

Z. Morley Mao
University of Michigan
Ann Arbor, MI, USA
zmao@umich.edu

Cristian Canton Ferrer
Meta AI
Seattle, WA, USA
ccanton@meta.com

Caner Hazirbas
Meta AI
New York, NY, USA
hazirbas@meta.com

## ABSTRACT

Textual prompt tuning has demonstrated significant performance improvements in adapting natural language processing models to a variety of downstream tasks by treating hand-engineered prompts as trainable parameters. Inspired by the success of textual prompting, several studies have investigated the efficacy of visual prompt tuning. In this work, we present Visual Prompt Adaptation (VPA), the first framework that generalizes visual prompting with test-time adaptation. VPA introduces a small number of learnable tokens, enabling fully test-time and storage-efficient adaptation without necessitating source-domain information. We examine our VPA design under diverse adaptation settings, encompassing single-image, batched-image, and pseudo-label adaptation. We evaluate VPA on multiple tasks, including out-of-distribution (OOD) generalization, corruption robustness, and domain adaptation. Experimental results reveal that VPA effectively enhances OOD generalization by 3.3% across various models, surpassing previous test-time approaches. Furthermore, we show that VPA improves corruption robustness by 6.5% compared to strong baselines. Finally, we demonstrate that VPA also boosts domain adaptation performance by relatively 5.2%. Our VPA also exhibits marked effectiveness in improving the robustness of zero-shot recognition for vision-language models.

## CCS CONCEPTS

• **Computing methodologies → Machine learning approaches**.

---

*This project was mainly conducted during Jiachen Sun's internship at Meta AI.

---

## KEYWORDS

Test-Time Adaptation, Out-of-Distribution Generalization, Corruption Robustness, Domain Adaptation

## 1 INTRODUCTION

Visual recognition, a crucial component in multimedia systems, plays an essential role in various applications. As technology evolves and the demand for intelligent multimedia systems increases, the importance of effective and robust visual recognition techniques cannot be overstated. Although various deep neural networks achieve state-of-the-art (SOTA) performance on test sets drawn from the same distribution as the training set [9], these expertly-trained models may struggle to generalize when faced with distribution shifts, leading to substantial performance drops [23, 26]. These shifts encompass common corruption [24], adversarial attacks [5, 52, 53, 55], conceptual changes [23, 30], and even out-of-distribution (OOD)[1] variations [26], and can emerge in numerous real-world applications such as autonomous driving [51, 57, 64] and facial recognition systems [59], where accurate and robust performance is critical. Therefore, addressing the vulnerabilities to distribution shifts is essential for enhancing the robustness and generalization capabilities of deep learning models.

Numerous architectural improvements and training techniques have been proposed to address the challenges associated with achieving robustness against various domain variations [12, 21, 45]. For instance, recent advances in Vision Transformer (ViT) architectures [12] have demonstrated significant improvements on many out-of-distribution (OOD) generalization and corruption robustness benchmarks [23, 26]. Pretraining and fine-tuning strategies, such as CLIP [45] and WiSE [62], have further enhanced the generalization

---

[1]In this study, we refer out-of-distribution (OOD) to *covariate* shift but not *concept* shift of the test/validation dataset.

performance of ViT models [12]. In addition to these general methods, a plethora of specialized training recipes have been presented to address specific objectives, including OOD generalization [61, 63], corruption robustness [8, 25], domain adaptation [25]. However, it remains challenging to address these generalization problems solely during the training phase, as a single training recipe cannot encompass all underlying distributions. Test-time updates serve as valuable complements, focusing on tailored adaptation for specific test data [61]. Such a scheme is particularly important for multimedia systems, where content may come from unseen domains.

Humans typically begin with their existing knowledge and extrapolate from it when learning a new skill. Prompting is a similar paradigm that aids machine learning models in adapting to various contexts or even new tasks through specific textual input. This approach has gained popularity in the field of natural language processing (NLP) [41]. Recent studies have shown that prompt tuning can enhance model generalization across different domains, where predefined prompts evolve as trainable parameters [66, 67]. Test-time Prompt Tuning (TPT) is a pioneering technique that leverages textual prompting during testing to improve the generalization capability of vision-language models [45, 49]. Prompt tuning is efficient in adapting a pretrained model, as it does not modify the original model parameters. In addition to prompting in NLP, recent studies have explored *visual* prompting during the training phase [3, 29], yielding substantial improvements on numerous vision benchmarks. However, there is still a scarcity of research examining the application of visual prompting in online test-time adaptation, indicating an area ripe for further exploration.

**Our Contributions**. In this paper, we propose visual prompt adaptation (VPA) to bridge the gap between visual prompting and online test-time adaptation, drawing inspiration from the success of textual prompting in NLP. VPA is a simple yet effective framework that generalizes existing prompt designs and adaptation objectives. Given a pretrained model, we attach additive and prependitive prompts to the frozen model during the adaptation phase. VPA is highly storage-efficient. Rather than adapting all the parameters of the model, VPA requires only a small number of prompts to be stored. This efficiency greatly reduces the storage overhead while maintaining the ability to effectively adapt the model to new contexts and tasks, making VPA a practical and appealing solution for real-world applications. In contrast to the pixel-space prompts in [3] and the randomized initialization for embedding-space prompts in [29], we design a straightforward but intuitive paradigm using zero attention to initialize our prompts, ensuring that the original performance remains unaffected. We combine VPA with various adaptation settings, including batched-image adaptation (BIA), single-image adaptation (SIA), and pseudo-label adaptation (PLA). For BIA, we input a batch of images belonging to different classes into the model simultaneously and leverage self-entropy minimization as the adaptation objective. In SIA, we employ marginal entropy minimization as the objective by enriching a single image into a batch through various augmentations. Additionally, we employ confidence selection to identify images with high confidence, which allows for more effective adaptation. By focusing on these high-confidence images, the VPA framework can better leverage the information contained within them. Pseudo-labels have been shown to be effective in enhancing test-time adaptation performance. As such, we incorporate

a memory queue that stores pseudo-labels for historical data to assist incoming images during adaptation. Our VPA combines the strengths of VPA, BIA, SIA, and PLA to achieve both *fully* test-time and *storage-efficient* adaptation.

