



Diffusion-Enhanced Test-Time Adaptation with Text and Image Augmentation

Chun-Mei Feng¹ · Yuanyang He² · Jian Zou³ · Salman Khan^{4,5} · Huan Xiong^{3,4} · Zhen Li⁶ · Wangmeng Zuo³ · Rick Siow Mong Goh¹ · Yong Liu¹

Received: 1 August 2024 / Accepted: 26 September 2024

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2025

Abstract

Existing test-time prompt tuning (TPT) methods focus on single-modality data, primarily enhancing images and using confidence ratings to filter out inaccurate images. However, while image generation models can produce visually diverse images, single-modality data enhancement techniques still fail to capture the comprehensive knowledge provided by different modalities. Additionally, we note that the performance of TPT-based methods drops significantly when the number of augmented images is limited, which is not unusual given the computational expense of generative augmentation. To address these issues, we introduce IT³A, a novel test-time adaptation method that utilizes a pre-trained generative model for multi-modal augmentation of each test sample from unknown new domains. By combining augmented data from pre-trained vision and language models, we enhance the ability of the model to adapt to unknown new test data. Additionally, to ensure that key semantics are accurately retained when generating various visual and text enhancements, we employ cosine similarity filtering between the logits of the enhanced images and text with the original test data. This process allows us to filter out some spurious augmentation and inadequate combinations. To leverage the diverse enhancements provided by the generation model across different modalities, we have replaced prompt tuning with an adapter for greater flexibility in utilizing text templates. Our experiments on the test datasets with distribution shifts and domain gaps show that in a zero-shot setting, IT³A outperforms state-of-the-art test-time prompt tuning methods with a 5.50% increase in accuracy.

Keywords Test time adaptation · Multi-modal learning · Generative models

1 Introduction

Recent advancements have shown that pre-trained vision-language models like CLIP (Radford et al., 2021) perform exceptionally well on a range of downstream tasks, without requiring specific task-related training data (Zhou et al., 2022b, a; Li et al., 2022b; Ramesh et al., 2022). Although their success is attributed to well-crafted prompts, the limitations of hand-crafted prompts and prompt tuning stem from their reliance on the training data distribution within the current domain, making it difficult to generalize to new distributions, particularly in zero-shot settings (Mandal et al., 2019). To address this issue, a technique called test-time prompt tuning (TPT) (Shu et al., 2022) has been introduced, which adapts prompt embeddings for each test sample from an unseen domain in real-time, without requiring training data or annotations. This approach is more practical for dynamic real-world applications where acquiring extensive labeled data for a new target distribution is often problematic.

Communicated by Long Yang.

✉ Zhen Li
lizhen@cuhk.edu.cn
Chun-Mei Feng
fengcm.ai@gmail.com

¹ Institute of High Performance Computing (IHPC), Agency for Science, Technology and Research (A*STAR), Singapore, Singapore

² National University of Singapore, Singapore, Singapore

³ Harbin Institute of Technology, Harbin, China

⁴ Mohamed bin Zayed University of Artificial Intelligence (MBZUAI), Abu Dhabi, UAE

⁵ Australian National University, Canberra, ACT, Australia

⁶ Chinese University of Hong Kong, Shenzhen, China

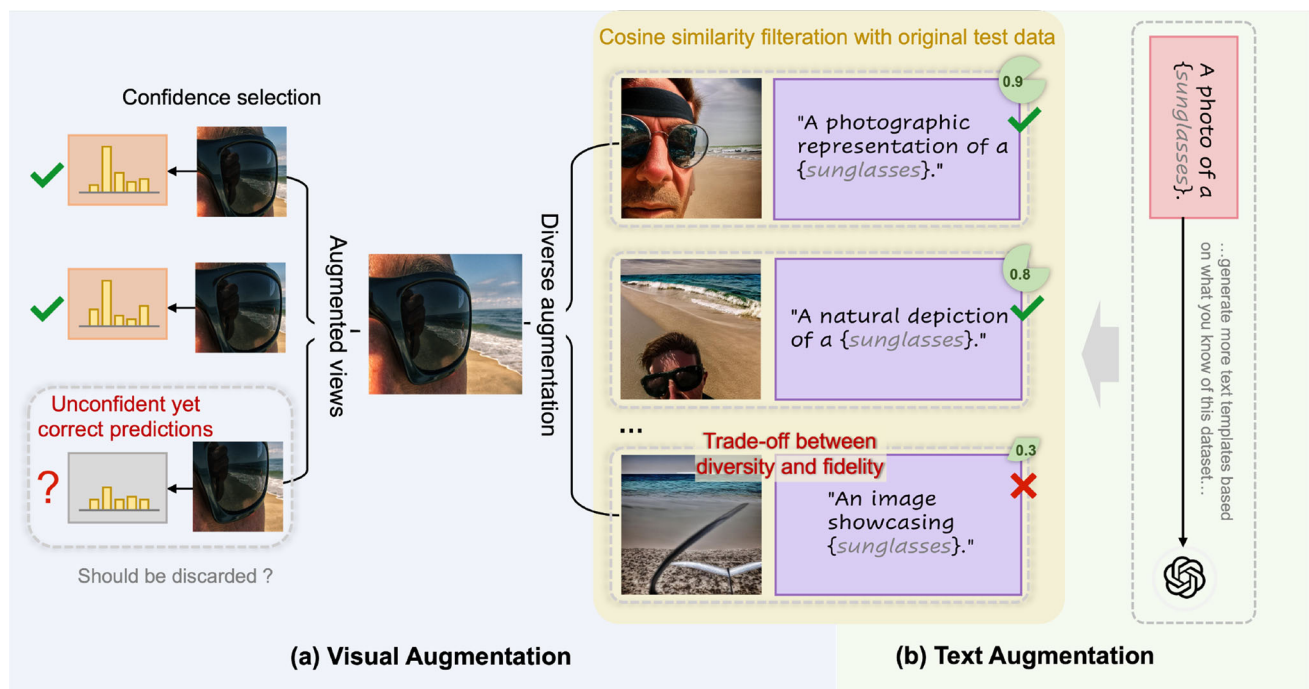


Fig. 1 **a** Visual augmentation of different views (TPT (Shu et al., 2022)) and data generated through diffusion, demonstrating a richer variety of visual appearances (DiffTPT (Feng et al., 2023b)). **b** Diverse text augmentation produced by GPT-4. Additionally, the augmentation

from different modalities, *i.e.*, images and text, will be combined into image-text augmentation pairs, which will then be filtered using cosine similarity to remove low-quality augmentations

The early practice of TPT involves combining confidence selection with entropy minimization for prompt tuning, utilizing various augmented views of each test sample (Shu et al., 2022). The augmentation method employed in Shu et al. (2022) relies on basic parametric transformations for addressing data scarcity (see Fig. 1a). However, these simple transformations fail to bring diversity in semantics into augmented views (Antoniou et al., 2017; Perez & Wang, 2017; Zhao et al., 2020; Shorten & Khoshgoftaar, 2019). The under-diversified augmented data may result in the learned prompt fitting only to the original image details other than the key semantics, thereby compromising its generalization capability. Moreover, the entropy-based confidence selection method proposed in Shu et al. (2022) does not sufficiently ensure prediction accuracy, as augmented samples with low-entropy predictions can still be misclassified, producing unrepresentative samples in the augmented pool.

Recent developments in image generation technology have significantly improved the handling of varied augmented data. Traditional image generation techniques, such as VAEs (Kingma & Welling, 2013) and GANs (Goodfellow et al., 2020), typically necessitate large datasets for effective training. In contrast, diffusion models have recently demonstrated exceptional capabilities in generating text-to-image outputs with impressive photo-realistic quality (Nichol et al., 2021; Ramesh et al., 2022; Saharia et al., 2022; Rom-

bach et al., 2022). Unlike the augmentation method utilized in Shu et al. (2022), data produced by diffusion models can display much greater diversity, leading to richer visual representations and enhancing the generalization of the learned prompts. However, while image generation models can produce visually different images, the information provided by unimodal data augmentation techniques remains limited. Furthermore, continuously augmenting images to improve performance is limited by computational resources. In other words, the performance of those methods significantly drops when the number of augmented images is constrained. This is not surprising given the high computational cost of generating image augmentations. Fortunately, with the advancement of large language models (LLMs), augmenting text using pre-trained language models to create image-text pairs with consistent content but varying styles can effectively supplement the insufficient information of the image modality.

In this work, we introduce a novel test-time adapter-based tuning method called IT³A, which enhances data diversity for test samples using pre-trained vision and language models, thereby improving the model's adaptability to unknown new test data. We also employ cosine similarity filtering to eliminate some spurious augmentations and inadequate combinations (see Figs. 1 and 3). For data diversity, IT³A utilizes both Stable Diffusion and GPT-4 for multi-modal data augmentation. Stable Diffusion is a text-to-image gen-

eration model that synthesizes images based on CLIP text features (Rombach et al., 2022). Instead of using CLIP text features, we leverage the CLIP image features of the test samples and input them into Stable Diffusion for image enhancement. This diffusion-based augmentation effectively generates a variety of images with rich visual appearance changes while retaining key semantics. We then adopt GPT-4 with specific instructions to generate multiple text templates that vary in style but maintain consistent semantics, pairing them with the augmented images. Please note that the GPT-4 can be substituted by any other open-source large language model, *e.g.*, LLaMA (Touvron et al., 2023), BELLE (BELLEGroup, 2023), Bloom (Le Scao et al., 2023), Vicuna (Chiang et al., 2023), and MOSS (Sun et al., 2024). To ensure prediction fidelity, we introduce cosine similarity-based filtering between the logits of the test data and the generated image-text pairs, which helps filter out spurious augmentations and inadequate combinations (see Fig. 3), allowing the model to strike a fair balance between diversity and fidelity. Experimental results indicate that, compared to state-of-the-art test-time prompt tuning methods, the IT³A approach improves zero-shot accuracy by an average of 5.50% (Feng et al., 2023b). In summary, our contributions are as follows:

To sum up, our contributions are as follows:

- We propose a novel test-time adaptation method, IT³A, that simultaneously leverages *text and image augmentations*, utilizing the diverse features generated by generative models for enhanced adaptation during testing.
- To ensure key semantics are faithfully preserved while generating *diverse* visual and textual augmentations, we employ cosine similarity filtering between the *logits* of the *original* test data and the *augmented images and texts*. This process removes spurious augmentation pairs, thereby improving the predictive accuracy of the enhanced images.
- Experimental results demonstrate that our IT³A method significantly outperforms state-of-the-art test-time tuning techniques.

