

Learning Disentangled Prompts for Compositional Image Synthesis

Kihyuk Sohn, Albert Shaw, Yuan Hao, Han Zhang,
Luisa Polania, Huiwen Chang, Lu Jiang, Irfan Essa
Google Research

{kihyuks, albertshaw, yuanhao, zhanghan, polania, huiwenchang, lujiang, irfanessa}@google.com

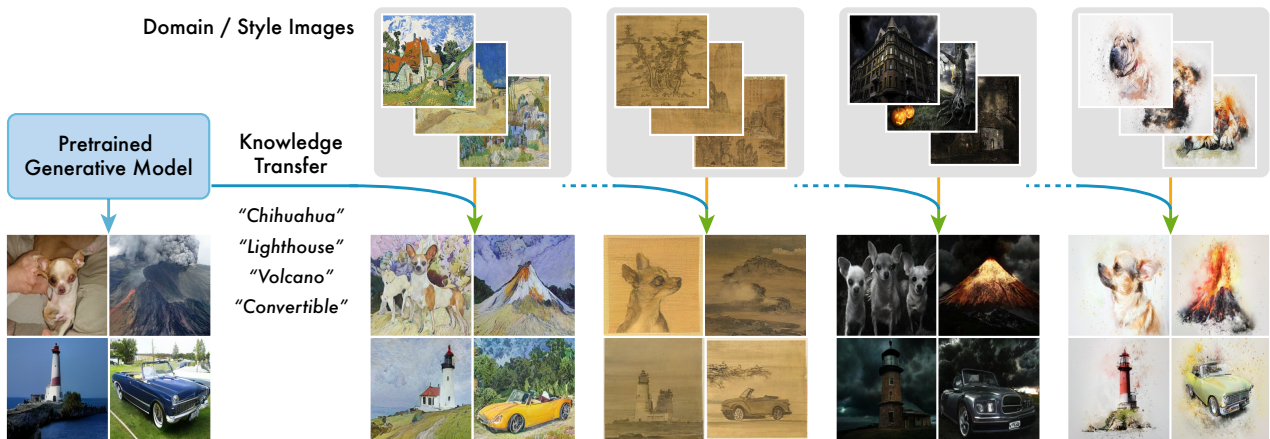


Figure 1. We demonstrate compositional generalization for image synthesis that composes the semantic knowledge (e.g., class) from the class-conditional generative vision transformers and the domain knowledge (e.g., style) from a few training images. For example, given a few (as low as 1) images of “Van Gogh house painting”, “Chinese painting”, “Haunted house”, and “Watercolor dog painting”, our method learns to compose their style with arbitrary object categories within the vocabulary of the pretrained generative model.

Abstract

We study domain-adaptive image synthesis, the problem of teaching pretrained image generative models a new style or concept from as few as one image to synthesize novel images, to better understand the compositional image synthesis. We present a framework that leverages a pretrained class-conditional generation model and visual prompt tuning. Specifically, we propose a novel source class distilled visual prompt that learns disentangled prompts of semantic (e.g., class) and domain (e.g., style) from a few images. Learned domain prompt is then used to synthesize images of any classes in the style of target domain. We conduct studies on various target domains with the number of images ranging from one to a few to many, and show qualitative results which show the compositional generalization of our method. Moreover, we show that our method can help improve zero-shot domain adaptation classification accuracy.

1. Introduction

Conditional image synthesis has witnessed remarkable

progress leveraging models that include: generative adversarial networks [16, 61, 2, 24, 25]; diffusion models [9, 42, 40]; and generative vision transformers [39, 10, 4, 57, 29]. Text-conditional image synthesis (or text-to-image), has received considerable attention. Text-to-image models, such as DALL-E [39], DALL-E 2 [38], GLIDE [32], Latent Diffusion Model [40], Imagen [42], Parti [60], show the power of compositional generalization, driven by free-form natural language prompts, for synthesizing novel images. For instance, in the synthesis of the “avocado armchair”, DALL-E invents a new type of armchair with color, shape and texture of avocado. While exact reasons responsible for such capabilities remain an active research area, it may be made possible, in part, due to a combination of an enormous amount of image-text paired training data in the scale of hundreds of millions to billions [44, 43], and a gigantic model with billions of parameters [60]. In addition, as argued in [47], the compositional power of text-to-image models is rather expected thanks to the built-in composable structure of text prompts. As such, an important but often overlooked question to answer is whether such compositional capability is

specific to models with text prompts or not.

In this paper, we study *the compositional generalization of image synthesis models trained on the standard ImageNet benchmark*. Since images in ImageNet [6] are labeled with a single salient object, synthesizing images by composing from multiple object classes may be infeasible. Instead, we study *zero-shot domain adaptive image synthesis (ZDAIS)*, a task of teaching pretrained class-conditional image generative models (*e.g.*, MaskGIT [4]) a new style from a few images of the target domain and applying it in synthesizing novel images by composing a known semantic (*e.g.*, class) and learned domain knowledge (*e.g.*, style). We illustrate in Fig. 1 the ZDAIS task and resulting synthesized images.

Solving ZDAIS requires explicit disentanglement of domain (*e.g.*, style) and semantic (*e.g.*, class) knowledge. This is in contrast with text-to-image models that implicitly learn compositionality via diverse text supervision that is unavailable in ZDAIS. This also differs from the conventional task of generative transfer learning (GTL) [56, 45, 33, 48] whose goal is to generate images of seen classes in the target domain, so no information disentanglement is required.

We argue that ZDAIS presents a new task to complement our understanding on compositional image synthesis. First, ZDAIS explores a compositional generalization on the public ImageNet dataset with an affordable computational budget. Second, it provides an understanding of compositionality without the text supervision, as ZDAIS represents the target domain using a few images. Third, we leverage various out-of-distribution ImageNet datasets that the community has established over years, such as an ImageNet-R [19], to diagnose zero-shot compositional generalization.

To this end, we present a solution built on the pretrained class-conditional MaskGIT [4], a non-autoregressive vision transformer for image synthesis that represents an image as a sequence of visual tokens. We adopt the generative visual prompt tuning [48] for knowledge transfer, while extending with a novel source class distilled visual prompt design to disentangle semantic and domain information using two kinds of learnable prompts. Subsequently, we present a domain adaptive classifier-free guidance to synthesize images improving not only visual fidelity but also domain adaptability. Our method is *simple*, introducing only one extra token to the pretrained generative transformer, and *efficient*, introducing less than 10k parameters trained within 10 minutes on a single GPU or TPU. Finally, our method is *effective*, as we show via thorough experimentation (Sec. 4).

We evaluate our method in three aspects. First, we conduct a qualitative study when only a few images are available from various target domains for training by visualizing synthesized images. We ablate design choices of our method to provide insights into how it works. Second, we report quantitative metrics on a newly proposed benchmark based on an ImageNet-R [19], which comes with many vi-

sual domains and their labels. Third, we make use of generated images by our method to solve zero-shot domain adaptation [37, 22] on the Office-home dataset [53]. Experimental results show our method is able to synthesize images of unseen classes in the target domain, confirming a compositional generalization of semantic and domain information.

2. Preliminary

As a preliminary, we discuss the Masked Generative Image Transformer (MaskGIT) [4] and the visual prompt tuning for generative transfer learning [48]. Sec. 5 includes a more comprehensive review of generative vision transformers and generative transfer learning.

Notation. Let \mathcal{X} be the input (*e.g.*, image) and \mathcal{Y} be the output (*e.g.*, label) domains, respectively. Let C be the number of classes. We use subscripts to denote domains, *i.e.*, src for the source and tgt for the target. Let D be the token embedding dimension, L be the number of transformer layers.

2.1. Masked Generative Image Transformer

Similarly to DALL-E [39] and Taming Transformer [10], MaskGIT [4] is a two-stage model for image synthesis consisting of a vector quantized (VQ) autoencoder [52, 10] and a transformer. Contrasting with those works, MaskGIT uses a non-autoregressive (NAR) transformer [8]. The VQ encoder converts image patches into a sequence of discrete tokens and VQ decoder maps discrete tokens back to an image. The NAR transformer learns to generate a sequence of discrete tokens given a condition (*e.g.*, class) if exists. For synthesis, the sequence of discrete tokens generated by the NAR transformer via scheduled parallel decoding [4] is fed into the VQ decoder, resulting in an image in a pixel space.

Since we are mostly interested in learning and using the NAR transformer, we refer [10, 4] for details of the VQ autoencoder. Let $\mathbf{z}(x) = [z_i(x)]_{i=1}^N$, $z_i(x) \in \{1, \dots, K\}$ be a token sequence, *i.e.*, an output of the VQ encoder for an image x . N is the number of tokens, corresponding to a spatial dimension of the encoder latent space, and K is the size of the token codebook. The NAR transformer is trained with the masked token modeling (MTM) loss [8, 4] as follows:

$$\mathcal{L}_{\text{MTM}}(\theta) = \mathbb{E}_{x, \mathbf{m}} \left[-\log \prod_{i: m_i=1} p_{\theta}(z_i | (t_c, \bar{\mathbf{z}})) \right] \quad (1)$$

where $t_c \in \mathbb{R}^{D \times 1}$ is a class token embedding for class c and θ is a transformer parameter. $\mathbf{m} = [m_i]_{i=1}^N$ is a sequence of binary random variables (*i.e.*, $m_i \in \{0, 1\}$) denoting which token to mask, and $\bar{\mathbf{z}} = [\bar{z}_i]_{i=1}^N$ is a masked token sequence, where $\bar{z}_i = z_i$ if $m_i = 0$ or a special token MASK otherwise. We omit x from $z(x)$ for brevity. Overall, \mathcal{L}_{MTM} optimizes to predict masked tokens values given the class token and the rest (unmasked) token values.

After obtaining $\theta^* = \arg \min_{\theta} \mathcal{L}_{\text{MTM}}(\theta)$, the scheduled parallel decoding [4] is used to generate a token sequence:

Require: $\mathbf{m} = [1]^N$, $\{n_1, \dots, n_T\}$, $\sum_{t=1}^T n_t = N$

- 1: **for** $t \leftarrow 1$ to T **do**
- 2: $\hat{z}_i \sim p_i \triangleq p_{\theta}(z_i | (t_c, \bar{\mathbf{z}}))$, $\forall i: m_i = 1$.
- 3: $z_i \leftarrow \hat{z}_i$, $m_i \leftarrow 0$, $i \in \arg \text{topk}_{i:m_i=1}^{n_t}(p_i)$.
- 4: **end for**

where $\arg \text{topk}^n$ returns a set of n indices with the highest values. $\{n_1, \dots, n_T\}$ is a mask schedule. See [4] for details.

2.2. Generative Visual Prompt Tuning

Visual prompt tuning [48] has been proposed as an effective way to transfer knowledge of pretrained generative vision transformers to the downstream generative tasks. For example, given a pretrained transformer parameter θ and a new task $(\mathcal{X}_{\text{tgt}}, \mathcal{Y}_{\text{tgt}})$, the visual prompt tuning introduces learnable prompt tokens $\{p_{\phi,c} \in \mathbb{R}^{D \times S}\}_{c \in \mathcal{Y}_{\text{tgt}}}$, where S is a number of prompt tokens. The parameter ϕ of the prompt token generator is learned by minimizing the MTM loss:

$$\phi^* = \arg \min_{\phi} \mathcal{L}_{\text{MTM}}(\phi) \quad (2)$$

$$\mathcal{L}_{\text{MTM}}(\phi) = \mathbb{E}_{x, \mathbf{m}} \left[-\log \prod_{i:m_i=1} p_{\theta}(z_i | (p_{\phi,c}, \bar{\mathbf{z}})) \right] \quad (3)$$

Similar decoding process is adopted, except that the class token embedding t_c is replaced with learned prompts $p_{\phi^*,c}$ to synthesize images belonging to \mathcal{X}_{tgt} instead of \mathcal{X}_{src} , conditioned on the class $c \in \mathcal{Y}_{\text{tgt}}$.

3. Zero-Shot Domain Adaptive Synthesis

In Sec. 3.1, we introduce the zero-shot domain adaptive image synthesis (ZDAIS) and describe the problem setting. Subsequently, we explain our method in Sec. 3.2.