We conduct extensive evaluations of VPA across three critical axes in real-world vision systems: OOD generalization (§ 4.1), corruption robustness (§ 4.2), and domain adaptation (§ 4.3). In our OOD generalization experiments, we are the first to evaluate visual prompting on a variety of large-scale datasets. We primarily use the ViT architecture fine-tuned on the ImageNet training set as our base model. Our results indicate that VPA enhances the average accuracy of ImageNet-scale OOD generalization benchmarks [23, 26] by 3.3%. In contrast, existing state-of-the-art methods, such as TENT [61] and DDA [16], struggle to perform effectively under challenging OOD scenarios. Moreover, VPA achieves a similar 2.6% improvement as MEMO [63] under SIA without updating the parameters of the frozen model. Notably, we also demonstrate that VPA enhances corruption robustness and domain adaptation performance by relative margins of 6.5% and 5.2%, respectively, compared to strong baselines [6, 32]. Moreover, we have shown that VPA could effectively improve the robustness of zero-shot recognition for vision-language models [45]. It is important to note that the goal of our study is not solely to pursue state-of-the-art results but to highlight the potential applications of visual prompting in test-time adaptation. Our promising results will encourage future research to develop new adaptation schemes using visual prompting.
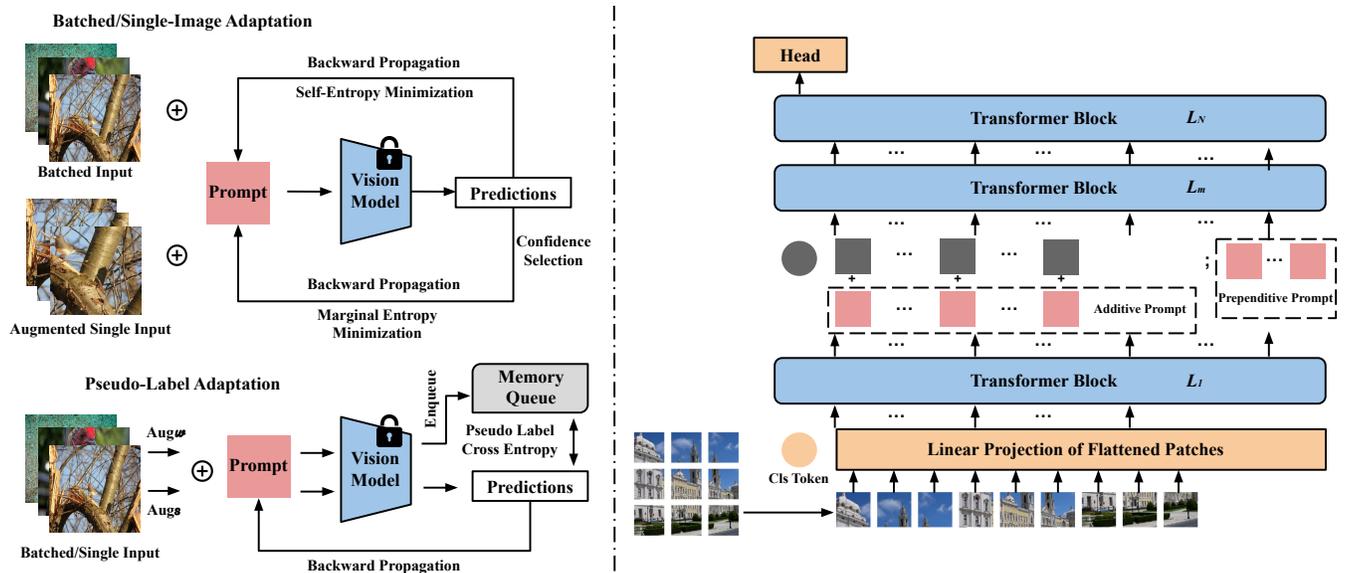
We summarize our main contributions as three-fold:
• We propose visual prompt adaptation (VPA), a *fully* test-time and *storage-efficient* adaptation framework that introduces both additive and prependitive adaptable tokens to improve the robustness of vision models.
• We conduct a rigorous taxonomy of VPA under different adaptation setups, including batched-image, single-image, and pseudo-label adaptation.
• We perform an extensive evaluation of VPA on various tasks, including out-of-distribution (OOD) generalization, corruption robustness, and domain adaptation. Our VPA consistently improves performance on these benchmarks by a significant margin.

## 2 RELATED WORK

In this section, we review topics related to our study, including prompting in foundation models, test-time adaptation, and out-of-distribution (OOD) robustness and domain adaptation.

**Prompting in Foundation Models**. Deep learning has made significant strides in natural language processing and computer vision tasks [36, 37]. In addition to architectural improvements [12, 22, 60], recent efforts have focused on fine-tuning foundational models with large-scale data to enable transfer learning across multiple downstream tasks [10, 45]. Prompting, which originated in language models, involves using human-engineered texts to improve the context-specific learning of a given task [33, 41]. Prompt tuning has been advanced for different goals in language models [19, 20, 38], while visual prompting has been explored for zero-shot recognition tasks [66, 67]. For instance, CoOp [67] and CoCoOp [66] both employ trainable prompts to improve zero-shot recognition performance, while TPT [49] leverages test-time adaptation of language

**Figure 1: Overview of Visual Prompt Adaptation Pipeline. Our VPA supports both batched- and single-image and pseudo-label adaptation settings as shown in the left figure. The visual prompt designs are illustrated in the right figure based on the ViT architecture.**

prompts to improve out-of-distribution robustness. Visual prompting has also been proposed to reprogram recognition tasks [14] and to enhance model performance on various downstream tasks [3]. Recent research has also introduced memory-efficient prompt tuning purely for vision models [29] to improve model generalization. In this work, we propose visual prompt adaptation that achieves *fully* test-time and *storage-efficient* adaptation to improve model generalization.

**Test-Time Adaptation**. Adapting machine learning models to different test domains has been applied to numerous tasks [35, 39, 40, 56, 58]. Among various adaptation techniques, we focus on test-time adaptation, which is particularly beneficial as it does not require label information from the test data [13, 39, 58]. *Source-free* adaptation allows models to adapt without any source-domain knowledge, adhering to real-world deployment constraints related to computation and privacy [13, 30]. *Fully* test-time adaptation is more rigorous, as it necessitates on-the-fly model updates without delaying inference.[42, 48] represent initial efforts towards achieving fully test-time adaptation, which involves updating or replacing the statistics of batch normalization (`BatchNorm`) layers[27] during inference. TENT accomplishes fully test-time adaptation by updating the model parameters in `BatchNorm` layers, using self-entropy minimization as its objective [61]. MEMO takes advantage of input augmentations to achieve single image adaptation, circumventing the batch-level adaptation requirement in TENT [63]. In contrast to optimizing model parameters, an alternative approach is to adapt input with minor modifications. DDA employs diffusion models to purify input data, although it is limited to specific types of corruption [16]. TPT introduces test-time adaptation via language prompting to enhance the OOD robustness of the CLIP model [49].

**Model Robustness against Distribution Shifts**. A trustworthy machine learning model should exhibit robust performance under data distribution shifts in real-world applications [54, 55]. Distribution shift refers to the differences between the underlying distributions of test and training data for a model trained on a specific dataset. Distribution shifts can naturally occur in the physical world due to environmental and conceptual variations [11, 31].In our study, we address three types of distribution shifts. Firstly, we consider natural variations such as object size, occlusion, and rendition changes as an OOD generalization problem, where the test set does not adhere to a specific pattern or concept. For instance, Hendrycks *et al.* proposed ImageNet-A [26], which serves as natural "adversarial" examples that challenge vision systems. ImageNet-R [23] was introduced to encompass a variety of patterns, including art, cartoons, deviantart, and graffiti *etc.* Secondly, we examine common corruptions of visual data that frequently occur in everyday life. For example, autonomous driving vehicles may encounter various weather changes, such as snow, fog, and rain [24]. Lastly, we explore the domain adaptation problem, where the test set differs from the training set but follows a specific pattern or concept, such as DomainNet [30] and VisDA-C [44]. Numerous methods have been investigated to enhance the OOD robustness of machine learning models at different stages. These approaches include pre-training techniques [21, 45, 65], finetuning methods [34, 62], and test-time strategies [16, 61, 63]. By exploring various techniques throughout the model's lifecycle, researchers aim to develop more robust models capable of handling distribution shifts problems. By tackling these distribution shifts, our study aims to improve the robustness and adaptability of machine learning models in real-world applications.

## 3 VPA: VISUAL PROMPT ADAPTATION

In this section, we introduce VPA, which leverages visual prompting for fully test-time adaptation. We first motivate our design. Next, we describe our visual prompt and adaptation designs in §3.1 and §3.2, respectively. Finally, we present our adaptation setups in §3.3.