Our initial findings were presented at ICCV 2023 (Feng et al., 2023b). This journal version offers three major enhancements: *Firstly*, we explore the potential of leveraging the diversity of pre-trained text generation models for test-time optimization; *e.g.*, using GPT-4 to generate various text templates based on its pre-trained knowledge. *Secondly*, to filter out potential spurious information in the generated image and text pairs, we move the cosine similarity between the original test images and the generated images to between the logits of the original test data and the augmented images and texts, thereby removing generated data with low quality. *Thirdly*, to fully leverage the diverse augmentations of images and text provided by generative models, we have replaced prompt

tuning in our conference version (Feng et al., 2023b) with an adapter. This change allows for more flexible use of text templates on the text encoder side. *Lastly*, in addition to comparisons with state-of-the-art methods, we also conducted various ablation studies to demonstrate the effectiveness of IT³A's improvements.

2 Related Work

2.1 Parameter-Efficient Fine-Tune

Large-scale pre-trained models have significantly boosted performance across numerous tasks in both natural language processing (Devlin et al., 2018; Radford et al., 2018) and computer vision (Jia et al., 2021; Chen et al., 2020; Jia et al., 2022; Feng et al., 2023a; Li et al., 2022a). These improvements are achieved by learning comprehensive representations and transferring the acquired knowledge to a variety of downstream applications. In recent years, a variety of parameter-efficient fine-tuning techniques, such as prompt tuning and adapters, have been developed to tailor pre-trained models for specific downstream tasks. One such technique, prompt tuning, allows pre-trained models to directly adapt to downstream tasks by incorporating a small set of trainable tokens into the input space. For example, CoOp (Zhou et al., 2022b) utilizes continuous prompt optimization, while CoCoOp (Zhou et al., 2022a) employs instance-wise prompt conditionalization, both methods aim to improve generalization to out-of-distribution data. Numerous studies employ adapters and non-parametric key-value cache approaches to fine-tune the CLIP model, enhancing its adaptability to specific target datasets. As an example, CLIP-Adapter incorporated a feature adapter to refine the CLIP model, enabling it to learn new features while preserving a straightforward architecture (Gao et al., 2021). Conversely, Tip-Adapter utilized a non-parametric key-value cache approach for training the adapter. This method bypasses backpropagation and enhances the model's adaptability to the target dataset (Zhang et al., 2021b). UPL reduces CLIP's dependence on labeled data by training an ensemble of prompt representations to enhance transfer performance without requiring target dataset labels (Huang et al., 2022). Nevertheless, its zero-shot generalization effectiveness is significantly influenced by the quality of the prompt design.

From a different perspective, Shu et al. (2022) introduced the concept of test-time prompt tuning (TPT) by generating varied augmented perspectives of individual test samples, which can be effectively utilized for zero-shot generalization of the base model (Shu et al., 2022). Yet, the data augmentation techniques employed in TPT (Shu et al., 2022) are hindered by excessively basic variations, while relying solely on entropy-based confidence selection may not con-

sistently ensure prediction accuracy. To enhance TPT (Shu et al., 2022), our conference version, DiffTPT (Feng et al., 2023b), proposes incorporating diffusion-based data augmentation and utilizing cosine similarity-based filtration to strike a better balance between data diversity and prediction accuracy. In light of the versatility offered by the adapter, we have opted to swap out prompt tuning in DiffTPT for adapter. This modification allows for the comprehensive utilization of the varied enhancements in both images and text facilitated by generative models.

2.2 Test-Time Adaptation

Adapting machine learning models to test samples presents a more difficult and realistic scenario, as it involves the absence of training data during the inference phase (Wang et al., 2020; Sun et al., 2020; Chen et al., 2022; Shanmugam et al., 2021). This approach addresses the issues of source data being inaccessible due to privacy reasons and allows for a single training session of the model, which can then be tailored to any unforeseen test distributions (Gao et al., 2022). An effective approach to creating a robust test-time objective is to decouple it from any particular training methodology. This can be achieved by either minimizing the entropy of the prediction probability distribution for the batch (Wang et al., 2020), or by eliminating the need for multiple test samples through the use of data augmentation techniques (Zhang et al., 2021a). For instance, an additional branch can be implemented to tailor the model to test samples by refining the objective during the testing phase (Sun et al., 2020; Liu et al., 2021). Wang et al. (2020) enhanced model confidence at test time by utilizing the model's own predictions for self-adjustment, thereby minimizing the generalization error on data exhibiting distribution shifts. An alternative method involves the explicit use of the Batch Normalization (BN) layer during test time to limit the parameters subject to optimization, thereby bolstering the model's robustness against distributional changes (Schneider et al., 2020).

Nonetheless, these techniques often face constraints, either due to the necessity of a substantial number of test samples for generating non-trivial solutions or due to limitations in the scalability of the model architecture. Later research focused on leveraging large-scale pre-trained models with parameter-efficient fine-tuning (Gao et al., 2022; Zhang et al., 2022b). For instance, TPT developed target-specific text prompts while keeping the main model parameters fixed. During testing, it generated various randomly augmented views and eliminated noisy augmentations that could result in inaccurate predictions by minimizing entropy (Shu et al., 2022). Nevertheless, entropy-based confidence selection (Shu et al., 2022) faces a limitation in effectively filtering out misclassified augmented samples that yield low-entropy

predictions. TDA adopted a training-free dynamic adapter to enable effective test time adaptation with vision-language models (Karmanov et al., 2024). However, this method necessitates retaining each sample in the testing data stream, which is not ideal for test-time training. To address this issue, we propose a filtration method based on cosine similarity between the logits of the original test data and the augmented images and texts (see Fig. 1). This approach aims to ensure that the augmented samples (including image and text) maintain consistent class semantics (*i.e.*, *prediction fidelity*) while introducing *diverse* information.

2.3 Image and Text Augmentation

Training models with synthetic images are gaining popularity and undergoing rapid development. In contrast to standard data augmentation methods, such as image manipulation (Shorten & Khoshgoftaar, 2019), image erasing (Zhong et al., 2020), and image mixup (Zhang et al., 2020; Hendrycks et al., 2019), image synthesis offers higher flexibility as these methods augment datasets with pre-defined transformations and cannot provide images with highly diverse content. Early image generation methods, including VAEs (Kingma & Welling, 2013) and GANs (Goodfellow et al., 2020), initially provided promising generated images (Brock et al., 2018), and have been widely applied to various vision tasks. In the latest advancements, there has been notable progress in the creation of diffusion models, aimed at producing images of superior quality with enhanced photo-realistic features compared to earlier image generation techniques (Ho et al., 2020; Nichol & Dhariwal, 2021; Saharia et al., 2022; Ramesh et al., 2022; Zhang et al., 2022a). New studies have illustrated the remarkable effectiveness of diffusion generative models in various applications. For instance, utilizing the latent space of powerful pretrained autoencoders has shown success in high-resolution image synthesis (Rombach et al., 2022), improving text-guided image synthesis (Nichol et al., 2021; Ramesh et al., 2022), establishing a diffusion-based prior for few-shot conditional image generation (Sinha et al., 2021), and implementing a probabilistic model for point cloud generation (Ho et al., 2022). The findings from these studies inspire us to enhance test data directly by integrating *diverse information while maintaining consistent semantics* using the diffusion model and enhancing the performance of test-time optimization.

Another line of approaches to data augmentation involves text augmentation. Large-scale pre-trained models not only generate images with richer visual appearance variations but also provide diverse textual representations through models like GPT and other large language models (LLMs). For example, recent research has utilized GPT-3.5 and GPT-4 to enhance data and has contrasted these approaches with

conventional cutting-edge methodologies for NLP augmentation (Piedboeuf & Langlais, 2023; Ubani et al., 2023). The findings from these studies indicated that the application of large language models (LLMs) as data amplifiers surpassed older methodologies. This was evident particularly in tasks like rephrasing existing texts (Dai et al., 2023) and creating new textual content in zero-shot scenarios using specific cues (Ubani et al., 2023). The *diverse textual representation generation* capabilities of LLMs also offer new perspectives for test-time optimization.

3 Methodology

3.1 Approach Overview

Although the augmentation technique presented in the study by Shu et al. (2022) has showcased notable successes in the realm of TPT, the effectiveness of this approach prominently relies on the range of diversity exhibited in the augmented images. Given that augmented perspectives frequently exhibit akin object and background visual compositions as the original test dataset, the model confronts overly simplistic alterations within the test ensemble, potentially instigating prompt overfitting. Furthermore, Shu et al. (2022) implement an entropy-driven confidence selection mechanism to discard augmented views displaying high entropy predictions. Essentially, the bulk of retained augmented visuals depict cropped variants of objects sourced from the initial test image (refer to Fig. 1a, augmented views). Consequently, the augmentation techniques outlined in Shu et al. (2022) give rise to subtle modifications in the augmented visuals, ultimately curbing the adaptability of learned textual prompts (Bansal & Grover, 2023).

In our conference version DiffTPT (Feng et al., 2023b), we employ a diffusion model on each test sample to generate diverse novel images, thereby capturing natural variations in appearance while retaining key semantics, effectively circumventing this issue. Consequently, diffusion-based data augmentation not only increases the quantity of original test samples but also maintains semantic consistency amidst distribution shifts. Furthermore, DiffTPT introduces cosine similarity-based filtering to eliminate potentially false enhancements that stable diffusion may introduce, preventing erroneous predictions.

However, we noticed that the diversity of a single modality is limited, even with powerful visual generative models. As such, we propose augmenting the original text template “a photo of a” while enhancing the visual features, refer to Fig. 2a. To fully utilize the diverse augmentations from generative models in both images and text, we have replaced prompt tuning in DiffTPT with an adapter. Then, we have shifted the cosine similarity computation from the

original test images to the logits of the original test data with augmented images and texts, thereby filtering out low-quality generated data, *i.e.*, augmented images and text, (refer to Fig. 2a). Next, we will delve into the details of test time optimization based on adapters, data augmentation for images/text, and cosine similarity-based filtering across multi-modalities.