3.1. Problem Formulation

We are given a class-conditional image generation model trained on the labeled source data $(\mathcal{X}_{\text{src}}, \mathcal{Y}_{\text{src}})$. In ZDAIS, we train a model on the labeled target data $(\mathcal{X}_{\text{tgt}}, \mathcal{Y}_{\text{tgt}})$, hoping to synthesize images of novel classes $c \in \mathcal{Y}_{\text{src}} \setminus \mathcal{Y}_{\text{tgt}}$ of a style consistent with \mathcal{X}_{tgt} . The setting is “zero-shot” in that no training data from $(\mathcal{X}_{\text{tgt}}, \mathcal{Y}_{\text{src}} \setminus \mathcal{Y}_{\text{tgt}})$ is given. Since the problem setting is relatively new, we clarify its setting by highlighting differences to existing problems below.

Relation to Generative Transfer Learning (GTL). Note that our setting is different from that of GTL [45, 48] whose goal is to synthesize an image of a class $c \in \mathcal{Y}_{\text{tgt}}$ of a target style \mathcal{X}_{tgt} when $(\mathcal{X}_{\text{tgt}}, \mathcal{Y}_{\text{tgt}})$ is available for training.

Relation to Neural Style Transfer (NST). While NST and ZDAIS synthesize an image belonging to \mathcal{X}_{tgt} , they solve different problems. NST [14, 15, 7] solves an image translation (*i.e.*, image-in, image-out), while ZDAIS solves a class-conditional image generation (*i.e.*, class-in, image-out).

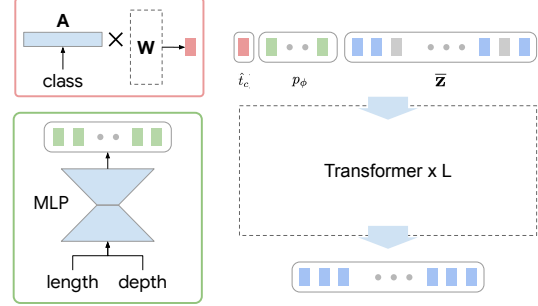


Figure 2. An overview of source class distilled visual prompt tuning. Source-class agnostic (green) and specific (red) prompts are trained to disentangle domain and class knowledge, while transformer and class embedding matrix (\mathbf{W}) are frozen.

One may argue that ZDAIS is simply a two step process of class-conditional image generation followed by NST. As shown in Secs. 4.2 and 4.3, the proposed method for ZDAIS shows a significant benefit over NST both on domain adaptive synthesis and domain adaptation by synthesis.

3.2. Proposed Method

We build our method on the class-conditional image synthesis model, MaskGIT [4]. We present a novel source-class distilled visual prompt tuning, which extends an idea of visual prompt tuning [48] to adapt to the target domain with a few training images (*e.g.*, 1~10 images), but with a novel component to compose with existing semantic (*e.g.*, class) knowledge of MaskGIT for ZDAIS.

3.2.1 Source Class Distilled Visual Prompt Tuning

As in Sec. 2.2, [48] has proposed to learn a class-conditional prompt $p_{\phi,c}$ for generative transfer that replaces the existing class token embedding t_c of the source domain. While this is a fine strategy to synthesize an image from $c \in \mathcal{Y}_{\text{tgt}}$, the resulting model cannot synthesize an image from classes outside of \mathcal{Y}_{tgt} (*e.g.*, $c \in \mathcal{Y}_{\text{src}} \setminus \mathcal{Y}_{\text{tgt}}$), as class and domain information are entangled in the learned prompt.

For ZDAIS, we want a representation where class and domain information are disentangled. Furthermore, to generalize to novel classes in $\mathcal{Y}_{\text{src}} \setminus \mathcal{Y}_{\text{tgt}}$, it is desirable to reuse class token embeddings of the source domain. To this end, we propose the following prompt design in Fig. 2:

$$p_{\phi,c} = \text{cat}(\hat{t}_c, p_{\phi}), \quad \hat{t}_c = \sum_{j=1}^{C_{\text{src}}} t_j a_{c,j}, \quad \forall c \in \mathcal{Y}_{\text{tgt}} \quad (4)$$

where $\mathbf{A} \triangleq [a_{c,j}] \in [0, 1]^{C_{\text{tgt}} \times C_{\text{src}}}$, $\sum_{j=1}^{C_{\text{src}}} a_{c,j} = 1$ is a learnable class affinity matrix between source and target classes. $\text{cat}(\cdot, \cdot)$ concatenates operands along the last dimension, *e.g.*, $p_{\phi,c} \in \mathbb{R}^{D \times (1+S)}$. As in Eq. (4), the proposed prompt has two types of tokens, class-specific \hat{t}_c and class-agnostic

p_ϕ . The learning objective of \mathbf{A} and ϕ is given as follows:

$$\phi^*, \mathbf{A}^* = \arg \min_{\phi, \mathbf{A}} \mathcal{L}(\phi, \mathbf{A}) \quad (5)$$

$$\mathcal{L}(\phi, \mathbf{A}) = \mathbb{E}_{x, m} \left[-\log \prod_{i: m_i=1} p_\theta(z_i | (\hat{t}_c, p_\phi, \bar{z})) \right] \quad (6)$$

Ideally, we want class-specific information to be captured by \hat{t}_c and the class-agnostic prompt p_ϕ learns only information about the domain. As shown in Appendix B.2, disentanglement seems to happen naturally when there are many images from various classes for the target domain, as p_ϕ is not sufficient to capture all information. However, in many cases we have only a handful of images from the target domain, and we end up under-utilizing \hat{t}_c , while p_ϕ learns all information without disentanglement. \hat{t}_c and p_ϕ need to be carefully designed for an information bottleneck (IB).

Enforcing Information Bottleneck. While there are principled ways to enforce IB [50, 1], we opt for a simple strategy by controlling the prompt capacity. First, for the class-agnostic prompt p_ϕ , we use $S=1$ with a low-dimensional (e.g., 8) bottleneck layer, to limit their expressive power. Second, we control the attention weight to the class-specific prompt \hat{t}_c to make sure that the transformer *does* use them. Note that the attention weight is computed by taking an inner product among token embeddings followed by the softmax operation.¹ For query (\mathbf{q}) and key (\mathbf{k}) token sequences, the attention weight is computed as follows:

$$\mathbf{q} = \bar{z}, \mathbf{k} = \text{cat}(\hat{t}_c, p_\phi, \bar{z}) \quad (7)$$

$$\begin{aligned} \text{Attn} &= \text{softmax}(\mathbf{q}^\top \mathbf{k}) \\ &= \text{softmax}(\text{cat}(\mathbf{q}^\top \hat{t}_c, \mathbf{q}^\top p_\phi, \mathbf{q}^\top \bar{z})) \end{aligned} \quad (8)$$

We propose to modify its computation as follows:

$$\text{softmax}(\text{cat}(\max\{\mathbf{q}^\top \hat{t}_c, \mathbf{q}^\top p_\phi\}, \mathbf{q}^\top p_\phi, \mathbf{q}^\top \bar{z})) \quad (9)$$

We use $\max\{\mathbf{q}^\top \hat{t}_c, \mathbf{q}^\top p_\phi\}$ in place of $\mathbf{q}^\top \hat{t}_c$. This ensures that attention weights to \hat{t}_c is always as large as those to p_ϕ . While simple, we find that our method effectively promotes the utilization of the class-specific token without introducing new hyperparameters. The pseudo-code for the prompt generator and the attention control is in Fig. 9.

Efficiency. Unlike [41, 26] based on the full fine-tuning of pretrained generative models, our method is efficient in terms of learnable parameters, train time and memory. For example, our method introduces less than 10k learnable parameters for p_ϕ and \mathbf{A} and can be trained on a single GPU within 1000 steps.

¹For presentation clarity, we ignore linear layers between token embeddings and the self-attention operator without loss of generality.

3.2.2 Domain Adaptive Classifier-Free Guidance

Classifier-free guidance (CFG) [21] has been extended for generative vision transformers [11, 60, 54] to improve the synthesis quality. The idea is to sample token values from both conditional and unconditional logits as follows:

$$\hat{z}_i \sim \text{softmax}(\rho_i) \quad (10)$$

$$\rho_i = (1 + \lambda)\rho_\theta(z_i | (t_c, p_\phi, \bar{z})) - \lambda\rho_\theta(z_i | (t_u, \bar{z})) \quad (11)$$

where t_c is a class token embedding of a class $c \in \mathcal{Y}_{\text{src}}$ and t_u is an unconditional token embedding for the source domain,² and λ is a guidance scale. We use ρ to denote logits, i.e., $p_\theta = \text{softmax}(\rho_\theta)$. As confirmed in Sec. 4.1.1, we find that it is crucial to drop both class and the prompt (p_ϕ) tokens when generating unconditional logits as it promotes to sample from the target instead of a source domain.

4. Experiments

We test the efficacy of our method for zero-shot domain adaptive image synthesis for various scenarios. In Secs. 4.1 and 4.2, we conduct a comprehensive study of ZDAIS when only a few training images from the target domain are available. Then, in Sec. 4.3, we show that synthesized images are used to train a classifier, achieving state-of-the-art zero-shot domain adaptation accuracy. An extended study with more training data is in Appendix B.2 due to a space limit.

4.1. Qualitative Evaluation of ZDAIS

We study ZDAIS with a few (e.g., 1~10) training images given from the target domain. We focus on the qualitative evaluation on various target domains in Sec. 4.1, and propose a new benchmark for a quantitative study in Sec. 4.2. **Dataset.** Similarly to [33], we collect 5~10 images for each target domain with a keyword such as ‘‘haunted house’’. The complete information on the images used is in Appendix D. **Setting.** A class-conditional MaskGIT [4] trained on the ImageNet [6] is used as a source model. As discussed in Sec. 3.2, we use $S=1$ for class-agnostic prompt. The bottleneck dimension is chosen between 2~32. Additional implementation details are in Appendix B.1.

While we collect images using the same keyword per target domain on purpose, we find that not all images belong to a single category. To this end, we assign each image its own class, i.e., $C_{\text{tgt}} = |\mathcal{Y}_{\text{tgt}}|$, instead of assigning them the same class. This change makes our method more flexible and applicable when target domain contains images of different object categories.

Evaluation. As no ground-truth images exist to compare with, we resort to visual inspection of synthesized images. After prompt tuning, we replace \hat{t}_c with a chosen ImageNet class for zero-shot image synthesis.

²Since we adopt publicly available MaskGIT pretrained model that is not trained with an additional unconditional token embedding, we simply average all class token embeddings to obtain t_u .

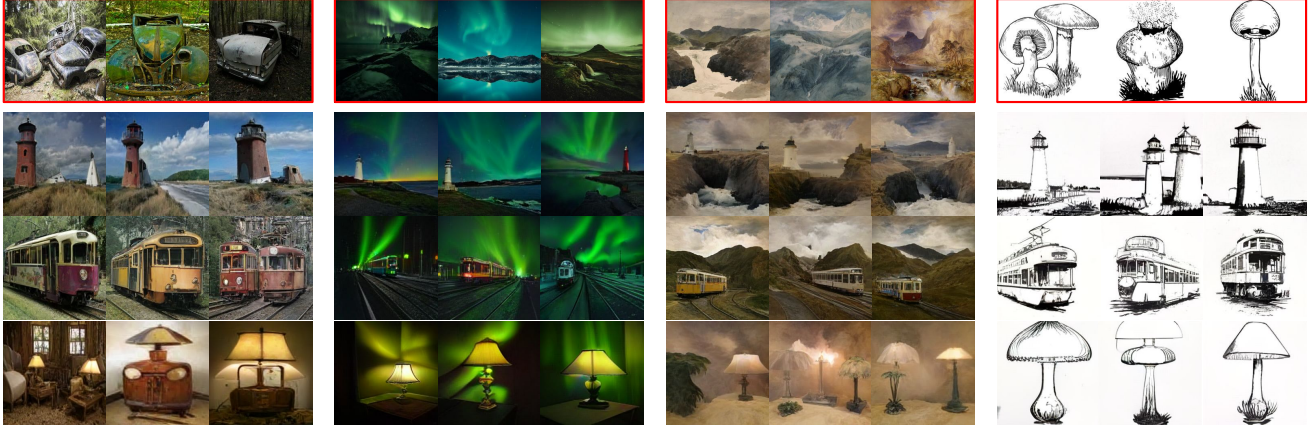


Figure 3. Zero-shot domain adaptive image synthesis from a few training images. Images inside red boxes are 3 training images (out of 3~10 in total) for each target domain, and the rest are synthesized. The class condition for each row is “lighthouse”, “tram” and “table lamp”. More images from diverse classes and target domains are in Fig. 10. Information on training images are in Appendix D.