**Why *Fully* Test-Time?** Although various training-phase (*i.e.*, pre-training and fine-tuning) methods have been proposed to improve model performance and generalization, there is no overarching combination that achieves the best performance. Therefore, test-time adaptation is a desirable complement. As briefly mentioned in §1 and §2, fully test-time adaptation updates the model without preventing inference nor accessing source domain information. Besides, it does not require any supporting dataset [28] and ensures that the adaptation solely relies on the current input (*i.e.*, the batched or single image). We believe such a setup is realistic as it accounts for protecting the privacy and intellectual property of the model and supports domain switches during inference.

**Why Visual Prompting?** Visual prompting has shown great potential in adapting and even reprogramming the model during the training phase [7, 29], making it memory and storage efficient. Visual prompt tuning is also beneficial in terms of faster convergence. Training a language model from scratch can be a time-consuming process; however, incorporating visual prompts can accelerate the training process by providing the model with relevant visual cues. This enables the model to learn and converge more quickly, resulting in improved efficiency and reduced training time. However, there are very few studies that have researched the application of visual prompting in test-time adaptation. Our study serves as a first step towards exploring the effectiveness of visual prompting in this area. Additionally, we integrate visual prompting into existing adaptation frameworks and demonstrate its superiority in improving model generalization and robustness.

### 3.1 Prompt Design

**Taxonomy of Visual Prompts**. In this study, we introduce a novel adaptation design that utilizes visual prompting. While prompting has been extensively studied in language models for various tasks, as discussed in § 2, there exist only a few visual prompt designs aimed at improving recognition performance in the training phase [7, 29]. We formally define these designs as *additive* and *prependitive* prompts and illustrate them using the architecture of the Vision Transformer (ViT) model. Consider a ViT model with $N$ layers, where an input image is divided into $m$ patches $\{I_i | 1 \leq i \leq m\}$. Each patch is then fed into the linear projection layer with positional encoding: $e_0^i = \text{Linear}(I_i)$. We denote the input to the $i$-th Transformer layer as $E_{i-1} = \{e_{i-1}^j | 1 \leq j \leq m\}$, and the $i$-th Transformer layer output as $[CLS_i; E_i] = L_i([CLS_{i-1}; E_{i-1}])$, where the classification head takes the final $CLS_N$ token for prediction: $y = \text{Head}(CLS_N)$. Each Transformer layer consists of a self-attention module, an MLP layer with LayerNorm, and residual connections. The additive prompting is defined as:

$$[CLS_i; E_i] = L_i([CLS_{i-1}; P_{i-1} + E_{i-1}]) \qquad (1)$$

where $P_i = \{p_i^j | 1 \leq j \leq m\}$ and + denotes element-wise addition. Similarly, our prependitive prompting is formulated as:

$$[CLS_i; Z_i; E_i] = L_i([CLS_{i-1}; P_{i-1}; E_{i-1}]) \qquad (2)$$

where the prompts $P$ play as additional tokens and $Z$ is the output corresponding to the input prompts. In the rest of this paper, we use $\oplus$ to denote the attachment of visual prompts *i.e.*, $P \oplus E$. The numbers of visual prompts and layers prompted are configurable.

While we use the ViT architecture to demonstrate our visual prompt designs, it's important to note that additive prompting is compatible with most model architectures, as it only modifies the numeric values without altering the input size. In contrast, prependitive prompting changes the input dimension, which is better suited for ViT-based models, as Transformer blocks are insensitive to the length of the input [12]. In contrast to randomized initialization, we design a **zero attention** scheme to initialize the prompt with zero tensors. This approach ensures that the initialization process does not impact the original performance of the frozen model.

### 3.2 Prompt Adaptation

In this section, we present the test-time adaptation procedure. Our study investigates two setups, namely *episodic* and *continual* adaptations. Episodic adaptation only applies to incoming data, and the model will be reset afterward. In contrast, continual adaptation lasts throughout the entire inference procedure.

A recent study by Goyal et al. [18] has demonstrated that self-entropy minimization is an almost-optimal objective function for episodic test-time adaptation on models trained with cross-entropy loss. Therefore, we adopt the objective of self-entropy minimization in our study. Let $f$ denote a well-trained classifier, and the self-entropy $H(\cdot)$ of a prediction is formulated as follows:

$$H(z, \tau) = - \sum_{i=1}^{c} \sigma(z/\tau)_i \log \sigma(z/\tau)_i \quad z = f(x \oplus P) \qquad (3)$$

Here $c$ is the number of classes, $\sigma(z)_i = \frac{\exp z(i)}{\sum_{j=1}^{c} \exp z(j)}$ is the softmax function, and $\tau$ is a tunable hyper-parameter that controls the softmax temperature. Self-entropy is an unsupervised loss function as it relies only on predictions and not on ground-truth information. However, since entropy reflects the prediction confidence, it can serve as an indicator of the model's performance on the supervised task [61]. As VPA is a general adaptation framework, we leverage both batched- and single-image adaptation settings in our work, as introduced below.

**Batched-Image Adaptation (BIA)**. Real-world machine learning services usually aggregate input data for batched predictions to save computation resources and time [1]. Therefore, we mainly focus on BIA in our study, whose objective is formulated as:

$$\hat{P} = \arg\min_{P} \frac{1}{K} \sum_{i=1}^{K} H(f(x_i \oplus P), \tau) \qquad (4)$$

where $K$ is the batch size. In BIA, VPA optimizes a visual prompt for all the test images in a given batch, which is the same setup as TENT [61], the SOTA method under BIA.

**Single-Image Adaptation (SIA)**. As self-entropy minimization requires batched inputs to function [61], we utilize the setups in MEMO [63] to expand a single image to a batch via augmentations.

**Table 1: OOD Generalization Evaluation Results (%) of VPA on ImageNet Variants. The OOD average accuracy is calculated from the evaluation of ImageNet-A, ImageNet-R, and ObjectNet datasets.**