3.2 Test-Time Adaptation

Pre-trained vision-language models such as CLIP are structured with dual encoders, including the image encoder $f(\cdot)$ and the text encoder $g(\cdot)$, offering a wealth of information for diverse downstream applications. In zero-shot classification scenarios, the prediction probabilities can be acquired by

$$p(y_i | \mathbf{x}) = \frac{\exp(\cos(\mathbf{w}_i, \mathbf{e}) / \tau)}{\sum_{j=1}^K \exp(\cos(\mathbf{w}_j, \mathbf{e}) / \tau)}, \quad (1)$$

where the image features denoted as \mathbf{e} generated by $f(\cdot)$ for the image \mathbf{x} , which in conjunction with the corresponding text feature \mathbf{w}_i are employed to calculate the cosine similarity $\cos(\mathbf{w}_i, \mathbf{e})$ pertaining to class i out of K classes. Furthermore, τ represents the temperature parameter. Different from the TPT (Shu et al., 2022) and DiffTPT (Feng et al., 2023b) which learn prompt embeddings, we alternatively tune the adapter to enable the use multiple augmented text templates other than initializing with single text template “a photo of a”. As such, we used CLIP-Adapter (Gao et al., 2024) *Adp* for better alignment of the augmented image-text pairs proposed by generative models, hence we have

$$p(y_i | \mathbf{x}) = \frac{\exp(\cos(\mathbf{w}_i, \mathbf{e}_{\text{adp}}) / \tau)}{\sum_{j=1}^K \exp(\cos(\mathbf{w}_j, \mathbf{e}_{\text{adp}}) / \tau)}, \quad (2)$$

where $\mathbf{e}_{\text{adp}} = \beta * \text{Adp}(\mathbf{e}) + (1 - \beta) * \mathbf{e}$, $\text{Adp}(\mathbf{e})$ denoting the adapted features generated by CLIP-Adapter (Gao et al., 2024) module. β is the coefficient in CLIP-Adapter (Gao et al., 2024) for mixing original and adapted features.

However, as basic models typically need to generalize to out-of-distribution samples, improvements are needed in the zero-shot generalization performance of CLIP. TPT (Shu et al., 2022) and DiffTPT (Feng et al., 2023b) both consider prompt tuning at test time, as it allows for modifying the context of class names to adapt to new test data samples. However, given that different modalities can offer a more diverse knowledge, our proposed IT³A involves multiple text templates, breaking the restriction of using one template “a photo of a” only. Here, we need to optimize the adapter based on a single test sample $\mathbf{x}_{\text{test}} \in \mathbb{R}^{C \times H \times W}$ during the testing phase. Formally, we have

$$\mathbf{a}^* = \arg \min_{\mathbf{a}} \mathcal{L}(\mathcal{F}, \mathbf{a}, \mathbf{x}_{\text{test}}), \quad (3)$$

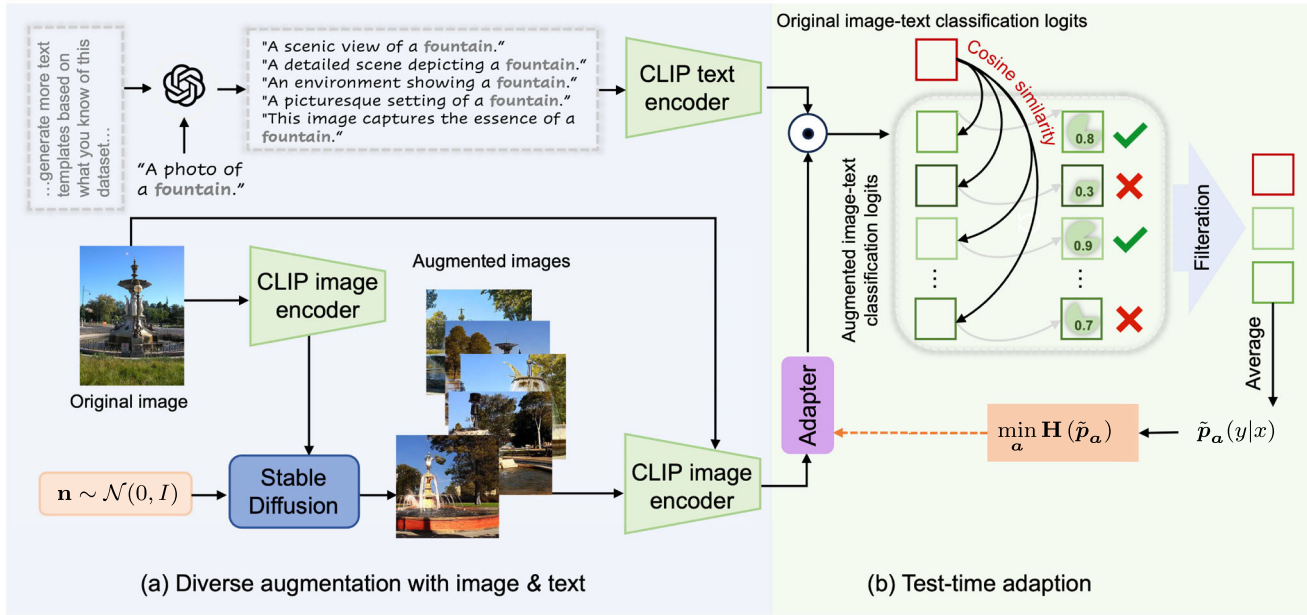


Fig. 2 Overview of our proposed IT^3A . First, **a** we utilize pre-trained generative models, *i.e.*, diffusion and GPT-4, to generate images and text data with *richer visual appearance variations and styles*. These are then randomly combined into different image-text pairs. Then, **b** we apply cosine similarity-based filtering on the classification logits of the

augmented image text pairs generated for a single test sample against their corresponding real test sample. This helps *remove spurious augmentations and inaccurate image-text pairs*, allowing our method to *balance diversity and fidelity*

where \mathcal{F} is the CLIP model consist of an image encoder $f(\cdot)$ and a text encoder $g(\cdot)$, \mathbf{a}^* denotes learnable adapter parameters of the CLIP-Adapter module. \mathcal{L} indicates the cross-entropy loss in the classification task.

To ensure the efficacy of adaptation during testing, we employ a combination of various augmented pairs, *i.e.*, image and text (*i.e.*, $N * M$ in total), alongside a mechanism for selecting different confidence. This can be formulated as:

$$\mathbf{a}^* = \arg \min_{\mathbf{a}} - \sum_{i=1}^K \tilde{p}_{\mathbf{a}}(y_i | \mathbf{x}_{\text{test}}) \log \tilde{p}_{\mathbf{a}}(y_i | \mathbf{x}_{\text{test}}), \quad (4)$$

$$\tilde{p}_{\mathbf{a}}(y_i | \mathbf{x}_{\text{test}}) = \frac{1}{\rho_H N M} \sum_{n,m=1}^{N,M} \mathbb{1}[\mathbf{H}_{n,m}] p_{\mathbf{a}}(y_i | \mathcal{M}(\cdot)), \quad (5)$$

where K indicates the number of classes. $\mathbf{H}_{n,m}$ is a mask representing $p_{\mathbf{a}}(y_i | \mathcal{M}(\cdot)) \leq \tau$ for filtering predictions by selection confidence levels in terms of self-entropy. $p_{\mathbf{a}}(y_i | \mathcal{M}(\cdot))$ represents the class probabilities for the n, m -th augmented pair $\mathcal{M}(\cdot)$ of original pair $(\mathbf{x}_{\text{test}}, \mathbf{t}_{\text{test}})$ from both the vision and language generative models with adapter model parameters \mathbf{a} . \mathbf{t}_{test} is the text description for image \mathbf{x}_{test} and category with label y_i . The threshold τ determines the selection of confidence levels that lead to a ρ_H fraction of all $N * M$ augmented pairs.

3.3 Diverse Data Augmentation

3.3.1 Diffusion-based Diverse Image Augmentation

Diffusion-driven image augmentation aims to create a range of varied and enriching augmented images. As illustrated in Fig. 2a, we begin by extracting latent features z_0 from the pre-trained CLIP encoder $f(\mathbf{x}_{\text{test}})$ of a given test image \mathbf{x}_{test} , followed by employing stable diffusion as the decoder to generate diverse augmented images. Here, we utilize the Stable Diffusion-V2 model as the generative framework, enabling the creation of a novel image $\mathcal{G}(g(\mathbf{t}), \mathbf{n})$ based on textual descriptions \mathbf{t} . $\mathbf{n} \sim \mathcal{N}(0, I)$ indicates the sampled noise. Given the absence of labels during test-time tuning, we opt to substitute $g(\mathbf{t})$ with the image encoder from the CLIP model, denoted as $f(\mathbf{x}_{\text{test}})$. Consequently, the synthetic image can be produced by utilizing

$$\mathcal{D}_n(\mathbf{x}_{\text{test}}) = \mathcal{G}(f(\mathbf{x}_{\text{test}}), \mathbf{n}_n), \quad (6)$$

where n -th augmented image is represented as $\mathcal{D}_n(\mathbf{x}_{\text{test}})$. The alignment capability of CLIP in associating images with text leads to the efficacy of diffusion-driven data augmentation in creating a varied set of augmented images. Figure 3 depicts our diverse and informative augmentations.

By incorporating the augmented images $\mathcal{D}_n(\mathbf{x}_{\text{test}})$, modifications to the paired augmentation $\mathcal{M}(\cdot)$ as outlined in Eq. (5) can be achieved through adaptation:

$$\mathcal{M}(\cdot) = (\mathcal{D}_n(\mathbf{x}_{\text{test}}), \mathcal{T}_m(\mathbf{t}_{\text{test}})) \quad (7)$$

where $\mathcal{T}_m(\mathbf{t}_{\text{test}})$ denotes the m -th augmented text. Please note that the augmented data from different views in TPT (Shu et al., 2022) are both incorporated with diffusion-based ones to take advantage of their complementary merits. IAs a consequence, we can still use Eq. (4) and Eq. (5) for test time adaptation.

3.3.2 LLM-based Diverse Text Augmentation

For text augmentation, we employ generative Large Language Models (LLMs), to obtain diverse text templates. As is shown in Fig. 2a, from a given single text template \mathbf{t}_{test} , we instruct the LLM to generate augmented text templates. Specifically, we employ GPT-4 (Achiam et al., 2023) as the generative language model, which can generate new template $\mathcal{K}(\mathbf{t}_{\text{test}}, \mathbf{s})$, with natural language LLM instruction \mathbf{s} . Thus, the augmented template can then be generated with

$$\mathcal{T}_m(\mathbf{t}_{\text{test}}) = \mathcal{K}(\mathbf{t}_{\text{test}}, \mathbf{s}), \quad (8)$$

where $\mathcal{T}_m(\mathbf{t}_{\text{test}})$ denotes the m -th augmented template. Below is an example of the \mathbf{s} .

"I am using the CLIP model for image classification on the ImageNet dataset. I am currently using the text template 'a photo of a {classname}'. Please generate 8 text templates for better classification performance. The classname should be the last word of each text template."

3.4 Filtration with Cosine Similarity

Although data augmentation based on generative models is effective in producing a variety of augmented data, it may introduce some spurious augmentations (see Fig. 3), resulting in low data fidelity and degraded performance during learning. These false augmentations stem partly from spurious augmentation generated by the diffusion model itself and partly from inadequate combinations of text and generated images (see the blue box in Fig. 3). Therefore, it is essential to balance the diversity of augmented data with the fidelity of predictions.