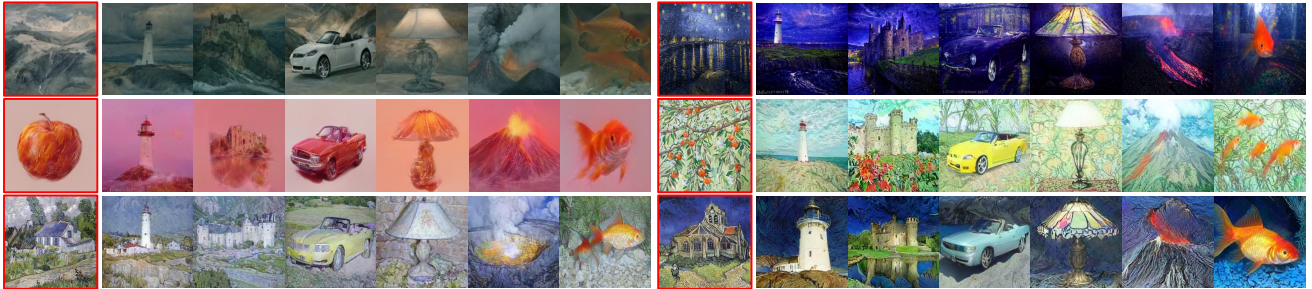


Figure 4. Zero-shot domain adaptive image synthesis from a single training image. Images inside red boxes are a **single (and only)** training image for each target domain, and the rest are synthesized. Class conditions are “lighthouse”, “castle”, “convertible”, “table lamp”, “volcano” and “goldfish”. More images from diverse classes are in Fig. 11. Information on training images are in Appendix D.

Results. We show synthesized images by our method along with target training images (in green boxes) in Figs. 3 and 10. The keyword for each target domain and the number of images are as follows: “abandoned cars” (10), “aurora mountain” (10), “inkpainting mountain” (3), “sketch mushroom” (4). We select “lighthouse”, “castle”, “tram”, “convertible” and “table lamp” as a class condition among 1000 ImageNet classes. It is clear that the model is able to disentangle class (*e.g.*, mushroom) and domain (*e.g.*, sketch) knowledge and transfer the domain knowledge to novel classes (*e.g.*, lighthouse, convertible) unseen from the target domain.

Moreover, Fig. 4 shows results on ZDAIS using a *single* training image from the target domain. We observe that our method not only adapts color or tone of an input image, but also the texture or overall mood of a training image well.

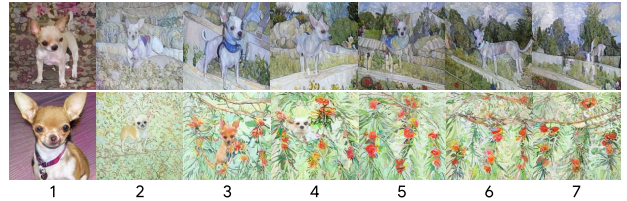
4.1.1 Ablation Study

Below, we conduct a series of ablation studies to better understand how the proposed method works. Due to the space limit, we provide extended study in Appendix B.3.

Bottleneck dimension. In Fig. 5, we show synthesized images from the models with varying bottleneck dimensions



(a) (top) “Van Gogh house painting” of *Chihuahua* and (bottom) “watercolor dog painting” of *Castle*. 6~10 target domain images are used.



(b) (top) “Van Gogh house painting” and (bottom) “oil painting appletree” of *Chihuahua*. 1 target domain image is used.

Figure 5. Ablation on the bottleneck dimension of class-agnostic prompt. The number for each column is the bottleneck dimension.

of the class-agnostic prompt generator. In Fig. 5a, we use 10 target training images and the bottleneck dimension is chosen between 2~128. We see that our method behaves some-

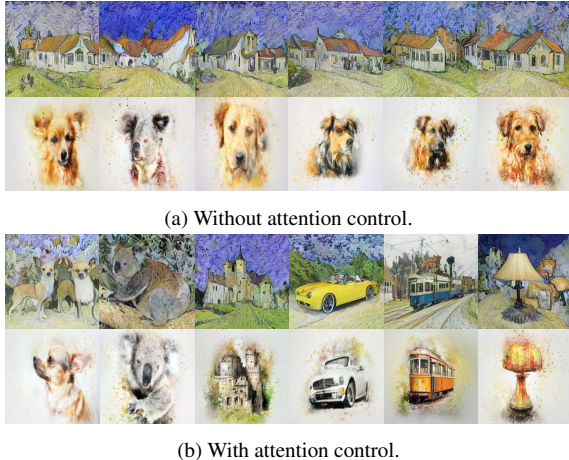


Figure 6. ZDAIS using 6~10 target domain images from “Van Gogh house painting” and “watercolor dog painting”. Classes for synthesis are (supposed to be) *Chihuahua*, *Koala*, *Castle*, *Convertible*, *Tram* and *Table lamp*.

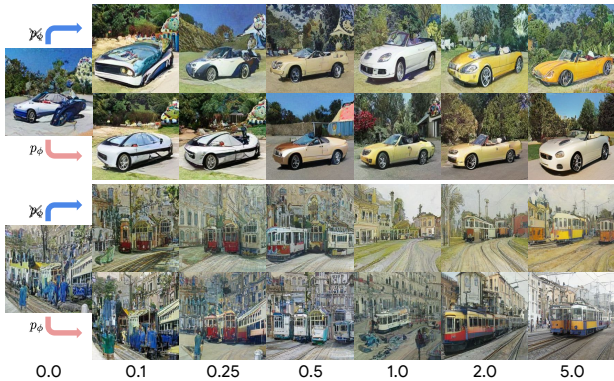


Figure 7. Synthesis with domain adaptive classifier-free guidance with varying guidance scale λ from 0 (*i.e.*, no guidance) to 5.0. Images in the first and third rows (blue arrow) are using our proposed formulation in Eq. (11), while those in the second and fourth rows (red arrow) are using p_ϕ for unconditional logits.

what sensibly to the bottleneck dimension as discussed in Sec. 3.2.1. When it is too small, it does not fully adapt to the target domain. On the other hand, when it is large, the class-agnostic prompt overpowers the class-specific prompt, and start to generate only target in-distribution images.

Similar behavior is observed in Fig. 5b where we conduct an experiment using 1 target training image with varying dimension from 1~7. When the target images are less, it is recommended to use even fewer bottleneck units.

Attention control (Eq. (9)). In Fig. 6a, we find that the class-agnostic prompt dominates without attention control and the model fails to synthesize images from novel classes, except for those close to the object category of the target domain training images (*e.g.*, chihuahua, koala). On the other hand, the attention control promotes the usage of the class-specific prompt and transfer the domain knowledge effec-

tively across novel classes as in Fig. 6b.

DA-CFG (Sec. 3.2.2). We visualize ZDAIS with varying guidance scales in Fig. 7. The model is trained on 8 images of “Van Gogh house painting”. While we see the transfer of target domain knowledge already happening with no guidance, the object shapes are seriously distorted. With a larger guidance scale, the model synthesizes more articulated images of an object.

Moreover, we confirm the importance of dropping class-agnostic prompt p_ϕ for domain adaptive synthesis. In Fig. 7 second and fourth rows, we show images synthesized without dropping p_ϕ from the second term of Eq. (11). While the model generates more articulated images with larger guidance scales, it starts to generate images in the source domain, losing adaptation capability.

Does accurate class prediction matter? As our method is designed to disentangle the domain knowledge using the class information from the source domain, it seems plausible to assume that the target images should be recognized as one of classes of the source domain. We check the predicted classes of images in Fig. 4 from the class affinity matrix \mathbf{A}^3 and report top 5 prediction. As in Tab. 1, most images have relevant classes in their top 5 prediction (*e.g.*, Van Gogh painting in the fourth row is predicted as gondola, whose typical images contain a boat in a river). This confirms that our method works as designed, but also suggests the limitation if the model fails to recognize object or concept in the image. We share such examples in Appendix B.3.1.

4.2. Quantitative Evaluation of ZDAIS

To test the efficacy of our method at scale, we need quantitative evaluation. Unfortunately, there is no benchmark to support our study yet. To this end, we propose a new benchmark using ImageNet-R [19]. Below, we explain the dataset and settings, followed by results.

Dataset. ImageNet-R [19] is built to systematically study the robustness of the image classification models. As such, the dataset contains images from diverse domains, such as painting, sketch, deviant art, embroidery, or toy, while sharing categories with the ImageNet. Images are labeled with both the class and domain. We test on 10 domains.

Setting. We pick 5 images from 5 classes for each domain and use them for visual prompt tuning. The list of images is in Tab. 5. We experiment with two settings. First, a single class-agnostic prompt is learned to model all 5 images (“single”). Second, we train one class-agnostic prompt per image and synthesize by randomly sampling prompt tokens (“ensemble”). Implementation details are in Appendix B.1.

Evaluation. Images from each domain of the ImageNet-R are used as ground-truths to compute FIDs. Since the number of images in each domain varies from 550 (origami) to

³This is done by ranking scores from c^{th} row of \mathbf{A} : $\mathbf{a}_c = [a_{j,c}]$.







image	top 5 classes predicted by A	image	top 5 classes predicted by A	image	top 5 classes predicted by A
	Komodo dragon, Irish Wolfhound, cassette, mountain , volcano		gondola , king penguin, rapeseed, drilling rig, paddle		plastic bag, Granny Smith apple , bee eater, carved pumpkin , lifeboat
	rose hip , monarch butterfly, lorikeet, apron, canoe		greenhouse , maypole, entertainment center , Groenendael dog, monastery		barn , entertainment center , chameleon, steam locomotive, castle

Table 1. Top 5 ImageNet classes by the class affinity matrix **A**. Relevant classes to the images are bold-faced.

Synthesis	art	cartoon	deviant art	embroidery	graffiti	graphic	origami	painting	sketch	toy	average
Source	92.0	134.4	89.0	166.3	166.5	138.8	155.5	78.4	94.0	97.7	121.3
Source + StyTr ² [7]	70.8	99.6	70.7	158.9	114.2	111.5	141.8	52.6	84.6	82.2	98.7
Target, in-dist [48]	168.0	118.3	173.8	176.1	168.0	182.0	142.9	145.1	163.1	133.3	157.1
Target, zero-shot (single)	95.5	85.9	78.2	121.3	114.0	109.3	116.7	64.0	47.5	84.2	91.6
Target, zero-shot (ensemble)	67.0	67.4	68.8	115.7	89.0	100.7	106.1	65.1	36.8	74.9	79.1

Table 2. FIDs (lower the better) on 10 target domains of the ImageNet-R [19]. “Source” refers to the class-conditional synthesis of the original MaskGIT trained on the ImageNet, “Target, in-dist” refers to the synthesis conditioned on learned prompt and target domain in-distribution classes, and “Target, zero-shot” refers the proposed zero-shot synthesis conditioned on learned prompt and source classes.

4634 (sketch), FIDs are not comparable between domains. We generate 3 types of images – source, target in-domain, and target zero-shot, for FID. Source refers to the synthesis by the pretrained MaskGIT, thus representing the distribution of an ImageNet [6]. Target in-domain refers to the synthesis by conditioning on the target domain in-distribution classes used for visual prompt tuning, which is consistent with a typical setting for generative transfer learning [48]. Finally, target zero-shot refers to the ZDAIS. Source classes are used as a condition only if there exist ground-truth images belonging to the same class in the target domain. We test 3 times and synthesize 20k images for each run.

Furthermore, we evaluate the two step process of class-conditional image generation followed by NST, as introduced in Sec. 3.1, to show the difference between ZDAIS and NST and how well our proposed method solves ZDAIS problem. For NST, we employ state-of-the-art StyTr² [7].