| Accuracy (↑) | Method | ImageNet | ImageNet-V2 | ImageNet-A | ImageNet-R | ObjectNet | OOD Average |
|---|---|---|---|---|---|---|---|
| Source | CLIP-ViT-LPFT | 81.2 | 71.1 | 49.3 | 71.1 | 52.3 | 57.6 |
| Episodic BIA | TENT (Norm Layer) | 81.3 | 71.3 | 49.6 | 71.8 | 52.6 | 58.0 |
| | TENT (Cls Token) | 81.2 | 71.0 | 49.4 | 71.5 | 52.3 | 57.7 |
| | TENT (All Parameters) | 81.2 | 71.2 | 49.7 | 71.6 | 52.4 | 57.9 |
| | DDA | 77.2 | 65.2 | 38.5 | 65.4 | 46.5 | 50.1 |
| | **Additiv VPA** | **81.3** | **71.4** | **50.4** | **72.0** | **52.8** | **58.4** |
| | **Prependitive VPA** | **81.3** | 71.3 | 50.1 | **72.0** | 52.5 | 58.2 |
| Episodic SIA | MEMO | **81.3** | 72.3 | 52.0 | 72.2 | **52.9** | 59.0 |
| | **Additive VPA** | 81.2 | 72.3 | 49.5 | **72.5** | 52.3 | 58.1 |
| | **Prependitive VPA** | 81.2 | **72.9** | 52.4 | 72.6 | 52.8 | **59.3** |
| Source | CLIP-ViT-WiSE | **79.8** | 70.5 | 49.7 | 71.9 | 52.4 | 58.0 |
| Episodic BIA | TENT (Norm Layer) | 79.6 | 70.7 | 49.8 | 72.2 | 52.8 | 58.3 |
| | TENT (Cls Token) | 79.7 | 70.5 | 50.0 | 72.2 | 52.4 | 58.2 |
| | TENT (All Parameters) | **79.8** | 70.6 | 50.3 | 72.5 | 52.5 | 58.4 |
| | DDA | 70.1 | 62.2 | 41.4 | 64.8 | 46.0 | 50.7 |
| | **Additive VPA** | **79.8** | **71.2** | **52.1** | **72.5** | **52.8** | **59.2** |
| | **Prependitive VPA** | **79.8** | 71.0 | 51.2 | 72.4 | 52.5 | 58.7 |
| Episodic SIA | MEMO | **80.1** | 72.0 | 53.9 | 72.6 | 53.0 | 59.8 |
| | **Additive VPA** | 80.0 | 72.1 | 50.3 | 72.3 | 52.5 | 58.4 |
| | **Prependitive VPA** | 80.0 | **72.5** | **54.2** | **72.7** | **53.2** | **60.0** |
| Source | CLIP-ResNet50×4 | 78.9 | 67.5 | 36.7 | 64.0 | 49.5 | 50.1 |
| Episodic BIA | BN | 78.3 | 67.5 | 27.2 | 55.1 | 40.5 | 40.9 |
| | TENT (Norm Layer) | 78.1 | 67.5 | 27.3 | 55.2 | 40.4 | 41.0 |
| | TENT (All Parameters) | 79.1 | 68.2 | 37.2 | 64.4 | 49.9 | 50.5 |
| | DDA | 69.0 | 61.1 | 24.6 | 56.3 | 41.0 | 40.6 |
| | **Additive VPA** | **79.1** | **68.6** | **37.9** | **65.0** | **49.9** | **50.9** |

The objective of VPA under SIA is formulated as:

$$\hat{P} = \arg\min_{P} H\left(\frac{1}{\eta K} \sum_{i=1}^{K} f(S(\mathcal{A}_i(x_i), \eta) \oplus P, \tau)\right) \quad (5)$$

$$S(x_i, \eta) = x_i \cdot \mathbb{1}[H(z_i, \tau) \leq \arg\text{top-}\eta K\{H(z, \tau)\}]$$

where $\mathcal{A}$ denotes a random augmentation function, $K$ is the augmented batch size, and $S()$ is a confidence selection function to pick augmented images with high confidence with a percentile of $\eta$, following the setting in [49]. The intuition behind SIA is to use marginal entropy minimization on the augmented input to optimize the prompt and enhance generalization.

During both BIA and SIA, the visual prompt $P$ is optimized by computing the gradient of the entropy loss *w.r.t.* $P$ (*i.e.*, $\frac{\partial H}{\partial P}$) during the backward pass. As self-entropy only relies on the network's predictions without labels or source domain information, and our visual prompt $P$ is independent of model parameters, VPA achieves fully test-time adaptation. Additionally, besides episodic test-time adaptation, we also explore the application of continual online learning in test-time adaptation.

**Pseudo-Labeling Adaptation (PLA).** Pseudo-labeling has been widely adopted in semi- and self-supervised learning. In this work, we adopt pseudo-labeling for test-time adaptation, which requires more setup than BIA and SIA. To implement this approach, we use a *memory queue* $\mathcal{M}$ with size $s$ that stores the final $CLS_N$ token, along with the prediction of historic data $z_i$ for reference when processing an incoming batch, *i.e.*, $\mathcal{M} = \{CLS_{Ni}, z_i\}_1^s$. During

adaptation, we generate reference labels for the incoming test data using its $k$ nearest-neighbor ($k$NN) predictions on the $CLS_N$ token before feeding into the head classifier. We then average the $k$NN predictions to produce the eventual pseudo label. We apply weak and strong augmentations to every incoming data sample inspired by FixMatch [50]. Specifically, we obtain the soft predictions for the weakly and strongly augmented samples, denoted as $z_W$ and $z_S$, respectively. We then generate the pseudo label for the incoming data based on the soft voting mechanism in our memory queue:

$$\hat{z}_i = \frac{1}{k} \sum_{j \in k\text{NN}_i} z_{Wj} \quad (6)$$

We apply cross-entropy minimization as our objective, utilizing a temperature hyper-parameter on the generated soft pseudo label:

$$H(z_S, \tau) = -\sum_{i=1}^{c} \sigma(\hat{z}/\tau)_i \log \sigma(z_S)_i \quad z = f(x \oplus P) \quad (7)$$

$$\hat{P} = \arg\min_{P} H(f(x_i \oplus P), \tau) \quad (8)$$

Finally, we dynamically update the $CLS_N$ token along with its prediction $z_W$ into our memory queue for next-round adaptation. Similarly, the visual prompt $P$ is optimized by computing the gradient of the cross-entropy loss.

**Table 2: Corruption Robustness Evaluation Results (%) of VPA on ImageNet-C with the Highest Severity Level.**

| Error Rate (↓) | Method | Gauss. | Shot | Impulse | Defocus | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixelate | JPEG | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | ViT | 52.5 | 51.7 | 51.5 | 56.0 | 69.2 | 51.0 | 56.1 | 46.3 | 50.3 | 45.6 | 24.9 | 69.1 | 55.1 | 34.2 | 33.7 | 49.1 |
| Episodic SIA | MEMO | 48.5 | 47.8 | 45.8 | 58.2 | 68.0 | 53.5 | 57.5 | 42.7 | 48.2 | 41.8 | 22.2 | 68.9 | 52.0 | 32.1 | 27.9 | 47.7 |
| | **Prependitive VPA** | 48.2 | 47.5 | 45.4 | 58.0 | 67.6 | 53.2 | 56.0 | 41.5 | 48.3 | 41.7 | 22.5 | 67.5 | 52.2 | 31.9 | 27.5 | **47.3** |
| | TENT (Norm Layer) | 47.7 | 49.0 | 45.5 | 59.4 | 68.1 | 49.0 | 55.9 | 49.5 | 48.6 | 42.0 | 22.3 | 62.1 | 52.1 | 32.3 | 27.5 | 47.4 |
| Continual BIA | **Additive VPA** | 48.5 | 49.1 | 45.7 | 58.9 | 67.5 | 48.8 | 55.8 | 56.1 | 49.7 | 42.1 | 22.1 | 61.8 | 52.5 | 32.3 | 27.7 | 47.9 |
| | **Prependitive VPA** | 47.0 | 48.5 | 44.2 | 56.8 | 65.7 | 48.2 | 55.5 | 48.2 | 48.0 | 40.1 | 21.8 | 61.4 | 51.3 | 31.2 | 27.2 | **46.5** |
| | AdaContrast | 45.8 | 44.7 | 44.5 | 47.2 | 57.8 | 41.8 | 46.0 | 35.2 | 39.8 | 34.8 | 22.8 | 47.5 | 40.2 | 28.5 | 29.5 | 40.4 |
| Continual PLA | CFA | 43.1 | 42.0 | 41.9 | 45.6 | 51.1 | 40.1 | 43.4 | 33.6 | 35.9 | 32.3 | 21.0 | 41.2 | 35.7 | 28.3 | 29.8 | 37.6 |
| | **Prependitive VPA** | 46.7 | 44.7 | 43.9 | 42.0 | 44.5 | 38.9 | 43.0 | 31.0 | 33.2 | 28.5 | 22.9 | 37.1 | 31.8 | 28.4 | 30.0 | **36.6** |