To achieve this, we use cosine similarity-based filtering to remove the above false augmentations, *i.e.*, low-quality images generated by the diffusion model, and inadequate combinations of text and generated images. In specific, we calculate the cosine similarity between the classification logits of the test sample pair $(\mathbf{x}_{\text{test}}, \mathbf{t}_{\text{test}})$ and each augmented

image-text pair $\mathcal{M}(\cdot) = (\mathcal{D}_n(\mathbf{x}_{\text{test}}), \mathcal{T}_m(\mathbf{t}_{\text{test}}))$. We then introduce a mask \mathcal{C} to identify augmented data with a similarity exceeding ε . Formally, $\mathcal{C}_{n,m} = \cos(\mathbf{l}_0, \mathbf{l}_{n,m}) > \varepsilon$, where \mathbf{l}_0 and $\mathbf{l}_{n,m}$ represent the classification logits of the test sample pair $(\mathbf{x}_{\text{test}}, \mathbf{t}_{\text{test}})$ and each augmented image-text pair $\mathcal{M}(\cdot)$, respectively. We note that ε is the threshold parameter that leads to a ρ_C percentage of the augmented image-text pairs. Formally, the test-time adaptation in Eq. (5) can be further modified as

$$\tilde{p}_a(y_i | \mathbf{x}_{\text{test}}) = \frac{1}{Z} \sum_{n,m=1}^{N,M} \mathbb{1}[\mathbf{H}_{n,m}] \cdot \mathbb{1}[\mathcal{C}_{n,m}] p_a(y_i | \mathcal{M}(\cdot)), \quad (9)$$







where $\frac{1}{Z} = \frac{1}{\rho_H \rho_C N M}$. As a result, we can generate a substantial number of augmented samples with greater diverse data, while retaining essential semantics to refine the adapter during test-time.

4 Experiments

4.1 Experimental Setup

Implementation Details. Our experiments are conducted with 32GB NVIDIA Tesla V100 GPUs and 40GB NVIDIA A100 GPUs, each run requiring one GPU. For CLIP-Adapter (Gao et al., 2024), the initial weights are randomly initialized, and the adapter model is fine-tuned based on a single test image. By default, the dimensionality reduction for the adapter is set to 4. DiffTPT (Shu et al., 2022) enhances each test image to create variations via Stable Diffusion-V2 and through diverse augment views (Shu et al., 2022). The number of variations is set to 7 for both stable diffusion and augment view. For our method, we generate 7 new images and further enhance the image-text pairs with 7 different text templates generated by GPT-4. The adapter undergoes optimization over 4 steps during the test phase using the AdamW optimizer, with the initial learning rate, ρ_H , and ρ_C set to 0.005, 0.3, and 0.8, respectively.

Datasets. We use two Scenarios to evaluate our proposed method, *i.e.*, \mathcal{S}_1 : Natural Distribution Shifts and \mathcal{S}_2 : Cross-Dataset Generalization. For \mathcal{S}_1 , following (Shu et al., 2022), we use four out-of-distribution (OOD) datasets including **ImageNet** (Deng et al., 2009), *i.e.*, **ImageNet-V2** (Recht et al., 2019), **ImageNet-A** (Hendrycks et al., 2021b), **ImageNet-R** (Hendrycks et al., 2021a), and **ImageNet-Sketch** (Wang et al., 2019). These datasets vary in image style and data domains, allowing us to evaluate the robustness of our method against natural distribution shifts. For \mathcal{S}_2 , we utilize 10 diverse datasets covering various species of plants and animals, scenes, textures, food,

Dataset / Classname Aircraft. / 707-320		Original pair	Diverse & informative augmented pairs			Spurious augmentation	Inadequate combinations
							
		A photo of a {707-320}.	This photograph captures an aircraft identified as a {707-320}.	An aerospace design of a {707-320}.	A side view of a {707-320}.	A detailed image of a {707-320}.	A flying example of a {707-320}.







Dataset / Classname Stanford cars / Chevrolet express van 2007							
		A photo of a {Chevrolet express van 2007}.	A detailed image of a {Chevrolet express van 2007}.	This photograph showcases a {Chevrolet express van 2007}.	An automotive design of a {Chevrolet express van 2007}.	A stylish photo of a {Chevrolet express van 2007}.	A sporty look at a {Chevrolet express van 2007}.

Fig. 3 Examples of the **original pairs** (orange box) from a single test data include **diverse and informative augmented pairs** (green box), as well as **spurious augmentations** and **inadequate combinations** (blue

box). Spurious augmentations and inadequate combinations are filtered using the cosine similarity of the predicted logits from the image-text pairs

transportation, human actions, satellite images, and general objects: **Flower102** (Nilsback & Zisserman, 2008), **Oxford-Pets** (Parkhi et al., 2012), **SUN397** (Xiao et al., 2010), **DTD** (Cimpoi et al., 2014), **Food101** (Bossard et al., 2014), **StanfordCars** (Krause et al., 2013), **Aircraft** (Maji et al., 2013), **UCF101** (Soomro et al., 2012), **EuroSAT** (Helber et al., 2019), and **Caltech101** (Fei-Fei et al., 2004). To explore cross-dataset generalization, ImageNet serves as the source dataset, while the other 10 datasets are used as target datasets for evaluation. In our experiments, we randomly select 1,000 test images from all classes to evaluate each method.

Baselines. To assess our proposed method, we employ three groups of methodologies: **a)** TPT (Shu et al., 2022), which is a state-of-the-art test-time prompt tuning technique optimized using multiple augmented views, **b)** traditional PEFT methods for CLIP, specifically CoOp (Zhou et al., 2022b) the few-shot prompt tuning baseline that adjusts a fixed prompt for each downstream dataset, CoCoOp (Zhou et al., 2022a) the enhanced few-shot prompt tuning baseline that creates input-conditional prompts via a lightweight neural network, and CLIP-Adapter (Gao et al., 2024), a flexible method which help model adapt to new dataset at feature level, and **c)** zero-shot CLIP, using the default prompt “a photo of a”. Adhering to the procedures of previous works (Zhou et al., 2022b, a; Shu et al., 2022; Gao et al., 2024), all baselines

are trained on ImageNet with 16-shot examples, 4 learnable prompt tokens for CoOp/CoCoOp, 2-layer linear adapter for CLIP-Adapter (Gao et al., 2024) and subsequently tested on OOD benchmarks. In TPT and DiffTPT, default template “a photo of a” and CoOp/CoCoOp pretrained weights are used for initialization. We note that such methods to initialize learnable prompts restrict the use of multiple augment templates. IT³A adopts CLIP-Adapter (Gao et al., 2024) as a more flexible backbone.

4.2 Comparison with State-of-the-Arts

4.2.1 Natural Distribution Shifts

Table 1 provides an overview of the performance assessments for various competitive approaches within Scenario 1, utilizing different backbones such as ResNet-50 and ViT-B/16. In this context, CLIP represents the zero-shot CLIP output with the standard prompt “a photo of a”. “—&CoOp” and “—&CoCoOp” refer to the implementation of test-time prompt tuning techniques on CoOp (Zhou et al., 2022b) or CoCoOp (Zhou et al., 2022a), respectively. These methods are fine-tuned with 16-shot training samples per category on ImageNet. Our proposed IT³A along with CLIP-Adapter and IT³A initialized with few-shot pretrained CLIP-Adapter

Table 1 Top 1 accuracy % of state-of-the-art baselines under Scenario 1 (S_1)