Results. We report FIDs in Tab. 2. We see that the proposed method (“Target, zero-shot”) achieves the lowest FID scores for most target domains. This confirms our previous findings in Sec. 4.1 that our method synthesizes not only good quality images, but also more faithful to each of the target domains than the MaskGIT without domain adaptive synthesis (“Source”). Note that this is not achievable with existing GTL methods [56, 45, 48]. While NST (“Source + StyTr²”) improves FIDs via style transfer to each of the target domains, they still fall short of our method. Moreover, the “ensemble” approach results in improved FIDs, as it retains the generation diversity better, which is important for certain domains (e.g., art, cartoon, graffiti). Qualitative difference between two methods are shown in Fig. 8. More images by our method are in Fig. 13.

4.3. Zero-shot Domain Adaptation by Synthesis

We test the quality of ZDAIS on zero-shot domain adaptation (ZSDA) [37, 22], a challenging task as there is no

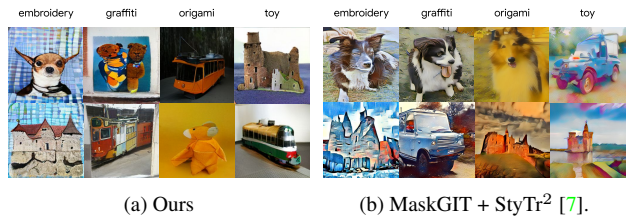


Figure 8. Synthesis on selected target domains of ImageNet-R.

data available in the target domain for the task of interest to train a classifier. We approach this by using synthesized images from our method conditioned on the classes for the task of interest as a training data. We conduct experiments on Office-home dataset [53] following the protocol in [22], which we further detail in Appendix B.4.

Similarly to previous sections, we use a class-conditional MaskGIT trained on ImageNet as a source model and apply our method using images from each split of target domains for zero-shot image synthesis. Following [22], we fine-tune ResNet-50 [17] trained on the ImageNet on the combined source and synthesized images.

Results are shown in Tab. 3. Compared to those trained on the source dataset only, classifiers trained on the combined source and synthesized images show consistent and significant improvement. In particular, we see 4~12% gain in accuracy when the target domain is “Clipart”, the most different from other domains, implying our method synthesizes images belonging to the target domain and classes. Visualization in Fig. 14 further confirms our claim. Moreover, our simple approach is shown to outperform state-of-the-art ZSDA method [22] by a large margin.

5. Related Work

To our knowledge, there is no previous work tackling the exact ZDAIS setting. We review related works of similar capabilities and highlight the differences.

Source domain Target domain	Art			Clipart			Product			Real			average
	Clipart	Product	Real	Art	Product	Real	Art	Clipart	Real	Art	Clipart	Product	
Source only	60.9	75.1	86.5	66.8	70.9	76.5	64.2	57.9	83.7	81.8	66.7	88.9	73.3
Source + Synth (w/o adapt)	61.8	84.3	87.8	77.4	83.2	87.2	77.4	62.8	89.3	80.7	63.4	88.9	78.7
Source + Synth (Ours)	69.0	84.8	88.2	80.9	84.2	86.9	80.9	69.8	89.1	82.7	70.7	89.0	81.4
[22]	71.0	76.5	85.1	62.1	68.7	75.1	64.4	69.2	82.0	77.9	76.2	88.5	74.7

Table 3. Zero-shot domain adaptation accuracy on Office-home dataset [53]. Numbers for [22] are taken from their paper. The best and those within standard error are bold-faced.

Generative Transfer Learning aims to build a generative model for target tasks, which often have limited data for training, by transferring knowledge from the source model trained on a large dataset. Prior works [56, 45, 33, 59] have been focusing on transferring from generative adversarial network (GAN) [16, 2] and recently extended to other generative models [48]. While they are shown to generate images belonging to the target domain, these methods are limited to synthesizing images of in-distribution whose training data is given. On the contrary, we study the composition of the semantic (*e.g.*, object categories) from the source model and the style depicted by a few out-of-distribution images from unseen classes in the target domain.

Closest to ZDAIS, cross-domain adaptation is studied in [59] that transfers the style of a single target domain image of an irrelevant class. Their demonstration has been limited to synthesizing images of a certain class (*e.g.*, church) while we show generation from diverse classes in the source class-conditional image generation model. We also make a direct comparison in Fig. 15, showing that our method learns to transfer higher-level concepts than color, tone or texture.

Zero-shot Image Synthesis with Text Guidance [39, 38, 32, 40, 42, 60] is relevant to ours. The text-to-image models generate images by specifying the style in the text prompt, *e.g.*, “watercolor painting”, “Van Gogh starry night”. They are shown effective for well known styles describable in text, but we point out two differences. First, they are trained from a large amount of training images of these styles. For example, one can easily find more than hundred images of “Van Gogh starry night” or “watercolor painting” with variations from the LAION-400M dataset [44].⁴ In ZDAIS, we show that learning composable token of such style is done from one image that is unseen in source model training. Second, there exist many modes in the style distribution (*e.g.*, Fig. 12), which may not be easily articulated with a text phrase. In practice, to capture such details with text-to-image models, one requires an extensive prompt engineering [34, 30, 20] and careful tuning [27, 49, 58, 36, 13, 5]. Nonetheless, as pointed out in [13], text guidance may be inherently ambiguous, whereas such ambiguity could be resolved for our method as it is guided by reference images.

Closest to ours include textual inversion [12] and model

⁴We browse images using above keywords in [clip-retrieval UI](#).

fine-tuning [41, 26] of text-to-image diffusion models [40, 42]. [12] proposes to learn a new vocabulary representing reference images by textual inversion. In DreamBooth [41], a token is learned with a ground-truth class prior, showing a better decomposition of style and semantic information. While they leverage the compositional capability of text-to-image models, we propose novel techniques, such as attention control or domain adaptive CFG, to enhance disentanglement for class-conditional generative models. Moreover, unlike [41], our method does not require model fine-tuning, making it not only efficient but also amenable to compose learned tokens. Our method is applicable to wide range of conditional generative models, *e.g.*, text-conditioned generative transformers [60, 54].

Neural Style Transfer (NST) [14, 15, 7, 35, 51] translates a content image to a certain style characterized by the style reference image while maintaining the contents (*e.g.*, structure, identity). As such, NST methods learn to disentangle content and style information with a set of losses for structure, identity preservation and appearance transfer, whose signals are given from discriminative models like VGG [46] or DINO [3]. Our method makes use of the semantic knowledge of pretrained generative transformers to disentangle target style, and is trained only with a masked token modeling loss. While both methods are used to synthesize images that transfers appearance or style, two methods have different goals and thus are not directly comparable. We provide a visual comparison in Fig. 16 for interested readers.

6. Conclusion

Towards better understanding of compositional generalization of models for image synthesis, we study the task of zero-shot domain adaptive image synthesis (ZDAIS) in the space of a class-conditional MaskGIT [4]. Through an extensive empirical study, we show that compositional generalization is achievable even without large amounts of training images nor rich text supervision. Moreover, we demonstrate that synthesized images by state-of-the-art generative models helps the classification on the challenging zero-shot domain adaptation setting. The proposed techniques, such as source class distilled visual prompt tuning or domain adaptive classifier-free guidance, are generic, and we expect them to be readily applicable to other family of conditional generative models, such as text-to-image [60, 54].

References

- [1] Alessandro Achille and Stefano Soatto. Emergence of invariance and disentanglement in deep representations. *The Journal of Machine Learning Research*, 19(1):1947–1980, 2018. [4](#)
- [2] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2018. [1](#), [8](#)
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. [8](#)
- [4] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. *arXiv preprint arXiv:2202.04200*, 2022. [1](#), [2](#), [3](#), [4](#), [8](#), [20](#)
- [5] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance. In *European Conference on Computer Vision*, pages 88–105. Springer, 2022. [8](#)
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [2](#), [4](#), [7](#), [20](#), [26](#)
- [7] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11326–11336, 2022. [3](#), [7](#), [8](#), [19](#), [22](#)
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. [2](#)
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. [1](#)
- [10] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12873–12883, 2021. [1](#), [2](#)
- [11] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. *arXiv preprint arXiv:2203.13131*, 2022. [4](#)
- [12] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. [8](#)
- [13] Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022. [8](#)
- [14] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015. [3](#), [8](#), [22](#)
- [15] Golnaz Ghiasi, Honglak Lee, Manjunath Kudlur, Vincent Dumoulin, and Jonathon Shlens. Exploring the structure of a real-time, arbitrary neural artistic stylization network. *arXiv preprint arXiv:1705.06830*, 2017. [3](#), [8](#), [19](#), [22](#)
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. [1](#), [8](#)
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [7](#)
- [18] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020. [12](#)
- [19] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021. [2](#), [6](#), [7](#), [17](#), [20](#), [23](#), [24](#), [25](#)
- [20] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. [8](#)
- [21] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. [4](#)
- [22] Won Young Jhoo and Jae-Pil Heo. Collaborative learning with disentangled features for zero-shot domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8896–8905, 2021. [2](#), [7](#), [8](#), [22](#)
- [23] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. *arXiv preprint arXiv:2203.12119*, 2022. [20](#)
- [24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. [1](#)
- [25] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. [1](#)
- [26] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. [4](#), [8](#)
- [27] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2426–2435, 2022. [8](#)

- [28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 20
- [29] José Lezama, Huiwen Chang, Lu Jiang, and Irfan Essa. Improved masked image generation with token-critic. In *European Conference on Computer Vision*, pages 70–86. Springer, 2022. 1
- [30] Vivian Liu and Lydia B Chilton. Design guidelines for prompt engineering text-to-image generative models. In *CHI Conference on Human Factors in Computing Systems*, pages 1–23, 2022. 8
- [31] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 20
- [32] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 1, 8
- [33] Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A Efros, Yong Jae Lee, Eli Shechtman, and Richard Zhang. Few-shot image generation via cross-domain correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10743–10752, 2021. 2, 4, 8
- [34] Jonas Oppenlaender. A taxonomy of prompt modifiers for text-to-image generation. *arXiv preprint arXiv:2204.13988*, 2022. 8
- [35] Jeeseung Park and Younggeun Kim. Styleformer: Transformer based generative adversarial networks with style vector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8983–8992, 2022. 8
- [36] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. 8
- [37] Kuan-Chuan Peng, Ziyang Wu, and Jan Ernst. Zero-shot deep domain adaptation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 764–781, 2018. 2, 7
- [38] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1, 8
- [39] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 1, 2, 8
- [40] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 1, 8
- [41] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242*, 2022. 4, 8
- [42] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 1, 8
- [43] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022. 1
- [44] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 1, 8
- [45] Mohamad Shahbazi, Zhiwu Huang, Danda Pani Paudel, Ajad Chhatkuli, and Luc Van Gool. Efficient conditional gan transfer with knowledge propagation across classes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12167–12176, 2021. 2, 3, 7, 8
- [46] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 8
- [47] Gautam Singh, Sungjin Ahn, and Fei Deng. Illiterate dall-e learns to compose. In *International Conference on Learning Representations*, 2022. 1
- [48] Kihyuk Sohn, Yuan Hao, José Lezama, Luisa Polania, Huiwen Chang, Han Zhang, Irfan Essa, and Lu Jiang. Visual prompt tuning for generative transfer learning. *arXiv preprint arXiv:2210.00990*, 2022. 2, 3, 7, 8, 20, 23
- [49] Xuan Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. Dual diffusion implicit bridges for image-to-image translation. *arXiv preprint arXiv:2203.08382*, 2022. 8
- [50] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000. 4
- [51] Narek Tumanyan, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Splicing vit features for semantic appearance transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10748–10757, 2022. 8
- [52] Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *NeurIPS*, 2017. 2
- [53] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5018–5027, 2017. 2, 7, 8, 18, 26
- [54] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi

- Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description. *arXiv preprint arXiv:2210.02399*, 2022. 4, 8
- [55] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32, 2019. 20, 23
- [56] Yaxing Wang, Abel Gonzalez-Garcia, David Berga, Luis Herranz, Fahad Shahbaz Khan, and Joost van de Weijer. Minegan: effective knowledge transfer from gans to target domains with few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9332–9341, 2020. 2, 7, 8
- [57] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual synthesis pre-training for neural visual world creation. In *European Conference on Computer Vision*, pages 720–736. Springer, 2022. 1
- [58] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. Tedigan: Text-guided diverse face image generation and manipulation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2256–2265, 2021. 8
- [59] Ceyuan Yang, Yujun Shen, Zhiyi Zhang, Yinghao Xu, Jiapeng Zhu, Zhirong Wu, and Bolei Zhou. One-shot generative domain adaptation. *arXiv preprint arXiv:2111.09876*, 2021. 8, 18, 22
- [60] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gungjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 1, 4, 8
- [61] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 5907–5915, 2017. 1

```

1 import flax.linen as nn
2 import jax.numpy as jnp
3
4 class ClassAgnosticTokenGenerator(nn.Module):
5     d_embed: int # Bottleneck dimension (P)
6     d_token: int # Token dimension (D)
7
8     @nn.compact
9     def __call__(self, ids: jnp.ndarray):
10         """Calls class agnostic token generator.
11
12         Args:
13             ids: B x 1, values are all zeros.
14
15         Return:
16             jnp.ndarray of dimension B x 1 x D.
17         """
18         MLP_c = nn.Embed(1, self.d_embed)
19         MLP_t = nn.Dense(self.d_token)
20
21         return MLP_t(nn.LayerNorm(MLP_c(ids)))
22
23
24 class ClassDistilledTokenGenerator(nn.Module):
25     n_class: int # Number of target class (C_{tgt})
26
27     @nn.compact
28     def __call__(self, cls_ids: jnp.ndarray, W: jnp.ndarray):
29         """Calls class distilled token generator.
30
31         Args:
32             cls_ids: B x 1, class labels.
33             W: C_src x D, a matrix for class embedding tokens.
34
35         Return:
36             jnp.ndarray of dimension B x 1 x D.
37         """
38         Affn = nn.Embed(self.n_class_tgt, W.shape[0])
39
40         token = jax.nn.softmax(Affn(cls_ids)) * W # B x D
41         return token[:, None, :]
42
43
44 def compute_attention_weight(query, class_token, domain_token):
45     """Computes attention weights with attention control.
46
47     Args:
48         query: tensor of size D x sequence length
49         class_token: tensor of size D x 1
50         domain_token: tensor of size D x 1
51
52     Return:
53         tensor of size (sequence length) x (sequence length + 2)
54     """
55     key = jnp.concatenate((class_token, domain_token, query), axis=1)
56     attn_weight = jnp.einsum('dq,dk->qk', query, key)
57     attn_weight = jnp.concatenate((jnp.maximum(attn_weight[:, 0], attn_weight[:, 1])
58         , attn_weight[:, 1:]), axis=1)
59     return attn_weight

```

Figure 9. An example code for the token generator and the attention control of Eq. (9) in Flax-ish [18] format.

Dataset	Link
Van Gogh House painting	https://www.rawpixel.com/image/3868934/illustration-image-art-vincent-van-gogh-person https://www.rawpixel.com/image/3865273/illustration-image-art-vincent-van-gogh-house https://www.rawpixel.com/image/3866294/illustration-image-art-vincent-van-gogh-house https://www.rawpixel.com/image/3868302/illustration-image-art-vincent-van-gogh-house https://www.rawpixel.com/image/3864574/illustration-image-art-vincent-van-gogh https://www.rawpixel.com/image/3864611/illustration-image-art-vincent-van-gogh-house https://www.rawpixel.com/image/537424/free-illustration-image-van-gogh-factory https://www.rawpixel.com/image/537422/free-illustration-image-van-gogh-cottage
Van Gogh starry night	https://search-production.openverse.engineering/image/ec5f5215-1307-457f-bc14-3e17d4fa4735 https://search-production.openverse.engineering/image/278dd07a-0b55-416b-blf5-05798f64cf34
Haunted house	https://search-production.openverse.engineering/image/17662d45-aa50-474e-83b6-bd170eda9bd9 https://search-production.openverse.engineering/image/f010aca4-14d0-4464-a706-b66ed7fb8569 https://search-production.openverse.engineering/image/ce8bd8d7-a509-4ece-a9e2-eb8581d0fb00 https://search-production.openverse.engineering/image/2f1alalee-c033-4781-bad0-3e12212a2361 https://search-production.openverse.engineering/image/89e86034-fe58-4f32-8259-9dbc4bac8ebc https://search-production.openverse.engineering/image/089c1cb0-f070-4dlf-9fc5-a82a7939d411 https://www.rawpixel.com/image/5906049/photo-image-public-domain-house-halloween https://www.rawpixel.com/image/5964852/free-public-domain-cc0-photo https://www.rawpixel.com/image/6051791/free-public-domain-cc0-photo
Abandoned cars	https://search-production.openverse.engineering/image/9e826ffa-bac1-4892-b78a-04eded1cefef https://search-production.openverse.engineering/image/545bb05c-cafc-42fc-89cc-34b0e51593a2 https://search-production.openverse.engineering/image/b1e4bc4e-982e-4792-a513-bdea0a5f72cb https://search-production.openverse.engineering/image/b103ce3b-c74c-41c4-97e0-8c3a7563bd59 https://www.rawpixel.com/image/3090749/free-photo-image-abandoned-car-vintage https://www.rawpixel.com/image/5941829/free-public-domain-cc0-photo https://www.rawpixel.com/image/4026101/oldsmobile-route-66 https://search-production.openverse.engineering/image/fbb2922a-a03c-44e7-9571-6901bba27957 https://search-production.openverse.engineering/image/f579052d-d3c8-4a59-8217-97b9f4b6ae55 https://search-production.openverse.engineering/image/f8db9513-f0a6-40e0-9e23-e6e34b5e04f3
Watercolor dogs	https://pixy.org/6458158/ https://pixy.org/5790662/ https://pixy.org/6379346/ https://pixy.org/5792209/ https://pixy.org/5792070/ https://pixy.org/5787524/ https://pixy.org/5782945/ https://pixy.org/6487827/ https://pixy.org/6557843/ https://pixy.org/5788893/
chinese inkpainting	https://search-production.openverse.engineering/image/382e0452-14a5-4895-af2f-917dfec8f40c https://search-production.openverse.engineering/image/22a3f474-128d-4a69-8d4e-98f904a77170 https://search-production.openverse.engineering/image/3aa5eee8-3a49-4e8e-9df9-f6ca2e0dc01c https://search-production.openverse.engineering/image/0793e761-3437-4591-a0e9-5a8f4d53ea7f https://search-production.openverse.engineering/image/293879fc-9f71-410f-b96b-161493e14377 https://search-production.openverse.engineering/image/04edb78b-b848-430b-9166-c61efab08a3c
aurora mountain	https://unsplash.com/photos/CgoRzWX4CDg https://unsplash.com/photos/U_diPCXCxBu https://unsplash.com/photos/-OkHUsepnzw https://unsplash.com/photos/pDeagUYN-Pk https://unsplash.com/photos/ZJDM1s6ppY8 https://unsplash.com/photos/Hn8N4I4eHA0 https://unsplash.com/photos/fpaSXDuoHkc https://unsplash.com/photos/58X3XfxxeVU https://unsplash.com/photos/uWmWoH9maR4 https://unsplash.com/photos/wlyDuFs-kGY
inkpainting mountain	https://unsplash.com/photos/6fv0MEf3FUE https://unsplash.com/photos/NspHfyZnMBE https://unsplash.com/photos/Vc8GBqapdfs
oilpainting apple	https://unsplash.com/photos/LqTdG9SjpyA
oilpainting appletree	https://unsplash.com/photos/K3QvdUkcQp4
sketch mushroom	https://freesvg.org/mushrooms-vector-illustration https://freesvg.org/psm-v07-d144-common-meadow-mushroom https://freesvg.org/1517757743 https://freesvg.org/puffball
watercolor images	https://unsplash.com/photos/KRzt15I6xac https://unsplash.com/photos/0pJPixfGFVo https://unsplash.com/photos/YiFVwDcgu8 https://unsplash.com/photos/9dnNnTrHxMI https://unsplash.com/photos/TygrVhOtrE https://unsplash.com/photos/8D-0K6JUAE https://unsplash.com/photos/6NSVtoSYwV0 https://unsplash.com/photos/-KfLa4I4eTo https://unsplash.com/photos/-IAS_N85ada https://unsplash.com/photos/X2QwsspYk_0 https://unsplash.com/photos/TAZga9MibgA https://unsplash.com/photos/wvD0zZnRbcw https://unsplash.com/photos/6dY9cFY-qTo

Table 4. Image sources for experiments in Sec. 4.1.

Domain	File name	Domain	File name
art	n01443537/art_10.jpg n01855672/art_0.jpg n02129604/art_4.jpg n02814860/art_15.jpg n03452741/art_0.jpg	graphic	n01443537/graphic_3.jpg n01855672/graphic_1.jpg n02129604/graphic_1.jpg n02980441/graphic_6.jpg n03452741/graphic_0.jpg
cartoon	n01443537/cartoon_28.jpg n01855672/cartoon_3.jpg n02129604/cartoon_4.jpg n02814860/cartoon_22.jpg n03452741/cartoon_2.jpg	origami	n01443537/origami_9.jpg n01855672/origami_1.jpg n02129604/origami_1.jpg n02814860/origami_1.jpg n03452741/origami_3.jpg
deviantart	n01443537/deviantart_16.jpg n01855672/deviantart_11.jpg n02129604/deviantart_23.jpg n02814860/deviantart_20.jpg n03452741/deviantart_0.jpg	painting	n01443537/painting_11.jpg n01855672/painting_0.jpg n02129604/painting_13.jpg n02814860/painting_14.jpg n03452741/painting_7.jpg
embroidery	n01443537/embroidery_1.jpg n01855672/embroidery_0.jpg n02129604/embroidery_7.jpg n02814860/embroidery_5.jpg n07614500/embroidery_10.jpg	sketch	n01443537/sketch_14.jpg n01855672/sketch_10.jpg n02129604/sketch_0.jpg n02814860/sketch_1.jpg n03452741/sketch_21.jpg
graffiti	n01443537/graffiti_1.jpg n01855672/graffiti_0.jpg n02129604/graffiti_9.jpg n02814860/graffiti_0.jpg n07614500/graffiti_1.jpg	toy	n01443537/toy_14.jpg n01855672/toy_4.jpg n02129604/toy_4.jpg n02814860/toy_1.jpg n03452741/toy_7.jpg

Table 5. Image sources for ImageNet-R experiments in Sec. 4.2. For reader’s convenience, n01443537 is for “goldfish”, n01855672 for “goose”, n02129604 for “tiger”, n02814860 for “lighthouse”, n03452741 for “grand piano”, n02980441 for “castle”, and n07614500 for “ice cream”.