## 3.3 Adaptation Setups

In our study, we primarily utilize the ViT-B/16 model architecture, which is now considered a standard benchmarking model [12]. Additionally, we use ResNet [22] to demonstrate the general effectiveness of additive VPA. We attach additive prompts to the 1st and 6th Transformer layers in ViT, which are in total $196 \times 2 = 392$ tokens. Similarly. we insert 50 adaptable prompts into every other layer of the ViT-B architecture, resulting in 300 tokens, for prependitive prompting. By default, we set the value of $\tau$ to 1.0 in episodic adaptation, and we analyze its impact in Section 4.4 to demonstrate that an optimal $\tau$ can lead to further improvements. We follow TENT [61] and set 10 adaptation steps for each batch. We also ablate the number of steps in Section 4.4. We empirically set the learning rate to 4.0 and 0.001 with SGD [47] for additive and prependitive prompting, respectively. We use a batch size of 64 for all experiments in Section 4. For SIA, we utilize random cropping [2] as the augmentation function $\mathcal{A}$ and $\mathcal{W}$. In the experimentation of PLA, we set $\tau = 0.07$ and the memory queue size to 1% of the test dataset by default. We utilize $k = 11$ similarly to AdaContrast [6]. Moreover, we leverage random cropping and RandAugment [8] as the weak ($\mathcal{W}$) and strong ($\mathcal{S}$) augmentations, respectively. Most importantly, we compare both *episodic* and *continual* learning settings in TENT [61] with our VPA in this study.

## 4 EXPERIMENTS AND RESULTS

This section reports on the experimental results of VPA and several other baseline methods across multiple benchmarks. As previously noted in § 1, we commence by assessing the performance of VPA in the presence of challenging distribution shifts (*i.e.*, OOD generalization) in § 4.1. Following this, we delve into exploring the potential of VPA in enhancing robustness against common corruptions in § 4.2. Lastly, we evaluate the performance of VPA on the domain adaptation task, alongside other baseline methods, in § 4.3.

## 4.1 Evaluation of OOD Generalization

**Experimental Setups**. We select models pre-trained with CLIP [45], and leverage two SOTA robust fine-tuning methods (*i.e.*, LPFT [34] and WiSE [62]) to train them on the ImageNet training set. To study model robustness to realistic OOD data that naturally occurs in the physical world, we leverage **ImageNet-A** [26], **ImageNet-R** [23], and **ObjectNet** [4]. ImageNet-A consists of 7,500 test images denoted as "natural adversarial examples" that are misclassified by a collection of standard models overlapped with 200 ImageNet

categories. ImageNet-R collects 30,000 images of 200 ImageNet categories with artistic renditions. ObjectNet is a large real-world test set for object recognition with control where object backgrounds, rotations, and imaging viewpoints are random. We chose the ObjectNet subset that overlaps 113 classes with ImageNet in our study. We also use **ImageNet-V2** [46], a robustness benchmark with mild distribution shift, to further validate our results. As these challenging datasets do not follow specific distribution patterns, research has shown that continual learning may not be effective in improving robustness [6, 61]. Therefore, we evaluate VPA under episodic BIA and SIA using different prompt types and measure the in-distribution (ID) and OOD accuracy as the metrics for our evaluation.

*4.1.1 Evaluation on Foundation Models.* We present our large-scale evaluation of models fine-tuned with different methods in this study. Specifically, we utilize additive and prependitive prompting for VPA under BIA and SIA, respectively. As ViT models do not use BatchNorm layers, we default to adapting the LayerNorm layers for TENT. We also experiment with adapting the classification token $CLS_N$ and all model parameters for additional comparisons, following the settings in [32]. Table 1 shows the experimental results of different methods under BIA and SIA. We observe that TENT (Norm) only achieves slight improvements against natural distribution shifts compared to the source-only baseline. While the LayerNorm layer is a linear module that is preferred by the adaptation assumption in TENT [61], it is independent of the input data. Additionally, natural OOD data does not follow a clear distributional pattern, unlike synthesized corruptions. Therefore, the benefits of linear module adaptation do not transfer to our setting [61]. Similarly, applying TENT to other model parameters does not show sensible improvements over the baseline either, which is consistent with prior studies [32]. In contrast, our VPA achieves visible enhancements over the source-only method. Specifically, VPA relatively improves OOD robustness by 3.8% on average, while maintaining or improving the (near) ID accuracy on ImageNet and ImageNet-V2 for all chosen pre-trained models. Conversely, TENT degrades the original performance on (near) ID data in some cases. We also evaluate VPA and MEMO under SIA and find that both methods achieve around 2.7% improvements over the source-only method (Table 1). While MEMO adapts the whole backbone, VPA achieves more efficient prompt-level adaptation.

Our results show that, unlike under BIA where additive and prependitive prompting achieve similar adaptation performance, the prependitive VPA performs better under SIA. As discussed in

**Table 3: Domain Adaptation Evaluation Results (%) of VPA on DomainNet-126 (S: Sketch, R: Real, C: Clipart, I: Infograph, Q: Quickdraw, and P: Painting).**

| Accuracy (↑) | Method | S→R | S→C | S→I | S→Q | S→P | R→S | R→C | R→I | R→Q | R→P | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | CLIP-ViT-LPFT | 67.5 | 69.2 | 35.5 | 16.4 | 52.8 | 55.2 | 68.6 | 44.4 | 9.7 | 60.9 | 48.0 |
| Episodic SIA | MEMO | 67.8 | 69.5 | 35.4 | 16.8 | 53.2 | 55.5 | 68.9 | 44.8 | 10.1 | 61.3 | 48.3 |
| | **Prependitive VPA** | 68.1 | 70.1 | 36.5 | 17.3 | 53.6 | 55.8 | 69.5 | 45.0 | 10.6 | 61.5 | **48.8** |
| Continual BIA | TENT (Norm Layer) | 67.7 | 69.0 | 35.8 | 16.6 | 53.0 | 55.3 | 68.9 | 44.5 | 9.9 | 61.0 | 48.2 |
| | **Prependitive VPA** | 68.7 | 69.8 | 36.3 | 16.8 | 53.2 | 55.8 | 69.0 | 44.8 | 10.1 | 61.3 | **48.6** |
| Continual PLA | AdaContrast | 70.2 | 72.0 | 36.5 | 18.0 | 54.2 | 58.4 | 71.8 | 45.5 | 10.5 | 63.3 | 50.0 |
| | CFA | 69.8 | 72.2 | 35.9 | 17.8 | 54.4 | 58.7 | 71.6 | 45.7 | 10.3 | 63.0 | 49.9 |
| | **Prependitive VPA** | 70.8 | 73.0 | 37.0 | 17.8 | 55.3 | 58.8 | 71.9 | 46.4 | 10.4 | 63.8 | **50.5** |

§ 3.2 and illustrated in Figure 1, SIA realizes test-time adaptation via augmentations over a single image. The workflow of SIA is as follows: 1) expand a single image as a mini-batch through augmentations for self-entropy minimization, 2) discard the expanded batch after adaptation, and 3) use the original single image and adapted model/prompt for inference. This setting is designed for *model-level* adaptation, as the model is trained to be insensitive to input augmentations. However, *additive prompts* are directly added on top of the input with a small magnitude, which makes their effectiveness sensitive to any change in the images. The prompt is adapted to the augmented batch but added to the original image in SIA, resulting in reduced effectiveness. On the other hand, prependitive prompting does not directly modify the semantics of the input, and the optimization over a single image is easier to converge than BIA with an appropriate prompt length. Our experiments show that prependitive prompting achieves a 3.3% improvement under SIA, while additive prompting does not provide a tangible enhancement compared to the source-only baseline method.