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sk.	Average	OOD Avg.
CLIP-RN50	56.70(<i>bs.</i>)	23.80(<i>bs.</i>)	50.20(<i>bs.</i>)	54.40(<i>bs.</i>)	33.70(<i>bs.</i>)	43.76(<i>bs.</i>)	40.53(<i>bs.</i>)
TPT	56.80(0.10) ↑	23.80(0.00) ↑	50.30(0.10) ↑	54.30(0.10) ↓	33.60(0.10) ↓	43.76(0.00) ↑	40.50(0.03) ↓
DiffTPT	58.00(1.30) ↑	31.40(7.60) ↑	51.80(1.60) ↑	56.50(2.10) ↑	35.80(2.10) ↑	46.70(2.94) ↑	43.88(3.35) ↑
IT³A	59.30 (2.60) ↑	33.90 (10.10) ↑	54.30 (4.10) ↑	59.90 (5.50) ↑	38.90 (5.20) ↑	49.26 (5.50) ↑	46.75 (6.22) ↑
CoOp	62.00(<i>bs.</i>)	25.00(<i>bs.</i>)	54.60(<i>bs.</i>)	54.70(<i>bs.</i>)	36.00(<i>bs.</i>)	46.46(<i>bs.</i>)	42.58(<i>bs.</i>)
TPT&CoOp	62.00(0.00) ↑	25.00(0.00) ↑	54.90(0.30) ↑	54.90(0.20) ↑	36.40(0.40) ↑	46.64(0.18) ↑	42.80(0.22) ↑
DiffTPT&CoOp	<u>63.00</u> (1.00) ↑	<u>33.70</u> (8.70) ↑	<u>55.70</u> (1.10) ↑	<u>57.60</u> (2.90) ↑	34.60(1.40) ↓	<u>48.92</u> (2.46) ↑	<u>45.40</u> (2.82) ↑
CoCoOp	58.20(<i>bs.</i>)	26.50(<i>bs.</i>)	53.10(<i>bs.</i>)	55.90(<i>bs.</i>)	35.90(<i>bs.</i>)	45.92(<i>bs.</i>)	42.85(<i>bs.</i>)
TPT&CoCoOp	58.20(0.00) ↑	26.50(0.00) ↑	53.20(0.10) ↑	55.90(0.00) ↑	35.80(0.10) ↓	45.92(0.00) ↑	42.85(0.00) ↑
DiffTPT&CoCoOp	58.30(0.10) ↑	26.30(0.20) ↓	53.30(0.20) ↑	56.00(0.10) ↑	<u>35.90</u> (0.00) ↑	45.94(0.02) ↑	42.88(0.03) ↑
CLIP-Ap.	60.50(<i>bs.</i>)	24.10(<i>bs.</i>)	53.30(<i>bs.</i>)	54.70(<i>bs.</i>)	34.50(<i>bs.</i>)	45.42(<i>bs.</i>)	41.65(<i>bs.</i>)
IT³A&CLIP-Ap.	62.40 (1.90) ↑	33.80 (9.70) ↑	56.30 (3.00) ↑	60.60 (5.90) ↑	36.60 (2.10) ↑	49.94 (4.52) ↑	46.83 (5.18) ↑
CLIP-ViT-B/16	63.60(<i>bs.</i>)	47.20(<i>bs.</i>)	59.40(<i>bs.</i>)	72.60(<i>bs.</i>)	46.00(<i>bs.</i>)	57.76(<i>bs.</i>)	56.30(<i>bs.</i>)
TPT	63.60(0.00) ↑	47.40(0.20) ↑	59.50(0.10) ↑	72.70(0.10) ↑	45.90(0.10) ↓	57.82(0.06) ↑	56.38(0.08) ↑
DiffTPT	64.80(1.20) ↑	54.50 (7.30) ↑	60.10(0.70) ↑	74.30(1.70) ↑	47.50(1.50) ↑	60.24(2.48) ↑	59.10(2.80) ↑
IT³A	66.00 (2.40) ↑	51.30(4.10) ↑	60.70 (1.30) ↑	76.00 (3.40) ↑	49.00 (3.00) ↑	60.60 (2.84) ↑	59.25 (2.95) ↑
CoOp	68.30(<i>bs.</i>)	48.10(<i>bs.</i>)	61.90(<i>bs.</i>)	70.70(<i>bs.</i>)	45.50(<i>bs.</i>)	58.90(<i>bs.</i>)	56.55(<i>bs.</i>)
TPT&CoOp	68.30(0.00) ↑	48.20(0.10) ↑	62.00(0.10) ↑	70.70(0.00) ↑	45.60(0.10) ↑	58.94(0.04) ↑	56.60(0.05) ↑
DiffTPT&CoOp	<u>69.70</u> (1.40) ↑	<u>53.00</u> (4.90) ↑	62.30(0.40) ↑	72.60(1.90) ↑	46.50(1.00) ↑	<u>60.82</u> (1.92) ↑	<u>58.60</u> (2.05) ↑
CoCoOp	65.90(<i>bs.</i>)	48.90(<i>bs.</i>)	60.90(<i>bs.</i>)	74.50(<i>bs.</i>)	47.80(<i>bs.</i>)	59.60(<i>bs.</i>)	58.03(<i>bs.</i>)
TPT&CoCoOp	65.90(0.00) ↑	48.80(0.10) ↓	60.90(0.00) ↑	74.60 (0.10) ↑	47.80(0.00) ↑	59.60(0.00) ↑	58.03(0.00) ↑
DiffTPT&CoCoOp	66.90(1.00) ↑	48.70(0.20) ↓	61.80(0.90) ↑	74.50(0.00) ↑	<u>49.10</u> (1.30) ↑	60.20(0.60) ↑	58.53(0.50) ↑
CLIP-Ap.	67.40(<i>bs.</i>)	48.10(<i>bs.</i>)	<u>62.50</u> (<i>bs.</i>)	72.90(<i>bs.</i>)	47.10(<i>bs.</i>)	59.60(<i>bs.</i>)	57.65(<i>bs.</i>)
IT³A&CLIP-Ap.	68.60 (1.20) ↑	56.10 (8.00) ↑	63.20 (0.70) ↑	74.60 (1.70) ↑	50.30 (3.20) ↑	62.56 (2.94) ↑	61.05 (3.40) ↑

ImageNet-Sk. denotes the ImageNet-Sketch dataset, while **OOD Avg.** represents the average performance across out-of-distribution datasets. The abbreviation *bs.* signifies the baseline for each group, *i.e.*, CLIP-RN50 / CLIP-ViT-B-16, CoOp, CoCoOp, and CLIP-Adapter. Arrows (↑ and ↓) indicate enhancements and reductions compared to the baseline. For comprehensive analyses, refer to Sec. 4.2

Bold is the best performance and the underline is the second-best performance

weights is also included in the table. We note that our method aims to maintain effective performance with limited augmented images in a resource-efficient manner. Therefore, unlike the conference version of DiffTPT (Feng et al., 2023b), we only use 8-fold augmentation here. As demonstrated in the table, in general IT³A and IT³A&CLIP-Ap. outperform all other methods and their variants correspondingly. On CLIP-RN50, the average performance of IT³A improved by 5.50, and the average for OOD generation was improved by 6.22. Compared to the conference version of DiffTPT, the performances are as follows: DiffTPT: 46.70/43.88 *vs.* IT³A: **49.26/46.75**. On CLIP-Adapter, IT³A also showed significant improvements, with all its performances exceeding those of DiffTPT applied to CoOp (Zhou et al., 2022b) or CoCoOp (Zhou et al., 2022a). Both methods enhance in-domain accuracy on **ImageNet** data and generalization to OOD data. For ResNet-50-based in-domain average performance, DiffTPT&CoCoOp: 45.94 *vs.* IT³A&CLIP-Adapter: **49.94**; for the generalization test of OOD data, DiffTPT&CoCoOp: 42.88 *vs.* IT³A&CLIP-Adapter: **46.83**.

Our method also achieved similar improvements on ViT-B/16. Specifically, compared to DiffTPT: 60.24 → **60.60**, 59.10 → **59.25**, 60.20 → **62.56**, and 58.53 → **61.05**. Since TPT (Shu et al., 2022) use random resized cropping to augment test images, their generalization ability is limited. DiffTPT (Feng et al., 2023b) can only acquire visual diversity from a single image modality, limiting the knowledge it captures. Notably, we found that IT³A significantly improves the generalization test on OOD data. This supports our conclusion that IT³A enhances robustness by acquiring more knowledge through multi-modal augmentation (*i.e.*, **prediction fidelity**) and increasing the **data diversity** of the test samples.

Naturally, the results generated by CLIP are the lowest, as direct testing on new datasets is significantly impacted by domain shifts. Although CoOp (Zhou et al., 2022b) and CoCoOp (Zhou et al., 2022a) benefit from learnable prompts, these methods depend on training datasets and do not utilize prompt tuning at test time. CLIP-Adapter achieves performance gain from learnable adapters but also faces the same

Table 2 Top 1 accuracy % of state-of-the-art baselines under \mathcal{S}_2 .

Method	Flower	DTD	Pets	Cars	UCF101	
CLIP-RN50	61.60 _(bs.) ↓	38.50 _(bs.) ↓	84.70 _(bs.) ↓	55.70 _(bs.) ↓	58.6 _(bs.) ↓	
CoOp ₂₀₂₂ (Zhou et al., 2022b)	60.90	36.60	88.00	54.60	59.01	
CoCoOp ₂₀₂₂ (Zhou et al., 2022a)	63.90	40.70	88.50	53.50	59.60	
CLIP-Ap- ₂₀₂₄ (Gao et al., 2024)	62.90	40.10	84.80	55.10	58.80	
TPT ₂₀₂₂ (Shu et al., 2022)	59.20(2.40) ↓	39.10(0.60) ↑	83.30(1.40) ↓	54.90(0.80) ↓	59.6(1.00) ↑	
DiffTPT ₂₀₂₃ (Feng et al., 2023b)	57.10(4.50) ↓	38.60(0.10) ↑	84.20(0.50) ↓	57.30(1.60) ↑	63.70(5.10) ↑	
IT ³ A	59.40(1.20) ↓	42.30(3.80) ↑	85.40(0.70) ↑	57.30(1.60) ↑	63.40(4.80) ↑	
Method	Caltech11	Food101	SUN397	Aircraft	EuroSAT	Avg.
CLIP-RN50	85.20 _(bs.) ↓	75.90 _(bs.) ↓	60.00 _(bs.) ↓	15.50 _(bs.) ↓	19.70 _(bs.) ↓	55.54 _(bs.) ↓
CoOp ₂₀₂₂ (Zhou et al., 2022b)	86.10	78.20	59.00	16.10	22.80	56.24
CoCoOp ₂₀₂₂ (Zhou et al., 2022a)	87.70	78.50	59.60	15.40	30.50	57.79
CLIP-Ap- ₂₀₂₄ (Gao et al., 2024)	86.10	74.20	60.70	16.60	25.80	56.51
TPT ₂₀₂₂ (Shu et al., 2022)	84.30(0.90) ↓	75.80(0.10) ↓	60.90(0.90) ↑	17.00(1.50) ↑	22.80(3.10) ↑	55.69(0.15) ↑
DiffTPT ₂₀₂₃ (Feng et al., 2023b)	87.30(2.10) ↑	75.90(0.00) ↑	63.10(3.10) ↑	16.50(1.00) ↑	34.50(14.80) ↑	57.82(2.28) ↑
IT ³ A	87.60(2.40) ↑	74.20(1.70) ↓	60.90(0.90) ↑	18.30(2.50) ↑	40.00(21.30) ↑	58.88(3.34) ↑
Method	Flower	DTD	Pets	Cars	UCF101	
CLIP-ViT-B/16	66.50 _(bs.) ↓	41.90 _(bs.) ↓	88.60 _(bs.) ↓	66.80 _(bs.) ↓	63.70 _(bs.) ↓	
CoOp ₂₀₂₂ (Zhou et al., 2022b)	68.10	41.60	89.80	65.30	64.50	
CoCoOp ₂₀₂₂ (Zhou et al., 2022a)	65.70	42.00	90.00	60.80	61.20	
CLIP-Ap- ₂₀₂₄ (Gao et al., 2024)	67.40	43.20	89.20	65.90	64.60	
TPT ₂₀₂₂ (Shu et al., 2022)	66.50(0) ↑	43.10(1.20) ↑	86.80(1.80) ↓	66.50(0.30) ↓	67.80(4.10) ↑	
DiffTPT ₂₀₂₃ (Feng et al., 2023b)	67.20(0.70) ↑	43.50(1.60) ↑	85.90(2.70) ↓	65.90(0.90) ↓	66.50(2.80) ↑	
IT ³ A	69.90(3.40) ↑	44.50(2.60) ↑	88.80(0.20) ↑	66.90(0.10) ↑	69.00(5.30) ↑	
Method	Caltech11	Food101	SUN397	Aircraft	EuroSAT	Avg.
CLIP-ViT-B/16	91.90 _(bs.) ↓	85.40 _(bs.) ↓	64.10 _(bs.) ↓	24.00 _(bs.) ↓	40.60 _(bs.) ↓	63.35 _(bs.) ↓
CoOp ₂₀₂₂ (Zhou et al., 2022b)	91.80	83.80	64.60	17.60	32.00	61.91
CoCoOp ₂₀₂₂ (Zhou et al., 2022a)	90.80	85.50	64.00	16.00	44.80	62.08
CLIP-Ap- ₂₀₂₄ (Gao et al., 2024)	92.30	84.80	65.30	23.10	39.40	63.50
TPT ₂₀₂₂ (Shu et al., 2022)	91.50(0.40) ↓	86.20(0.80) ↑	66.20(2.10) ↑	21.20(2.80) ↓	37.00(3.60) ↓	63.28(0.07) ↓
DiffTPT ₂₀₂₃ (Feng et al., 2023b)	94.00(2.10) ↑	84.40(1.00) ↓	67.30(3.20) ↑	20.50(3.50) ↓	41.60(1.00) ↑	63.68(0.33) ↑
IT ³ A	93.80(1.90) ↑	84.50(0.90) ↓	68.80(4.70) ↑	25.40(1.40) ↑	39.70(0.90) ↓	65.13(1.78) ↑

Avg. represents the average performance of the Cross-Dataset Generalization. Arrows (↑ and ↓) indicate enhancements and reductions compared to the CLIP method, *i.e.*, CLIP-RN50 and CLIP-ViT-B/16. For comprehensive analyses, refer to Sec. 4.2

Bold is the best performance and the underline is the second-best performance

problem. This means said methods fail to consider zero-shot generalization in practical, real-world settings, which results in lower effective performance. Our findings confirm the initial hypothesis that enhancing test data with diverse synthetic data can boost zero-shot generalization performance.