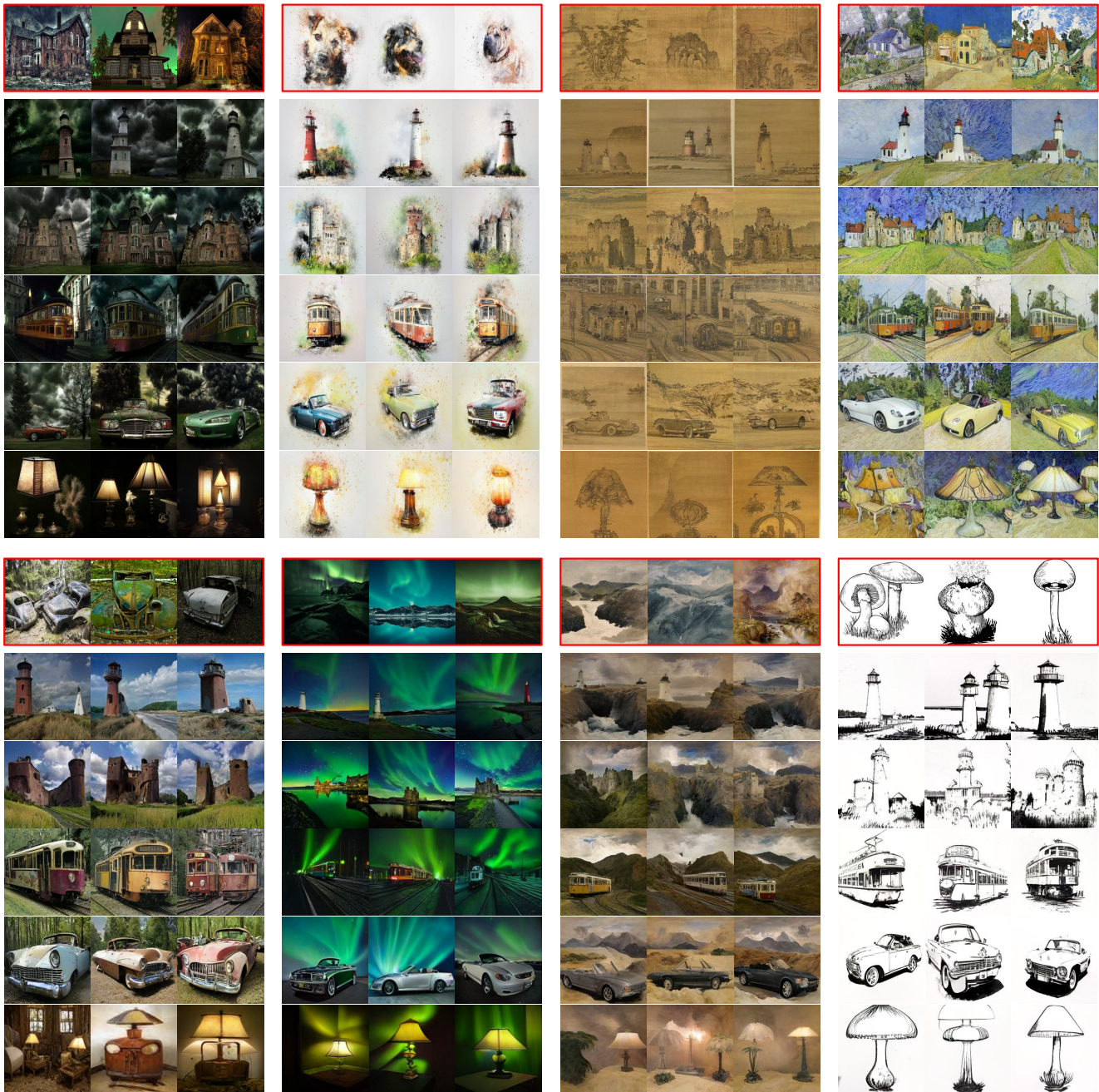


Figure 10. Zero-shot domain adaptive image synthesis from a few training images. Images inside red boxes are 3 training images (out of 6~10 in total) for each target domain, and the rest are synthesized. The class condition for each row is “light house”, “castle”, “tram”, “convertible”, and “table lamp”. We provide information on training images in Appendix D.



Figure 11. Zero-shot domain adaptive image synthesis from a single training image. Images inside red boxes are a **single (and only)** training image for each target domain, and the rest are synthesized. From left to right, top to bottom, class conditions are “lighthouse”, “castle”, “convertible”, “table lamp”, “volcano”, “goldfish”, “triumphal arch”, “teddy bear”, “tram”, “bolete”, “cliff”, “American alligator”.

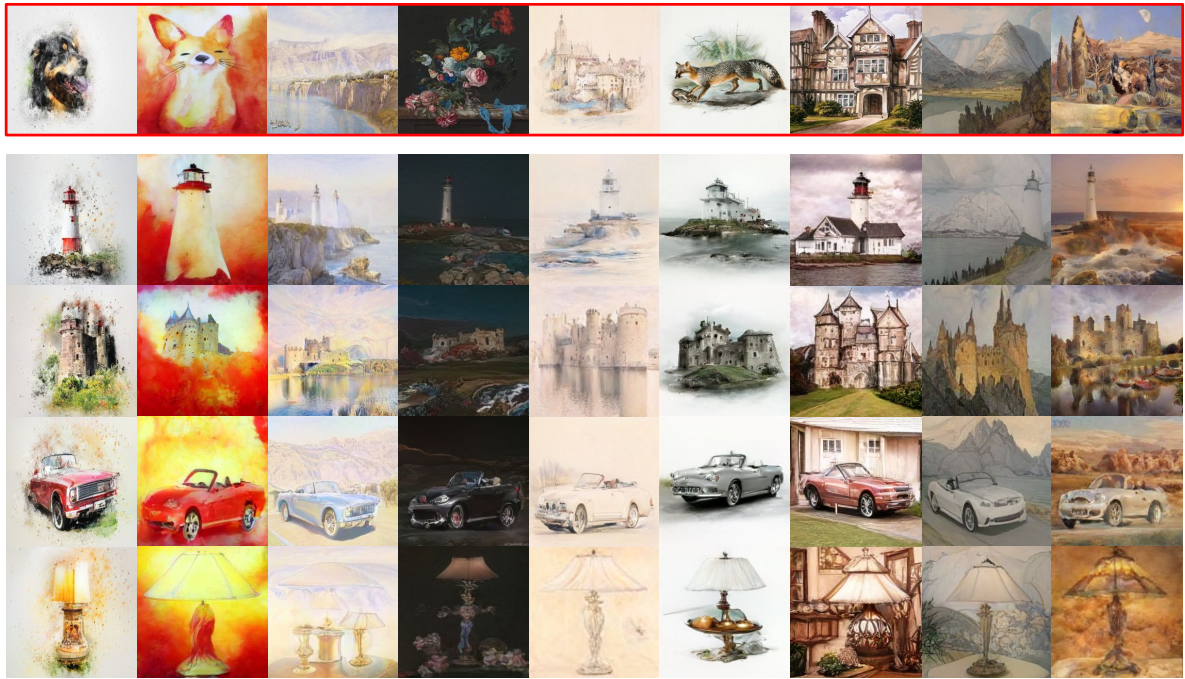


Figure 12. Many styles of “watercolor painting”. Our method is able to adapt to different styles of “watercolor painting” using a single reference image per target domain (in red boxes) and generalize to unseen concepts. This is in contrast with text-guided image synthesis methods where inherent ambiguity could occur due to ambiguity in natural language prompt.

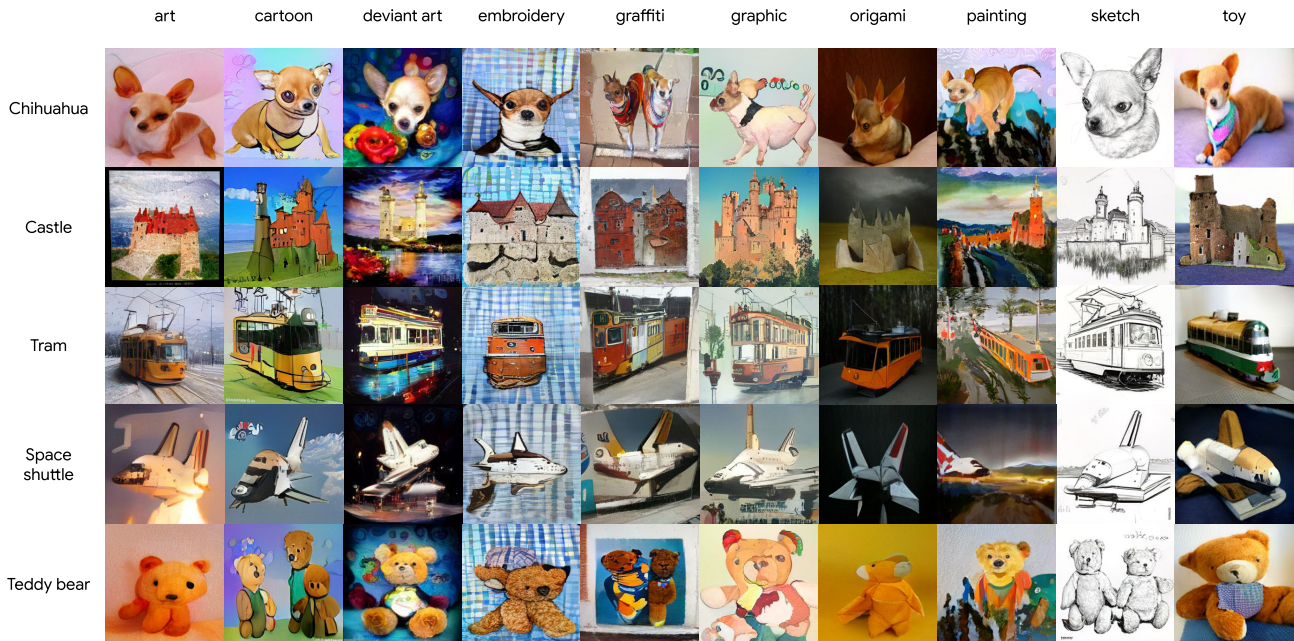


Figure 13. Zero-shot domain adaptive image synthesis on ImageNet-R [19].



Figure 14. Zero-shot domain adaptive image synthesis on Office-home [53]. From top to bottom, class condition for each two rows are “File cabinet”, “Backpack”, “Candle”, “Mop”, “Mug”, “Bed”.



Figure 15. Qualitative comparison with GenDA [59]. Images in the first row of each figure are taken from [59], and images in the second row of each figure are synthesized by our method. Images on the left (superman, Mona Lisa) are given as target domain images.

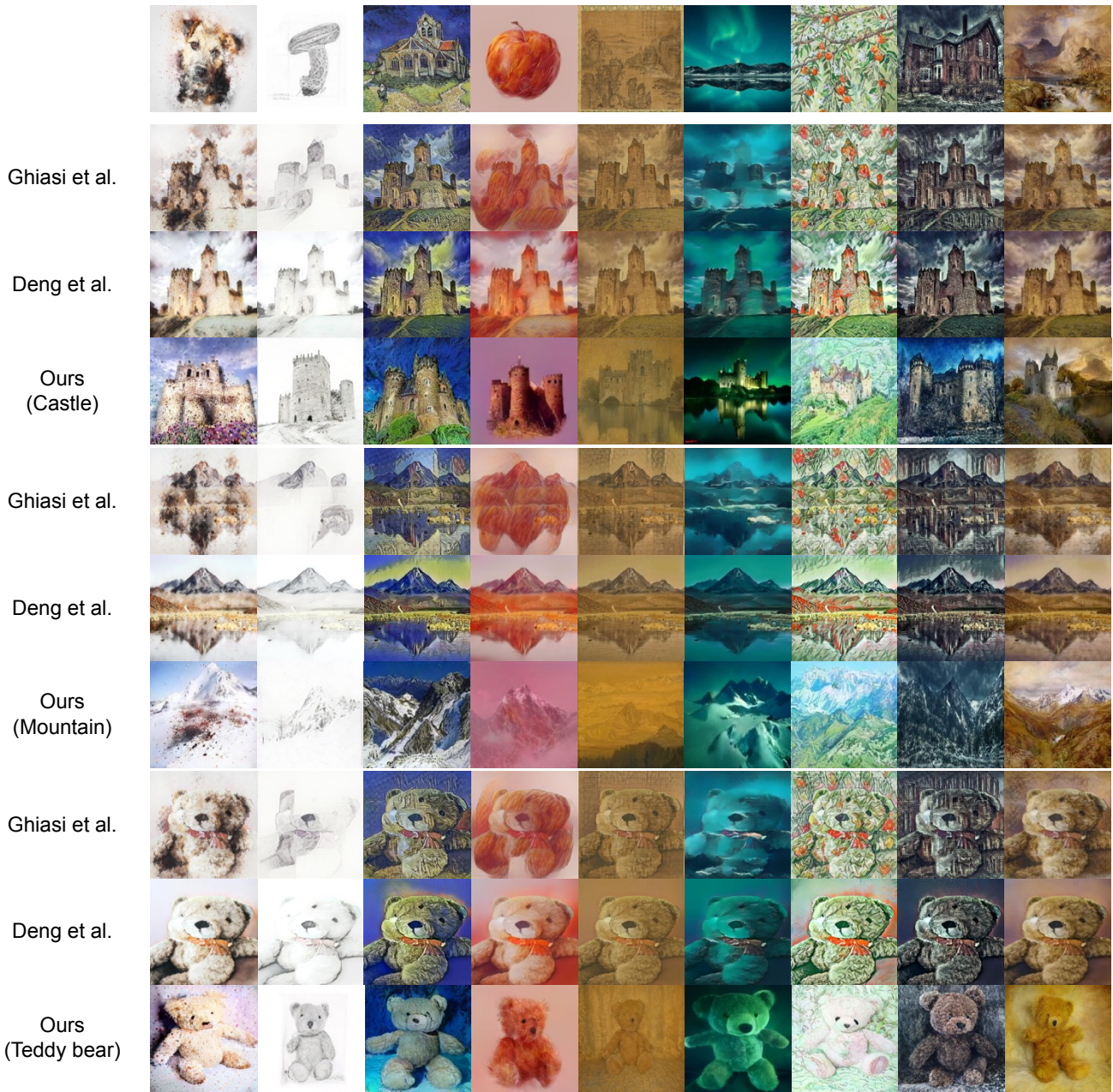


Figure 16. Comparison to neural style transfer methods of Ghiasi *et al.* [15] and Deng *et al.* [7]. Unlike neural style transfer methods whose input to the system consists of content and style images, our method does not take a content image as an input to enforce structural or identity-wise similarity between content and synthesized images. Instead, our method is class-conditional, *i.e.*, input to the system is composed of class label and style image. Conditioned classes used here are “castle”, “mountain”, and “teddy bear”.