We conduct another experiment of VPA on the ResNet architecture. Specifically, we use ResNet50x4 pre-trained with CLIP, and Table 1 presents the evaluation results, where VPA achieves similar improvement on the ResNet model [62]. Surprisingly, we find that TENT has an around 20% performance drop on ResNet50x4. TENT, by default, replaces the original statistics (*i.e.*, $\mu$ and $\sigma$) in BatchNorm layers with the statistics of the input data. This setting is useful when the input batch is from one specific domain (*e.g.*, synthesized corruptions). However, natural distribution shifts do not follow such assumptions, rendering significant performance degradation. In comparison, VPA consistently achieves better OOD robustness on different fine-tuning methods and model architectures.

## 4.2 Evaluation of Corruption Robustness

**Experimental Setups**. In this section, we evaluate the performance of VPA against common corruptions using the **ImageNet-C** dataset [24]. ImageNet-C is designed to assess the robustness and generalization capabilities of computer vision models by introducing 15 different corruptions at five severity levels to the original ImageNet validation dataset. These corruptions include various types of noise, blur, and distortion, making ImageNet-C a more realistic and challenging test of model robustness and generalization. We adopt the methodology from [6], focusing on the highest corruption severity level's performance. Our primary interest lies in the continual adaptation setting, as it has been shown to be
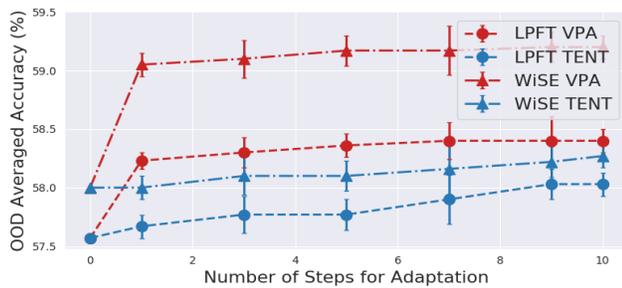
more effective in enhancing corruption robustness; this is because each corruption can be considered as being drawn from a similar distribution shift. Based on our experiments, which demonstrated that the prependitive prompt outperforms the additive design under SIA, we employ only prependitive prompts for this evaluation. As foundation models do not show visible improvements on corruption robustness benchmarks, we use the ViT model pretrained on ImageNet in this section. In addition, we compare VPA with MEMO, a SOTA episodic adaptation method. Error rate serves as our evaluation metric for assessing corruption robustness.

Table 2 presents the experimental results for the highest severity level, demonstrating that our prependitive VPA consistently achieves the best robustness improvement among all the baselines. Importantly, VPA attains the highest robustness under episodic SIA and provides greater storage efficiency compared to MEMO. This efficiency results from VPA adapting only the additional prompts, whereas MEMO adapts all the model parameters. For the evaluation under continual BIA, our prependitive VPA outperforms TENT and the additive design by 1.9% and 3.0%, respectively. This outcome may be associated with the nature of corruption benchmarks. Additive prompts are directly added to the embedding of the corrupted input, while prependitive prompts do not directly mix with the original corrupted embedding. Instead, they leverage the attention mechanism to interact with the original embedding, leading to a better robustness gain. Continual PLA achieves the largest gain, benefiting from the pseudo labels generated by the memory queue. Our prependitive VPA outperforms AdaContrast [6] and CFA [32] by relative margins of 10.4% and 2.7%, respectively.
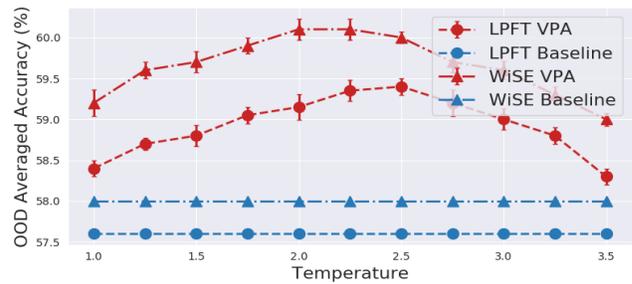
## 4.3 Evaluation of Domain Adaptation

**Experimental Setups**. In this section, we discuss our experiments and results related to the domain adaptation task. We employ the **DomainNet-126** dataset [43] for this purpose. DomainNet encompasses common objects from six domains (*i.e.*, sketch, real, clipart, infograph, quickdraw, and painting) and 345 categories. In our study, we empirically use the sketch (S) and real (R) images as the training sets and evaluate the adaptation performance on the remaining subsets. We leverage the ViT model pretrained with CLIP and fine-tuned by LPFT in this evaluation. For other setups, we follow the same configuration as used in the corruption robustness evaluation.

Table 3 displays the evaluation results for DomainNet-126, where we observe that VPA consistently achieves the best robustness

(a) Different Number of Adaptation Steps.



(b) Different Softmax Temperatures $\tau$.

**Figure 2: Ablation Studies of Visual Prompt Adaptation on OOD Generalization. We show that the first adaptation step contributes the most in VPA and additional improvements can be achieved with an optimal temperature hyperparameter $\tau$.**

**Table 4: Ablation Study of the Prompt Size in Episodic VPA.**

| Performance (%) | Prompt Size | OOD (↑) | Corruption (↓) | DA (↑) |
|---|---|---|---|---|
| Additive VPA | 196 | 58.5 | 48.3 | - |
| | 392 | **58.8** | **47.9** | - |
| | 588 | 58.7 | 48.1 | - |
| Prependitive VPA | 150 | 58.1 | 47.0 | 48.2 |
| | 300 | **58.5** | **46.5** | **48.6** |
| | 450 | 58.3 | 46.8 | 48.4 |

across various adaptation settings. DomainNet shares certain similarities with ImageNet-C, as each domain exhibits a specific pattern. On average, our VPA outperforms the source-only baseline by 5.2%. Furthermore, PLA-based VPA demonstrates more substantial improvements, with relative margins of 1.0% and 1.1% for AdaContrast and CFA, respectively. This highlights the effectiveness of our VPA approach in addressing domain adaptation challenges while maintaining robust performance across different settings.

## 4.4 Ablation Studies

Besides visual prompt designs, this section provides an empirical analysis of VPA on different hyper-parameter settings under different adaptation settings.

**Prompt Size**. We begin by conducting ablation studies on the prompt size utilized in additive and prependitive VPA. As outlined in § 3.3, we utilize 50 adaptable prompts in 6 layers as the default for our prependitive VPA and 196 adaptable prompts in 2 layers of the ViT-B model for our additive VPA, leading to 300 and 392 learnable tokens, respectively. We vary the number of prompts adapted and assess the VPA's performance on OOD generalization and corruption robustness. Specifically, we vary the number of prompted layers for the additive VPA. The results are presented in Table 4. The evaluation highlights the existence of an optimal point for the number of prompting tokens. Having too many prompts can make it difficult to optimize, while a relatively small number of prompts may restrict the capability of VPA.