4.2.2 Cross-Dataset Generalization

We evaluate the ability of our proposed method and several baseline models to generalize from ImageNet to 10 different fine-grained datasets by recording their quantitative perfor-

mances, as shown in Table 2. TPT (Shu et al., 2022) is deployed in a zero-shot manner, CoOp (Zhou et al., 2022b), CoCoOp (Zhou et al., 2022a), and CLIP-Adapter (Gao et al., 2024) are fine-tuned on ImageNet with 16-shot samples per category. Due to the diverse nature of these fine-grained datasets, each method displays varying performance levels on each dataset. However, our proposed IT³A still achieves the best performance, *i.e.*, raising the **Avg.** accuracy from 55.54 to **58.88** on CLIP-RN50, and from 63.03 to **65.13** on CLIP-ViT-B/16, all based on just 8-fold multi-modal augmentation. Notably, our method achieved performance gains of 1.06%

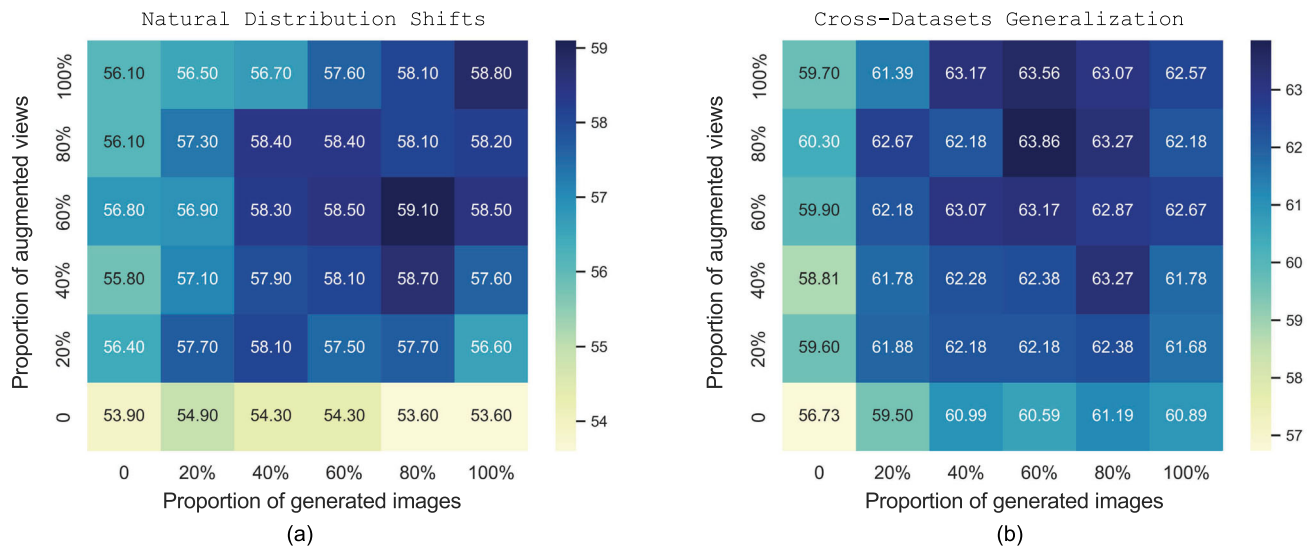


Fig. 4 Top-1 accuracy variation against different proportions of standard augmented views and diffusion-augmented images for scenarios (a) S_1 and (b) S_2

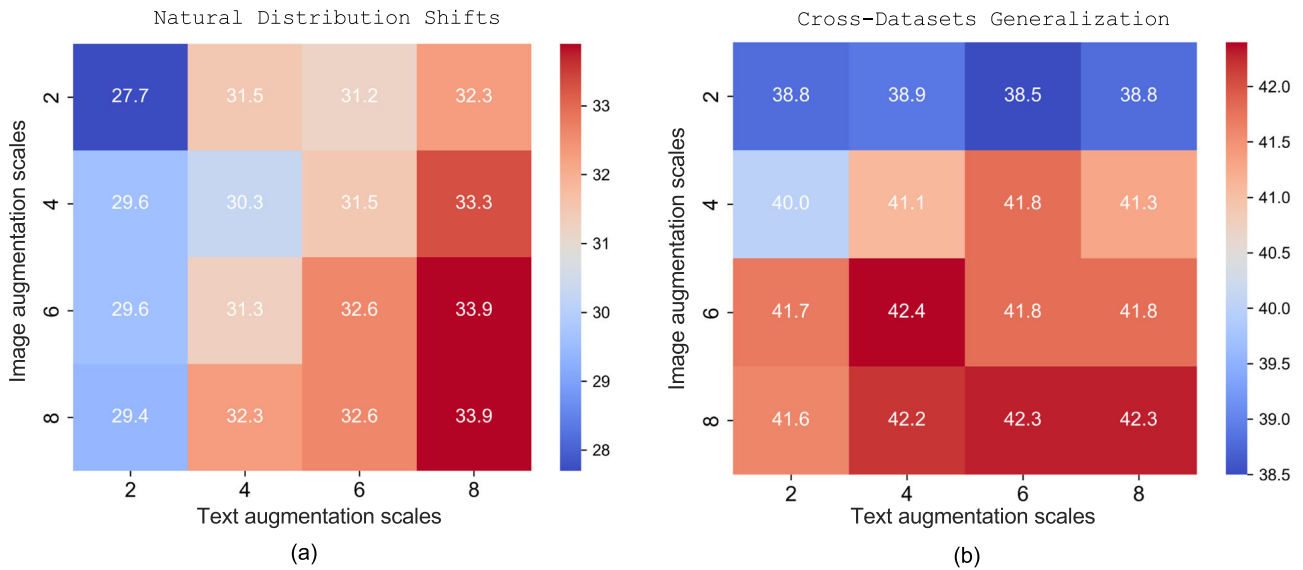


Fig. 5 Top-1 accuracy variation against different scales of text augmentation and image augmentation for **a** ImageNet-A under scenario S_1 and **b** DTD under scenario S_2

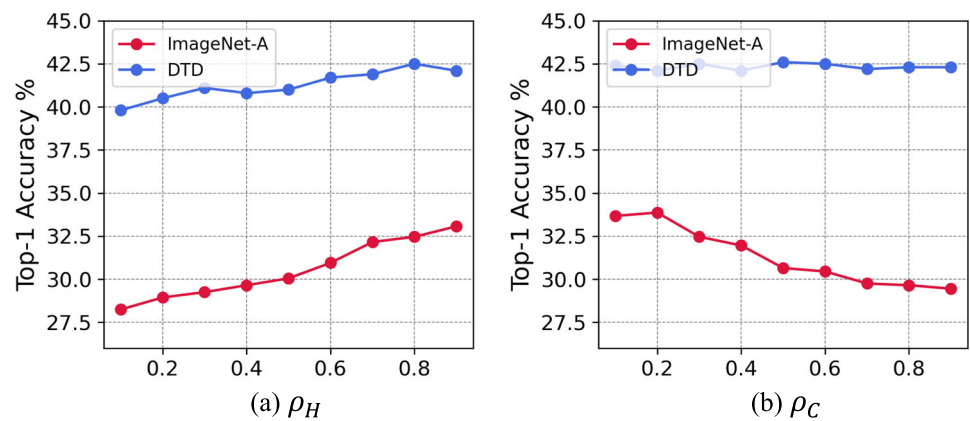
and 1.45% over DiffTPT (Shu et al., 2022) on ResNet-50 and ViT-B/16 backbones, respectively. This demonstrates that, among all competing methods, our approach is robust to natural distribution shifts even without training data, and significantly outperforms few-shot prompt tuning methods, *i.e.*, CoOp (Zhou et al., 2022b), CoCoOp (Zhou et al., 2022a), and CLIP-Adapter (Gao et al., 2024). Although TPT, DiffTPT, and our proposed IT³A exhibits some performance decline on a few datasets, this is primarily due to the limited augmentation used during the testing phase, *i.e.*, 8-fold, while in the previous version of TPT and DiffTPT, they augment the test with 64-fold.

4.3 Ablation Studies

4.3.1 Balancing Synthetic Data vs. Standard Augmentation

Given that our method for image domain augmentation leverages the complementary advantages of both the standard augmentation (Shu et al., 2022) and diffusion-based augmentation, it is essential to investigate how these two methods contribute to training the classifier. To address this, we assess the average performance of ResNet-50 across the two scenarios, *i.e.*, S_1 on ImageNet-R and S_2 on UCF101, inherit the conference version. For improved visualization,

Fig. 6 Top-1 accuracy analysis of the ratios ρ_H and ρ_C with regard to \mathcal{S}_1 and \mathcal{S}_2



we present a plot in Fig. 4 that illustrates mixed combinations of different ratios, with the x-axis representing the percentage of synthetic data generated through diffusion-based augmentation and the y-axis indicating the percentage of data obtained through standard augmentation. In the matrix of this figure, each cell \mathcal{V}_{ij} corresponds to the classification performance of DiffTPT using $i\%$ of synthetic data and $j\%$ of standard augmented data. Figure 4a demonstrates a significant improvement in accuracy for Natural Distribution Shifts as the quantity of standard augmented data increases while keeping the synthetic data level fixed. In contrast, the effects are even more pronounced when the proportion of synthetic data is increased while the amount of standard augmented data is held constant. Overall, increasing the amount of synthetic data leads to better performance in \mathcal{S}_1 . In Fig. 4b, we show the performance of the classifier for \mathcal{S}_2 , *i.e.*, Cross-Dataset Generalization. We observe that, while keeping the amount of synthetic data fixed, the effectiveness of the classifier increases significantly as the proportion of standard augmented data increases.