A. Extended Descriptions on Method

Here, we provide more detailed information on our proposed method. Especially, similarly to [23], we find that the deep prompt, where separate learnable prompts are in place for all transformer layers, is more adaptive to the target domain than a shallow prompt. This is in contrast with [48] where only the shallow prompt has been discussed. Fig. 17 illustrates deep prompt used in our method. Systematic study comparing different prompt designs for generative transfer is beyond the scope of this work.

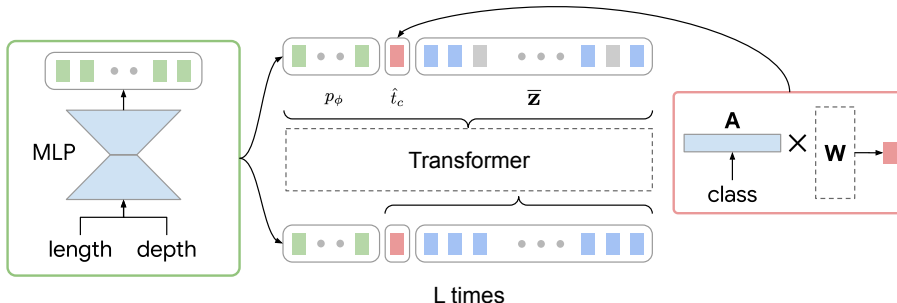


Figure 17. Generative visual prompt tuning with deep prompts. Separate learnable prompts (green) are provided for each transformer layer.

B. Experiments

B.1. Implementation Detail

Experiments in Secs. 4.1 and 4.2. All experiments are done using the class-conditional MaskGIT [4] trained on the ImageNet dataset, whose checkpoint is publicly available.⁵ We use a deep prompt, as explained in Appendix A, with $S = 1$. For optimization, we use an Adam optimizer [28] with the learning rate of 0.001 decaying with a cosine learning rate decay [31]. No weight decay nor learning rate warm-up is used. The parameters are updated for 1000 steps with the batch size of 128. We note that the number of training images are 1~10 for experiments of Sec. 4.1, and as such we repeat the training data by applying a stochastic data augmentation including random crop and horizontal flip. However, using a large batch size is not essential, and we empirically find that using batch size of 16 with 1000 optimization steps is sufficient for our method to work. The code will be released publicly for reproducible research.

B.2. ZDAIS with Many Training Images

Our prompt design in Sec. 3.2.1 is geared towards solving the problem when there are only a few training images provided from the target domain. As shown in Sec. 4.1.1, we need a stronger regularization, *e.g.*, reducing the bottleneck dimension, as there are fewer training images available. On the other hand, we hypothesize that the failure of disentanglement could be mitigated when there are many training images from diverse classes in the target domain.

Setting. We conduct experiments to confirm our hypothesis, where we use an entire data of ImageNet-R sketch [19] as a target domain training set, and test the ZDAIS performance on the ImageNet-sketch [55]. Note that ImageNet-sketch and ImageNet-R sketch are collected from two different groups and images in these two datasets do not necessarily overlap. Since ImageNet-R sketch contains 200 overlapping classes to the ImageNet [6], we exclude images from those classes when computing the FID. The total number of images used for visual prompt tuning is 4,634, *e.g.*, ≈ 23 images per class. In addition, we conduct experiments using 1~5 images per class (200~1,000 images in total, respectively) to see the impact of number of images per class on generation performance.

We train our models with varying class-agnostic token length $S = \{1, 16, 64, 128\}$ without the bottleneck layer. To accommodate class-agnostic prompt of length longer than 1, we extend an attention control mechanism in Eq. (9) as follows:

$$\text{softmax}(\text{cat}(\max\{\mathbf{q}^\top \hat{t}_c, \log \sum_{s=1}^S \exp(\mathbf{q}^\top p_{\phi,s})\}, \mathbf{q}^\top p_\phi, \mathbf{q}^\top \mathbf{q})) \quad (12)$$

⁵<https://github.com/google-research/maskgit>

where we replace $\mathbf{q}^\top p_\phi$ into $\log \sum_{s=1}^S \exp(\mathbf{q}^\top p_{\phi,s})$ from Eq. (9). This ensures sufficient attention to be given to the class-specific token even when the class-agnostic prompt has multiple tokens.

Result. We report results in Tab. 6. FIDs are computed between 20k synthesized images and 20k images sampled from classes that do not belong to ImageNet-R sketch dataset. Similarly to Tab. 2, “Source” refers to the synthesis of zero-shot classes using pretrained class-conditional MaskGIT without domain transfer, “Target, in-dist” refers to the in-distribution synthesis, *i.e.*, domain adaptive synthesis conditioned on the classes of ImageNet-R sketch, and “Target, zero-shot” refers to the zero-shot domain adaptive image synthesis. The FID measured by synthetic images of the source model is 77.0, implying that there exists a significant domain gap between the ImageNet and ImageNet sketch. When measured between target in-distribution synthetic images and the target zero-shot ground-truth images, we still get significantly lower FIDs around 30. This denotes that the FID is more affected by the domain than the class. Our proposed method for ZDAIS further reduces the FID, achieving 14.8 using 12 synthesis steps and 12.2 using 36 synthesis steps. This suggests that our method synthesizes images not only belonging to the target distribution, but also respecting the class semantics. We visualize synthesized images from the zero-shot classes (*i.e.*, which do not belong to classes of ImageNet-R sketch) in Fig. 19.

Different from experiments in Sec. 4.1 where the number of training images from the target domain is less than 10, we observe a successful disentanglement of domain and class without strongly regularizing the prompt capacity. On the contrary, we see an improvement in the domain adaptive zero-shot synthesis with longer class-agnostic prompts. This implies that when there are many training examples, it is still beneficial to have a powerful class-agnostic prompt to adapt to the target domain.

Another interesting observation is that the generation diversity is affected by the amount of training images. Fig. 20 shows synthesized images by the models trained with 1 or 5 images per class or full (≈ 23 images per class) training data. When comparing the most visually compelling results, we do not see too much differences between these models, and even synthesized images from the model trained with 1 images per class are highly realistic sketch images. However, the model trained with a few training images often sample from very narrow distribution, leading to a reduced diversity. As a result, as confirmed in Tab. 6, we get worse FIDs. The model performs similarly to the model trained on the full dataset in terms of a zero-shot FID when the number of training images per class is more than or equal to x .

B.3. Extended Ablation Study

B.3.1 Does accurate class prediction matter?

Our method is designed to disentangle semantic information and the target domain information using class-specific and class-agnostic prompts, respectively. As discussed in Sec. 4.1.1, our method is able to find reasonable semantic concepts within the vocabulary of the pretrained model. One question arises – what if our model fails to recognize semantic concept in the target images, or they cannot be described using the source classes? To this end, we conduct additional experiments with abstract painting images as target domain images.

Tab. 7 shows images used for the target domain and the class prediction by class affinity matrix \mathbf{A} . We show results in Fig. 18. To our surprise, our method learns to compose the style of the target domain to a certain degree with semantic knowledge from the source domain. However, we see some failure cases from the third row of Fig. 18. For example, synthesized images conditioned on the class “convertible” and “table lamp” are seriously entangled with the rotating pattern of the target domain training image. From Tab. 7, we find predictions being highly uncorrelated.

B.4. Zero-shot Domain Adaptation by Synthesis

B.4.1 Generation Details

We largely follow the implementation details in Appendix B.1. The class-conditional MaskGIT trained on the ImageNet is used as a source model, and we train a class-agnostic prompt of length $S = 128$ following Appendix B.2. The visual prompts are trained on the training set of each domain and split for 200 epochs.

For synthesis, we condition on the closest ImageNet classes we could find for each class of Office-home dataset. For example, for “Backpack” class of Office-home dataset, we choose “backpack (n02769748)” and “mailbag (n03709823)”. However, for some classes, there is not exactly matching classes between Office-home and ImageNet. To identify the closest classes, we train our model on the REAL WORLD domain of Office-home dataset, and match the classes using the class affinity matrix \mathbf{A} . This results in matching “Batteries” class to “oil filter” and “lighter”. Though they are not exactly matching, we find synthesized images are still useful to train a classifier. We list matching classes between Office-home and ImageNet in Tab. 8. Fig. 14 visualizes zero-shot domain adaptive synthesized images. We highlight how our model synthesizes images in the CLIPART domain, where the domain gap from the ImageNet might be the largest, by zero-shot transfer.

B.4.2 Classification Details

We mostly follow the same experimental setting and methodology as [22] with some exceptions. Following [22], we fine-tuned from an imagenet pretrained Resnet-50. We used a `jax` implementation using the weights from the standard `pytorch` model zoo model. We did a small hyperparameter search on our source-only classification results and used those hyperparameters all our classification training including Source + Synth training. We found this gave a much stronger source-only transfer baseline and in some cases even beat the domain adaptation results from [22]. We hyperparameter searched over learning rate in $\{0.01, 0.005, 0.001, 0.0005, 0.0001\}$ and selected $lr = 0.0005$. We used SGD with nesterov momentum. We found as is commonly known that this has resulted in a model with much better robustness than the ADAM trained model. We follow the learning rate schedule from [22] and cut the learning rate by 0.1 a $\frac{1}{3}$ and $\frac{2}{3}$ into the total training steps. Following [22] we also used a batch size of 8 and trained on one GPU. We chose to train all experiments for a consistent 6250 steps instead of 50 epochs in [22] since all the datasets have a wildly varying number of examples. Following [22] we reported the best accuracy on the target domain during each training run which we evaluated every 250 steps.

For our method where the classification model is fine-tuned on the combined source and synthesized images, we used the same hyperparameters for comparability, but added 1000 generated target domain examples from each class. We did a hyperparameter search over how often to sample from the source or the synth examples, $\{25\%, 50\%, 75\%, 100\%\}$ and used 75% of samples chosen were synthesized examples.

C. Comparison to Other Methods

C.1. Comparison to Cross-Domain Adaptation [59]

Cross-domain adaptation (CDA), briefly studied in [59], is relevant to ZDAIS. While CDA has been demonstrated on unconditional generation model of single category (*e.g.*, church generation model), our method makes use of class-conditional generation model to transfer the style of the target domain to diverse semantic classes.

In Fig. 15 we make a direct comparison to GenDA [59]. Results in the first row are taken from [59] and those in the second row are synthesized by our method. For synthesis, we condition with the church class (n03028079). When transferring from the “Superman” (Fig. 15a), GenDA generates yellow-ish church images, likely by learning color from the face region. On the other hand, our method generates cartoon-like church images with more abstract style. This implies that our method not only learns color or tone from the target domain training image, but higher-level concepts like cartoon or comic book.