**Adaptation Steps and Temperature $\tau$**. We then ablate the effect of adaptation steps and the temperature parameter of additive episodic adaptation. In Figure 2(a), the average OOD accuracy generally improves as the number of steps increases. Encouragingly, we find the first step of VPA contributes the most to the OOD robustness improvement, where the relative improvements are 1.2% and

**Table 5: OOD Generalization Evaluation Results (%) of VPA for Zero-Shot Recognition in the Vision-Language Model.**

| | ImageNet | ImageNet-V2 | ImageNet-A | ImageNet-R |
|---|---|---|---|---|
| CLIP | 66.7 | 60.9 | 47.9 | 74.0 |
| CLIP+TPT | 69.0 | 63.4 | 53.5 | 76.5 |
| CLIP+TPT+**VPA** | **69.1** | **63.7** | **53.9** | **77.0** |

1.9% for models fine-tuned with LPFT and WiSE, respectively. In contrast, one-step TENT adaptation shows no improvements over the source-only baseline, demonstrating the effectiveness of VPA. Figure 2(b) shows that there is a sweet point for the temperature parameter for OOD robustness improvement: We find that with an optimal temperature, the OOD robustness of the LPFT model could further improve by 1.5% on average. However, selecting an optimal $\tau$ requires an additional validation set with access to the label, so we do not tune the temperature $\tau$ in the central part of our evaluation. We leave this as a future study to automatically select the temperature parameter for VPA.

**Vision-Language Model**. We evaluate the performance of our proposed VPA combined with TPT [49] on the zero-shot recognition of the CLIP model. We adopt the SIA experimental setup from [49]. The results in Table 5 demonstrate that our VPA consistently enhances the zero-shot recognition performance on the challenging OOD generalization benchmarks by approximately 0.5%.

## 5 CONCLUDING REMARKS

To conclude, we propose VPA, a pioneering framework for generalizing visual prompting with test-time adaptation. VPA effectively improves OOD generalization, corruption robustness, and domain adaptation performance across diverse settings and tasks. The effectiveness of VPA highlights the potential of incorporating visual prompting in future research to address a wide array of adaptation challenges.

## ACKNOWLEDGMENTS

# REFERENCES

[1] 2022. Machine learning inference during deployment. https://learn.microsoft.com/en-us/azure/cloud-adoption-framework/innovate/best-practices/ml-deployment-inference#batch-inference.

[2] 2022. Random Cropping in Pytorch. https://pytorch.org/vision/main/generated/torchvision.transforms.RandomCrop.html.

[3] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. 2022. Exploring visual prompts for adapting large-scale models. *arXiv preprint arXiv:2203.17274* 1, 3 (2022), 4.

[4] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. 2019. ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2019/file/97af07a14cacba681feacf3012730892-Paper.pdf

[5] Nicholas Carlini and David Wagner. 2017. Towards Evaluating the Robustness of Neural Networks. In *2017 IEEE Symposium on Security and Privacy (SP)*. 39–57. https://doi.org/10.1109/SP.2017.49

[6] Dian Chen, Dequan Wang, Trevor Darrell, and Sayna Ebrahimi. 2022. Contrastive Test-Time Adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 295–305.

[7] Xinlei Chen and Kaiming He. 2021. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15750–15758.

[8] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. 2020. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*. 702–703.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[11] Akshay Raj Dhamija, Manuel Günther, and Terrance Boult. 2018. Reducing network agnostophobia. *Advances in Neural Information Processing Systems* 31 (2018).

[12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

[13] Cian Eastwood, Ian Mason, Christopher KI Williams, and Bernhard Schölkopf. 2021. Source-free adaptation to measurement shift via bottom-up feature restoration. *arXiv preprint arXiv:2107.05446* (2021).

[14] Gamaleldin F Elsayed, Ian Goodfellow, and Jascha Sohl-Dickstein. 2018. Adversarial reprogramming of neural networks. *arXiv preprint arXiv:1806.11146* (2018).

[15] Yulu Gan, Yan Bai, Yihang Lou, Xianzheng Ma, Renrui Zhang, Nian Shi, and Lin Luo. 2023. Decorate the newcomers: Visual domain prompt for continual test time adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 7595–7603.

[16] Jin Gao, Jialing Zhang, Xihui Liu, Trevor Darrell, Evan Shelhamer, and Dequan Wang. 2022. Back to the Source: Diffusion-Driven Test-Time Adaptation. *arXiv preprint arXiv:2207.03442* (2022).

[17] Yunhe Gao, Xingjian Shi, Yi Zhu, Hao Wang, Zhiqiang Tang, Xiong Zhou, Mu Li, and Dimitris N Metaxas. 2022. Visual Prompt Tuning for Test-time Domain Adaptation. *arXiv preprint arXiv:2210.04831* (2022).

[18] Sachin Goyal, Mingjie Sun, Aditi Raghunathan, and Zico Kolter. 2022. Test-time adaptation via conjugate pseudo-labels. *arXiv preprint arXiv:2207.09640* (2022).

[19] Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang. 2021. Ppt: Pre-trained prompt tuning for few-shot learning. *arXiv preprint arXiv:2109.04332* (2021).

[20] Xu Han, Weilin Zhao, Ning Ding, Zhiyuan Liu, and Maosong Sun. 2021. Ptr: Prompt tuning with rules for text classification. *arXiv preprint arXiv:2105.11259* (2021).

[21] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16000–16009.

[22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[23] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. 2021. The Many Faces of Robustness: A Critical Analysis of Out-of-Distribution Generalization. *ICCV* (2021).

[24] Dan Hendrycks and Thomas Dietterich. 2019. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261* (2019).

[25] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. 2019. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781* (2019).

[26] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. 2021. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15262–15271.

[27] Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*. PMLR, 448–456.

[28] Yusuke Iwasawa and Yutaka Matsuo. 2021. Test-time classifier adjustment module for model-agnostic domain generalization. *Advances in Neural Information Processing Systems* 34 (2021), 2427–2440.

[29] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual prompt tuning. *arXiv preprint arXiv:2203.12119* (2022).

[30] Youngeun Kim, Donghyeon Cho, Kyeongtak Han, Priyadarshini Panda, and Sungeun Hong. 2021. Domain adaptation without source data. *IEEE Transactions on Artificial Intelligence* 2, 6 (2021), 508–518.

[31] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. 2021. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*. PMLR, 5637–5664.

[32] Takeshi Kojima, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Robustifying Vision Transformer without Retraining from Scratch by Test-Time Class-Conditional Feature Alignment. *arXiv preprint arXiv:2206.13951* (2022).

[33] Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. 2016. Ask me anything: Dynamic memory networks for natural language processing. In *International conference on machine learning*. PMLR, 1378–1387.

[34] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. 2022. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054* (2022).

[35] Jogendra Nath Kundu, Naveen Venkat, R Venkatesh Babu, et al. 2020. Universal source-free domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4544–4553.

[36] Yann LeCun, Yoshua Bengio, et al. 1995. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks* 3361, 10 (1995), 1995.

[37] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436–444.

[38] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691* (2021).

[39] Rui Li, Qianfen Jiao, Wenming Cao, Hau-San Wong, and Si Wu. 2020. Model adaptation: Unsupervised domain adaptation without source data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9641–9650.

[40] Jian Liang, Dapeng Hu, and Jiashi Feng. 2020. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine Learning*. PMLR, 6028–6039.

[41] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586* (2021).

[42] Zachary Nado, Shreyas Padhy, D Sculley, Alexander D'Amour, Balaji Lakshminarayanan, and Jasper Snoek. 2020. Evaluating prediction-time batch normalization for robustness under covariate shift. *arXiv preprint arXiv:2006.10963* (2020).