4.3.2 Analysis of Ratio ρ_H and ρ_C

As mentioned in Sec. 3.4, ρ_H and ρ_C filter out the less informative “noisy” and spurious augmented pairs in overall generative augmentation by standard of self-entropy and cosine similarity. We evaluated the accuracy for various values of ρ_H and ρ_C in Fig. 6 across two scenarios to determine the amount of information that good test augmentations should retain. From Fig. 6a, it can be observed that on the DTD dataset under \mathcal{S}_2 , there is a trade-off between augmented data quantity and augmented data quality, with the highest accuracy achieved at a value of 0.8. Compared to \mathcal{S}_2 , \mathcal{S}_1 requires more data for extended amount of learning to bridge the gap between in-domain ImageNet distribution and OOD adversarial samples from ImageNet-A. Accordingly, the performance of IT³A improves with an increase in ρ_H on the ImageNet-A dataset. For ρ_C in Fig. 6b, larger values correspond to more data pairs, while smaller values

indicate higher data quality. The Fig. 6b illustrates a trade-off between the number of augmented data pairs and data quality.

4.3.3 Effect of the Generated Dataset Size

Since the primary contribution of our method lies in integrating multi-modal augmentation information, we investigate the impact of these two modalities, *i.e.*, image and text, on classifier training. In Fig. 5, we present mixed combinations of different modalities and augmentation scales, where the x-axis represents the scale of text augmentation and the y-axis represents the scale of image augmentation. In the matrix of this figure, each element \mathcal{V}_{ij} indicates the classification performance of IT³A with $i \times$ text augmentation and $j \times$ image augmentation. As illustrated in Fig. 5a, in \mathcal{S}_1 , the augmented textual information provides essential components for domain adaptation. Consequently, as the augmentation levels for both text and images increase, the model performance improves. In contrast, for \mathcal{S}_2 , the model relies less on textual data, resulting in reduced sensitivity to changes in text augmentation levels.

4.3.4 Steps of Prompt Updating

To evaluate the effectiveness of the learning updates, we have recorded the accuracy across different optimization steps for two scenarios in Fig. 7. As illustrated in the figure, the performance of IT³A on the DTD dataset under \mathcal{S}_2 continually improves from 27.3 to 33.8 with an increasing number of optimization steps. In contrast, on the ImageNet-A dataset under \mathcal{S}_1 , the performance of IT³A increases to 42.5 and then stabilizes as the number of optimization steps increases. This indicates that additional optimization steps do not provide further benefits to the classifier and only serve to increase inference time. Taking both performance and computational efficiency into account, we set the number of steps to 4 in our experiments.

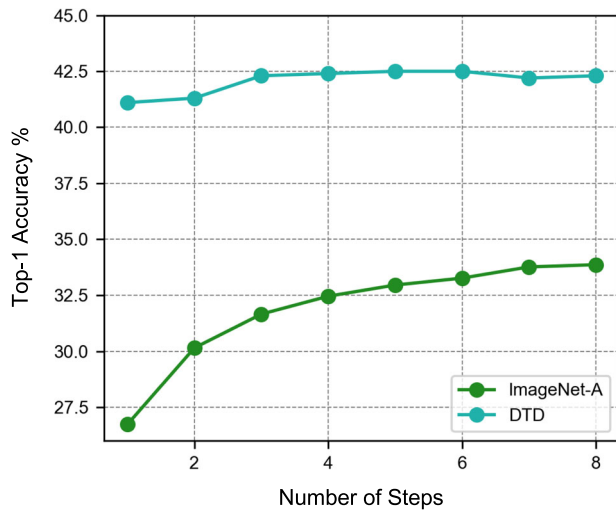


Fig. 7 Ablation studies on the learning steps under S_1 and S_2

Table 3 Investigation of the different configurations of the adapter under both S_1 and S_2

Datasets	1 layer	2 layers	4 layers	8 layers
ImageNet-A	32.60	33.90	32.90	32.00
DTD	42.00	42.30	42.70	42.00

4.3.5 Effect of the Configurations of the Adapter

Here, we investigate whether the different configurations will influence the model performance. Table 3 summarized the accuracies with different layers in the adapter on two scenarios, *i.e.*, ImageNet-A under S_1 and DTD under S_2 . As can be seen in the table, the performance for ImageNet-A dataset increases as the layers of adapter module increase from 1 layer to 2 layers, *i.e.*, from 32.60 to 33.90, but then decreases to 32.00 eventually as layers keep increasing to 8. Similarly, there is also a peak in performance in the middle of the layer range, *i.e.*, 42.70 at 4 layers. According to the results of different configurations of the adapter in this table, we set the layers to 2 in our method.

4.3.6 Inference Cost

Despite the fact that the original SD can be time-intensive, such as taking 6 seconds to infer 10 test images for TPT and 36 minutes with standard SD, recent advancements have led to the development of faster SD models. For example, ToMe (Bolya & Hoffman, 2023), two-stage distillation (Meng et al., 2023), and Consistency Model (Song et al., 2023). Notably, the Consistency Model can generate 10 images in just 0.5 seconds, compared to the original SD's 70 seconds. Moreover, efficiency can be further enhanced using

techniques like TensorRT and Memory Efficient Attention,¹ resulting in additional gains of 25% and 100% in inference speed, respectively. Compared to DiffTPT, IT³A only needs to generate $\frac{1}{M}$ the number of images compared to DiffTPT for the same number of image-text pairs. Also, the computational cost of generating multiple distinct text templates through instructions for GPT-4 is negligible. Therefore, for approximately the same amount of computational consumption, IT³A can generate M times of augment pairs of DiffTPT. In other words, the overall computational cost of IT³A is much lower than that of DiffTPT when comparing under the same amount of augment pairs. Additionally, compared to the prompt-learning method used in the conference version, DiffTPT (Feng et al., 2023b), IT³A requires less tuning time for a single test image, *i.e.*, IT³A: 0.33s vs. DiffTPT: 1.08s.

5 Conclusion

This paper proposes a multi-modal test-time optimization method that leverages enhanced data from pre-trained models in both image and text modalities. By combining the strengths of these modalities, the method improves the model's adaptability to unknown test data. To fully utilize the diversity provided by generative models in both vision and language, we have replaced prompt tuning from the conference version, DiffTPT, with adapters. This change allows for more flexible use of text templates on the text encoder. Additionally, using cosine similarity filtering between the original test data and the augmented images and text ensures that key semantics are faithfully preserved during various visual and textual augmentations. Experiments on test datasets with distribution shifts and unseen classes demonstrate that the IT³A method improves zero-shot accuracy by an average of 4.98% compared to the state-of-the-art TPT method. Our approach of multi-modal augmentation during test time can inspire developments in test time strategies for other multi-modal tasks, such as composed image retrieval, which is also a future direction we are currently exploring.

Acknowledgements This work was supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG2-TC-2021-003), Agency for Science, Technology and Research (A*STAR) through its AME Programmatic Funding Scheme Under Project A20H4b0141, A*STAR Central Research Fund "A Secure and Privacy Preserving AI Platform for Digital Health", and Agency for Science, Technology and Research (A*STAR) through its RIE2020 Health and Biomedical Sciences (HBMS) Industry Alignment Fund Pre-Positioning (IAF-PP) (grant no. H20C6a0032).

Data availability The authors declare that the data supporting the experiments in this study are available within the paper. The code is available at <https://github.com/chunmeifeng/DiffTPT>.

¹ <https://www.photoroom.com/tech/stable-diffusion-100-percent-faster-with-memory-efficient-attention>