C.2. Comparison to Neural Style Transfer

As we discuss in Sec. 5, ZDAIS could be relevant to the neural style transfer (NST) [14, 15, 7] in that both methods aim to generate images that belong to the target domain in some way. The critical difference is that the NST is a translation method, *i.e.*, image-in, image-out, while ZDAIS is a class-conditional generation method, *i.e.*, class-in, image-out. As such, NST evaluates structural similarity between input and output images, while we check semantic similarity between input class and output image for our method. Therefore, two methods are not directly comparable using the same evaluation criteria.

Nevertheless, in Fig. 16 we collect results from a few state-of-the-art neural style transfer methods [15, 7] and our method to highlight some qualitative differences. As mentioned, structural or semantic similarities are own characteristics of each method, so should not be counted as a factor for comparison. In terms of transferring the style, both methods seem to work well. We observe that synthesized images by our method are more harmonized than those by NST as there is no structural similarity constraint applied to our method. This is also seen from Fig. 14 or Fig. 13 where our model not only transfer artistic style, but also the high-level concepts, such as “Clipart”, “Origami”, or “Toy”. On the other hand, synthesized images by NST usually reflect the style more faithfully.

D. Information on Images for Target Domain

We provide the sources of images used in Tab. 4 for experiments in Sec. 4.1 and Tab. 5 for experiments in Sec. 4.2 for reproducible research.

Synthesis	Shot	Step	$S = 1$	$S = 16$	$S = 64$	$S = 128$
Source	–	12	77.0			
Target, in-dist [48]	≈ 23 (full)	12	30.2	31.4	31.5	31.5
Target, zero-shot (ours)	1	12	66.1	46.1	41.3	35.8
Target, zero-shot (ours)	2	12	40.2	26.6	24.8	22.1
Target, zero-shot (ours)	3	12	30.1	19.9	17.6	17.2
Target, zero-shot (ours)	4	12	26.5	17.8	16.6	15.7
Target, zero-shot (ours)	5	12	21.5	16.1	14.7	14.7
Target, zero-shot (ours)	≈ 23 (full)	12	18.9	15.1	14.4	14.8
Target, zero-shot (ours)	5	36	20.7	15.0	13.0	12.9
Target, zero-shot (ours)	≈ 23 (full)	36	18.4	13.7	12.6	12.2

Table 6. FID (lower the better) of models trained on ImageNet-R sketch [19] and tested on ImageNet-sketch [55]. The number of training images (“Shot”) and the number of decoding step (“Step”) are denoted.

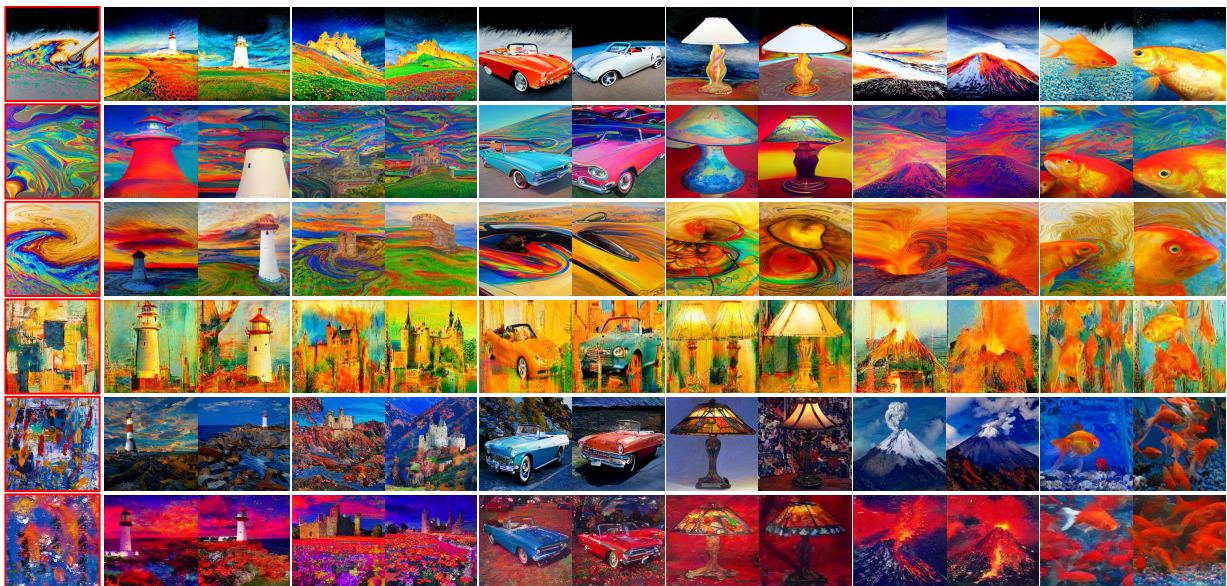


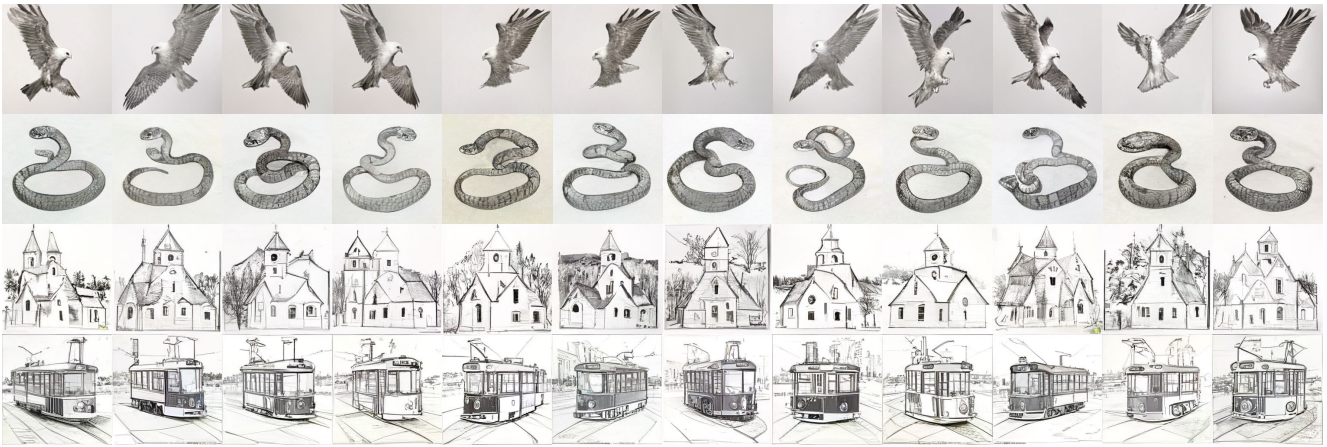
Figure 18. Zero-shot domain adaptive image synthesis from a single training image. Images inside red boxes are single training image for each target domain, and the rest are synthesized. Class conditions are “lighthouse”, “castle”, “convertible”, “table lamp”, “volcano” and “goldfish”. Information on training images are in Appendix D.

image	top 5 classes predicted by \mathbf{A}	image	top 5 classes predicted by \mathbf{A}	image	top 5 classes predicted by \mathbf{A}
	carousel, skunk feather boa, website, borzoi		carousel, cowboy boat, quilt, bubble, rock beauty fish		neck brace, gyromitra, poncho, tailed frog, clothes iron
	paintbrush, common gallinule, bookstore, stone wall, crate		paintbrush, goldfish, shower curtain, snowmobile, fox squirrel		paintbrush, ski, goldfish, feature boa, snorkel

Table 7. Top 5 ImageNet classes by the class affinity matrix \mathbf{A} . Relevant classes to the images are bold-faced.



Figure 19. Zero-shot domain adaptive image synthesis trained on ImageNet-R sketch [19]. Class conditions are “kite (n01608432)”, “banded gecko (n01675722)”, “water snake (n01737021)”, “African grey parrot (n01817953)”, “tiger beetle (n02165105)”, “Angora rabbit (n02328150)”, “church (n03028079)”, “cab (n02930766)”, “dining table (n03201208)”, “motor scooter (n03791053)”, “streetcar (n04335435)”, and “cliff (n09246464)”, which do not exist in the ImageNet-R sketch.



(a) 1 image per class.



(b) 5 images per class.



(c) Full (≈ 23 images per class).

Figure 20. Zero-shot domain adaptive image synthesis from ImageNet-R sketch [19] with varying amount of training images, (a) 1 or (b) 5 images per class, or (c) full, *i.e.*, approximately 23 images per class, in the target domain. Class conditions are “kite (n01608432)”, “water snake (n01737021)”, “church (n03028079)”, “streetcar (n04335435)”, which do not exist in the ImageNet-R sketch.

Office-home class	ImageNet class
Alarm Clock	analog clock (n02708093), digital clock (n03196217), stopwatch (n04328186)
Backpack	backpack (n02769748), mailbag (n03709823)
Batteries†	oil filter (n03843555), lighter (n03666591)
Bed	quilt (n04033995), four-poster (n03388549)
Bike	mountain bike (n03792782), unicycle (n04509417), bicycle-built-for-two (n02835271)
Bottle	pop bottle (n03983396), water bottle (n04557648), beer bottle (n02823428)
Bucket	bucket (n02909870), pot (n03991062), milk can (n03764736), caldron (n02939185)
Calculator	remote control (n04074963), hand-held computer (n03485407)
Calendar†	menu (n07565083)
Candles	candle (n02948072), torch (n04456115)
Chair	rocking chair (n04099969), throne (n04429376), barber chair (n02791124)
Clipboards†	menu (n07565083)
Computer	desktop computer (n03180011)
Couch	studio couch (n04344873)
Curtains	window shade (n04590129), shower curtain (n04209239)
Desk Lamp	table lamp (n04380533)
Drill	power drill (n03995372)
Eraser	rubber eraser (n04116512)
Exit Sign	street sign (n06794110)
Fan	electric fan (n03271574)
File Cabinet	file (n03337140)
Flipflops	sandal (n04133789)
Flowers	daisy (n11939491), yellow lady's slipper (n12057211)
Folder	envelope (n03291819)
Fork†	spatula (n04270147)
Glasses	sunglasses (n04356056)
Hammer	hammer (n03481172)
Helmet	crash helmet (n03127747), football helmet (n03379051)
Kettle	teapot (n04398044)
Keyboard	computer keyboard (n03085013), space bar (n04264628), typewriter keyboard (n04505470)
Knives	scabbard (n04141327), cleaver (n03041632), letter opener (n03658185)
Lamp Shade	lampshade (n03637318)
Laptop	laptop (n03642806), notebook (n03832673)
Marker	ballpoint (n02783161)
Monitor	screen (n04152593), monitor (n03782006)
Mop	swab (n04367480), broom (n02906734)
Mouse	mouse (n03793489)
Mug	coffee mug (n03063599), cup (n07930864)
Notebook	binder (n02840245)
Oven	rotisserie (n04111531), microwave (n03761084)
Pan	frying pan (n03400231), tray (n04476259), dutch oven (n03259280)
Paper Clip	safety pin (n04127249)
Pen	ballpoint (n02783161), fountain pen (n03388183)
Pencil	pencil sharpener (n03908714)
Postit Notes†	envelope (n03291819)
Printer	printer (n04004767)
Push Pin†	pinwheel (n03944341), syringe (n04376876)
Radio	radio (n04041544)
Refrigerator	refrigerator (n04070727)
Ruler	rule (n04118776)
Scissors†	screwdriver (n04154565), letter opener (n03658185), can opener (n02951585), corkscrew (n03109150)
Screwdriver	screwdriver (n04154565)
Shelf	plate rack (n03961711)
Sink	washbasin (n04553703)
Sneakers	running shoe (n04120489)
Soda	pop bottle (n03983396)
Speaker	loudspeaker (n03691459)
Spoon	ladle (n03633091), wooden spoon (n04597913)
TV	television (n04404412)
Table	dining table (n03201208)
Telephone	dial telephone (n03187595), pay-phone (n03902125)
ToothBrush†	paintbrush (n03876231)
Toys†	toyshop (n04462240), maraca (n03720891)
Trash Can	ashcan (n02747177)
Webcam†	projector (n04009552), tripod (n04485082)

Table 8. Class matching between Office-home [53] and ImageNet [6]. Office-home classes without exact match are denoted with †.