[43] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. 2019. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*. 1406–1415.

[44] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. 2017. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924* (2017).

[45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*. PMLR, 8748–8763.

[46] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. 2019. Do imagenet classifiers generalize to imagenet?. In *International Conference on Machine Learning*. PMLR, 5389–5400.

[47] Sebastian Ruder. 2016. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747* (2016).

[48] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. 2020. Improving robustness against common corruptions by covariate shift adaptation. *Advances in Neural Information Processing Systems* 33 (2020), 11539–11551.

[49] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and Chaowei Xiao. 2022. Test-time prompt tuning for zero-shot generalization in vision-language models. *arXiv preprint arXiv:2209.07511* (2022).

[50] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. 2020. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems* 33 (2020), 596–608.

[51] Jiachen Sun, Yulong Cao, Qi Alfred Chen, and Z. Morley Mao. 2020. Towards Robust LiDAR-based Perception in Autonomous Driving: General Blackbox Adversarial Sensor Attack and Countermeasures. In *29th USENIX Security Symposium (USENIX Security 20)*. USENIX Association, 877–894. https://www.usenix.org/conference/usenixsecurity20/presentation/sun

[52] Jiachen Sun, Yulong Cao, Christopher B Choy, Zhiding Yu, Anima Anandkumar, Zhuoqing Morley Mao, and Chaowei Xiao. 2021. Adversarially robust 3d point cloud recognition using self-supervisions. *Advances in Neural Information Processing Systems* 34 (2021), 15498–15512.

[53] Jiachen Sun, Karl Koenig, Yulong Cao, Qi Alfred Chen, and Z Morley Mao. 2020. On adversarial robustness of 3d point cloud classification under adaptive attacks. *arXiv preprint arXiv:2011.11922* (2020).

[54] Jiachen Sun, Akshay Mehra, Bhavya Kailkhura, Pin-Yu Chen, Dan Hendrycks, Jihun Hamm, and Z Morley Mao. 2022. A Spectral View of Randomized Smoothing Under Common Corruptions: Benchmarking and Improving Certified Robustness. In *European Conference on Computer Vision*. Springer, 654–671.

[55] Jiachen Sun, Weili Nie, Zhiding Yu, Z Morley Mao, and Chaowei Xiao. 2022. Pointdp: Diffusion-driven purification against adversarial attacks on 3d point cloud recognition. *arXiv preprint arXiv:2208.09801* (2022).

[56] Jiachen Sun, Qingzhao Zhang, Bhavya Kailkhura, Zhiding Yu, Chaowei Xiao, and Z Morley Mao. 2022. Benchmarking robustness of 3d point cloud recognition against common corruptions. *arXiv preprint arXiv:2201.12296* (2022).

[57] Jiachen Sun, Haizhong Zheng, Qingzhao Zhang, Atul Prakash, Z Morley Mao, and Chaowei Xiao. 2023. CALICO: Self-Supervised Camera-LiDAR Contrastive Pre-training for BEV Perception. *arXiv preprint arXiv:2306.00349* (2023).

[58] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A Efros, and Moritz Hardt. 2019. Test-time training for out-of-distribution generalization. (2019).

[59] Fatemeh Vakhshiteh, Ahmad Nickabadi, and Raghavendra Ramachandra. 2021. Adversarial attacks against face recognition: A comprehensive study. *IEEE Access* 9 (2021), 92735–92756.

[60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[61] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. 2020. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726* (2020).

[62] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. 2022. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7959–7971.

[63] Marvin Zhang, Sergey Levine, and Chelsea Finn. 2021. Memo: Test time robustness via adaptation and augmentation. *arXiv preprint arXiv:2110.09506* (2021).

[64] Qingzhao Zhang, Shengtuo Hu, Jiachen Sun, Qi Alfred Chen, and Z Morley Mao. 2022. On adversarial robustness of trajectory prediction for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15159–15168.

[65] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. 2021. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832* (2021).

[66] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16816–16825.

[67] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to prompt for vision-language models. *International Journal of Computer Vision* 130, 9 (2022), 2337–2348.

[68] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*. 2223–2232.

# APPENDICES

## A  DISCUSSION

In this paper, we have investigated the application of visual prompting in fully online test-time adaptation. Although prompting has been extensively studied in NLP tasks and has recently gained significant attention for improving the zero-shot performance of vision-language models, its exploration in vision systems remains limited. In NLP tasks, prompted embedding tends to lose semantics after tuning, unlike na"ive prompt engineering [67]. Similarly, explaining the operational principle of visual prompting remains challenging. As introduced in [14], a more general interpretation of visual prompting involves reprogramming a well-trained vision model to achieve any deterministic goal. Despite the inherent limitations of fully test-time visual prompting, we have explored its application across three critical aspects of real-world machine learning systems: OOD generalization, corruption robustness, and domain adaptation, demonstrating its effectiveness. Our experimental results generally show that prependitive VPA is more effective, whereas additive VPA is more universal to different model architectures. We believe visual prompting could also be applied to other tasks as ViT architectures advance and become more dominant. Another avenue for future research is prompt design. In this work, we have explored four combinations of prompting setups in § 3.1; additional design choices, including image-to-image models [68] and embedding-space prompting [29], could yield intriguing applications in the future. It is important to note that our study's focus is not solely on achieving state-of-the-art results but rather on examining how visual prompting performs within the fully test-time adaptation framework. We are aware of a concurrent, yet unpublished work, DePT [17], which investigates a similar topic. DePT primarily focuses on offline test-time adaptation, where the model is adapted offline. Offline adaptation offers more room for sophisticated prompting designs. We also noticed another concurrent work [15] that studies a similar problem. We believe that all three studies are complementary, each providing valuable insights into different aspects of test-time adaptation using visual prompting.

## B  ABLATION STUDIES

We further conducted two ablation studies on the augmentation and $k$NN soft majority voting in PLA.

**Table 6: Ablation Study on Augmentation Method in VPA.**

| PLA (%) | ImageNet-C ($\downarrow$) | DomainNet-126 ($\uparrow$) |
|---|---|---|
| RandAugment | 36.6 | 50.5 |
| AugMix | **35.7** | **51.1** |

In single-image adaptation (SIA), we employ random cropping as the preferred augmentation technique. The weak augmentation employed in pseudo-label adaptation also utilizes random cropping, whereas the strong augmentation is constructed based on RandAugment [8]. We here utilize AugMix [25] to evaluate the impact of distinct augmentation techniques. Our findings reveal that the usage of AugMix can further boost the performance of Visual Prompting Adaptation (VPA) within the pseudo-label adaptation framework by approximately 0.8%, as shown in Table 6. These findings illustrate the considerable potential for performance enhancement through the thoughtful selection and application of image augmentation techniques.

We default to leverage $k = 11$ in the soft majority voting. We further ablate the importance of $k$ in this experiment. As presented in Table 7, $k$ is indeed an essential hyper-parameter in pseudo-label adaptation. We find that $k \in [11, 15]$ generally achieves the highest performance gain.

**Table 7: Ablation Study on $k$NN Soft Majority Voting.**

| PLA (%) | ImageNet-C ($\downarrow$) |
|---|---|
| $k = 3$ | 44.1 |
| $k = 7$ | 39.0 |
| $k = 11$ | 36.6 |
| $k = 15$ | **36.4** |
| $k = 21$ | 38.2 |