References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. (2023). Gpt-4 technical report. arXiv preprint [arXiv:2303.08774](https://arxiv.org/abs/2303.08774)
- Antoniou, A., Storkey, A., & Edwards, H. (2017). Data augmentation generative adversarial networks. arXiv preprint [arXiv:1711.04340](https://arxiv.org/abs/1711.04340)
- Bansal, H., & Grover, A. (2023). Leaving reality to imagination: Robust classification via generated datasets. arXiv preprint [arXiv:2302.02503](https://arxiv.org/abs/2302.02503)
- BELLEGroup. (2023). Belle: Be everyone's large language model engine. <https://github.com/LianjiaTech/BELLE>.
- Bolya, D., & Hoffman, J. (2023). Token merging for fast stable diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 4598–4602.
- Bossard, L., Guillaumin, M., & Van Gool, L. (2014). Food-101—mining discriminative components with random forests. In: *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland*, September 6–12, 2014, Proceedings, Part VI 13, Springer, pp 446–461.
- Brock, A., Donahue, J., & Simonyan, K. (2018). Large scale gan training for high fidelity natural image synthesis. arXiv preprint [arXiv:1809.11096](https://arxiv.org/abs/1809.11096).
- Chen, D., Wang, D., Darrell, T., & Ebrahimi, S. (2022). Contrastive test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 295–305.
- Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, PMLR, pp 1597–1607.
- Chiang, W.L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J.E., et al. (2023). Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023) 2(3):6.
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., & Vedaldi, A. (2014). Describing textures in the wild. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 3606–3613.
- Dai, H., Liu, Z., Liao, W., Huang, X., Cao, Y., Wu, Z., Zhao, L., Xu, S., Liu, W., Liu, N., et al. (2023). Auggpt: Leveraging chatgpt for text data augmentation. arXiv preprint [arXiv:2302.13007](https://arxiv.org/abs/2302.13007).
- Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, IEEE, 248–255.
- Devlin, J., Chang, M.W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint [arXiv:1810.04805](https://arxiv.org/abs/1810.04805).
- Fei-Fei, L., Fergus, R., & Perona, P. (2004). Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *2004 conference on computer vision and pattern recognition workshop*, IEEE, pp 178–178.
- Feng, C.M., Li, B., Xu, X., Liu, Y., Fu, H., & Zuo, W. (2023a). Learning federated visual prompt in null space for mri reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 8064–8073.
- Feng, C.M., Yu, K., Liu, Y., Khan, S., & Zuo, W. (2023b). Diverse data augmentation with diffusions for effective test-time prompt tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp 2704–2714.
- Gao, P., Geng, S., Zhang, R., Ma, T., Fang, R., Zhang, Y., Li, H., & Qiao, Y. (2021). Clip-adapter: Better vision-language models with feature adapters. arXiv preprint [arXiv:2110.04544](https://arxiv.org/abs/2110.04544).
- Gao, P., Geng, S., Zhang, R., Ma, T., Fang, R., Zhang, Y., Li, H., & Qiao, Y. (2024). Clip-adapter: Better vision-language models with feature adapters. *International Journal of Computer Vision*, 132(2), 581–595.
- Gao, Y., Shi, X., Zhu, Y., Wang, H., Tang, Z., Zhou, X., Li, M., & Metaxas, D.N. (2022). Visual prompt tuning for test-time domain adaptation. arXiv preprint [arXiv:2210.04831](https://arxiv.org/abs/2210.04831).
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2020). Generative adversarial networks. *Communications of the ACM*, 63(11), 139–144.
- Helber, P., Bischke, B., Dengel, A., & Borth, D. (2019). Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7), 2217–2226.
- Hendrycks, D., Mu, N., Cubuk, E.D., Zoph, B., Gilmer, J., & Lakshminarayanan, B. (2019). Augmix: A simple data processing method to improve robustness and uncertainty. arXiv preprint [arXiv:1912.02781](https://arxiv.org/abs/1912.02781).
- Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F., Dorundo, E., Desai, R., Zhu, T., Parajuli, S., & Guo, M., et al. (2021a). The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp 8340–8349.
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., & Song, D. (2021b). Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 15262–15271.
- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 6840–6851.
- Ho, J., Chan, W., Saharia, C., Whang, J., Gao, R., Gritsenko, A., Kingma, D.P., Poole, B., Norouzi, M., Fleet, D.J., et al. (2022). Imagen video: High definition video generation with diffusion models. arXiv preprint [arXiv:2210.02303](https://arxiv.org/abs/2210.02303).
- Huang, T., Chu, J., & Wei, F. (2022). Unsupervised prompt learning for vision-language models. arXiv preprint [arXiv:2204.03649](https://arxiv.org/abs/2204.03649).
- Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., & Duerig, T. (2021). Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, PMLR, pp 4904–4916.
- Jia, M., Tang, L., Chen, B.C., Cardie, C., Belongie, S., Hariharan, B., & Lim, S.N. (2022). Visual prompt tuning. In: *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, Springer, pp 709–727.
- Karmanov, A., Guan, D., Lu, S., El Saddik, A., & Xing, E. (2024). Efficient test-time adaptation of vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 14162–14171.
- Kingma, D.P., & Welling, M. (2013). Auto-encoding variational bayes. arXiv preprint [arXiv:1312.6114](https://arxiv.org/abs/1312.6114)
- Krause, J., Stark, M., Deng, J., & Fei-Fei, L. (2013). 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pp 554–561.
- Le Scao, T., Fan, A., Akiki, C., Pavlick, E., Ilić, S., Hesslow, D., Castagné, R., Luccioni, A.S., Yvon, F., Gallé, M., et al. (2023). Bloom: A 176b-parameter open-access multilingual language model. arXiv preprint [arXiv:2211.05100](https://arxiv.org/abs/2211.05100).
- Li, H., Feng, C.M., Zhou, T., Xu, Y., & Chang, X. (2022a). Prompt-driven efficient open-set semi-supervised learning. arXiv preprint [arXiv:2209.14205](https://arxiv.org/abs/2209.14205)
- Li, L.H., Zhang, P., Zhang, H., Yang, J., Li, C., Zhong, Y., Wang, L., Yuan, L., Zhang, L., Hwang, J.N., et al. (2022b). Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 10965–10975.

- Liu, Y., Kothari, P., Van Delft, B., Bellot-Gurlet, B., Mordan, T., & Alahi, A. (2021). Ttt++: When does self-supervised test-time training fail or thrive? *Advances in Neural Information Processing Systems*, 34, 21808–21820.
- Maji, S., Rahtu, E., Kannala, J., Blaschko, M., & Vedaldi, A. (2013). Fine-grained visual classification of aircraft. arXiv preprint [arXiv:1306.5151](https://arxiv.org/abs/1306.5151).
- Mandal, D., Narayan, S., Dwivedi, S.K., Gupta, V., Ahmed, S., Khan, F.S., & Shao, L. (2019). Out-of-distribution detection for generalized zero-shot action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 9985–9993.
- Meng, C., Rombach, R., Gao, R., Kingma, D., Ermon, S., Ho, J., & Salimans, T. (2023). On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 14297–14306.
- Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., & Chen, M. (2021). Glide: Towards photorealistic image generation and editing with text-guided diffusion models. arXiv preprint [arXiv:2112.10741](https://arxiv.org/abs/2112.10741).
- Nichol, A.Q., & Dhariwal, P. (2021). Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, PMLR, pp 8162–8171.
- Nilsback, M. E., & Zisserman, A. (2008). Automated flower classification over a large number of classes. *2008 Sixth Indian Conference on Computer Vision (pp. 722–729)*. Graphics & Image Processing: IEEE.
- Parkhi, O.M., Vedaldi, A., Zisserman, A., & Jawahar, C. (2012). Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, IEEE, pp 3498–3505.
- Perez, L., & Wang, J. (2017). The effectiveness of data augmentation in image classification using deep learning. arXiv preprint [arXiv:1712.04621](https://arxiv.org/abs/1712.04621).
- Piedboeuf, F., & Langlais, P. (2023). Is chatgpt the ultimate data augmentation algorithm? *Findings of the Association for Computational Linguistics: EMNLP, 2023*, 15606–15615.
- Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. (2018). Improving language understanding by generative pre-training.
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. (2021). Learning transferable visual models from natural language supervision. In *International conference on machine learning*, PMLR, pp 8748–8763.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., & Chen, M. (2022). Hierarchical text-conditional image generation with clip latents. arXiv preprint [arXiv:2204.06125](https://arxiv.org/abs/2204.06125).
- Recht, B., Roelofs, R., Schmidt, L., & Shankar, V. (2019). Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, PMLR, pp 5389–5400.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 10684–10695.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S.K.S., Ayan, B.K., Mahdavi, S.S., Lopes, R.G., et al. (2022). Photorealistic text-to-image diffusion models with deep language understanding. arXiv preprint [arXiv:2205.11487](https://arxiv.org/abs/2205.11487).
- Schneider, S., Rusak, E., Eck, L., Bringham, O., Brendel, W., & Bethge, M. (2020). Improving robustness against common corruptions by covariate shift adaptation. *Advances in Neural Information Processing Systems*, 33, 11539–11551.
- Shanmugam, D., Blalock, D., Balakrishnan, G., & Guttat, J. (2021). Better aggregation in test-time augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp 1214–1223.
- Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of big data*, 6(1), 1–48.
- Shu, M., Nie, W., Huang, D.A., Yu, Z., Goldstein, T., Anandkumar, A., & Xiao, C. (2022). Test-time prompt tuning for zero-shot generalization in vision-language models. arXiv preprint [arXiv:2209.07511](https://arxiv.org/abs/2209.07511).
- Sinha, A., Song, J., Meng, C., & Ermon, S. (2021). D2c: Diffusion-decoding models for few-shot conditional generation. *Advances in Neural Information Processing Systems*, 34, 12533–12548.
- Song, Y., Dhariwal, P., Chen, M., & Sutskever, I. (2023). Consistency models. arXiv preprint [arXiv:2303.01469](https://arxiv.org/abs/2303.01469).
- Soomro, K., Zamir, A.R., & Shah, M. (2012). Ucf101: A dataset of 101 human actions classes from videos in the wild. arXiv preprint [arXiv:1212.0402](https://arxiv.org/abs/1212.0402).
- Sun, T., Zhang, X., He, Z., Li, P., Cheng, Q., Liu, X., Yan, H., Shao, Y., Tang, Q., Zhang, S., Zhao, X., Chen, K., Zheng, Y., Zhou, Z., Li, R., Zhan, J., Zhou, Y., Li, L., Yang, X., Wu, L., Yin, Z., Huang, X., Jiang, Y.G., & Qiu, X. (2024). Moss: An open conversational large language model. Machine Intelligence Research <https://github.com/OpenMOSS/MOSS>.
- Sun, Y., Wang, X., Liu, Z., Miller, J., Efros, A., & Hardt, M. (2020). Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, PMLR, pp 9229–9248.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. (2023). Llama: Open and efficient foundation language models. arXiv preprint [arXiv:2302.13971](https://arxiv.org/abs/2302.13971).
- Ubani, S., Polat, S.O., & Nielsen, R. (2023). Zeroshotdataaug: Generating and augmenting training data with chatgpt. arXiv preprint [arXiv:2304.14334](https://arxiv.org/abs/2304.14334).
- Wang, D., Shelhamer, E., Liu, S., Olshausen, B., & Darrell, T. (2020). Tent: Fully test-time adaptation by entropy minimization. arXiv preprint [arXiv:2006.10726](https://arxiv.org/abs/2006.10726).
- Wang, H., Ge, S., Lipton, Z., & Xing, E.P. (2019). Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems* 32.
- Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., & Torralba, A. (2010). Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, IEEE, pp 3485–3492.
- Zhang, L., Deng, Z., Kawaguchi, K., Ghorbani, A., & Zou, J. (2020). How does mixup help with robustness and generalization? arXiv preprint [arXiv:2010.04819](https://arxiv.org/abs/2010.04819).
- Zhang, M., Levine, S., & Finn, C. (2021a). Memo: Test time robustness via adaptation and augmentation. arXiv preprint [arXiv:2110.09506](https://arxiv.org/abs/2110.09506).
- Zhang, M., Cai, Z., Pan, L., Hong, F., Guo, X., Yang, L., & Liu, Z. (2022a). Motiondiffuse: Text-driven human motion generation with diffusion model. arXiv preprint [arXiv:2208.15001](https://arxiv.org/abs/2208.15001).
- Zhang, R., Fang, R., Zhang, W., Gao, P., Li, K., Dai, J., Qiao, Y., & Li, H. (2021b). Tip-adapter: Training-free clip-adapter for better vision-language modeling. arXiv preprint [arXiv:2111.03930](https://arxiv.org/abs/2111.03930).
- Zhang, T., Wang, X., Zhou, D., Schuurmans, D., & Gonzalez, J.E. (2022b). Tempera: Test-time prompting via reinforcement learning. arXiv preprint [arXiv:2211.11890](https://arxiv.org/abs/2211.11890).
- Zhao, S., Liu, Z., Lin, J., Zhu, J. Y., & Han, S. (2020). Differentiable augmentation for data-efficient gan training. *Advances in Neural Information Processing Systems*, 33, 7559–7570.
- Zhong, Z., Zheng, L., Kang, G., Li, S., & Yang, Y. (2020). Random erasing data augmentation. *Proceedings of the AAAI conference on artificial intelligence*, 34, 13001–13008.
- Zhou, K., Yang, J., Loy, C.C., & Liu, Z. (2022a). Conditional prompt learning for vision-language models. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp 16816–16825.

Zhou, K., Yang, J., Loy, C. C., & Liu, Z. (2022). Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9), 2337–2348.

manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted