

Few-shot Image Generation via Cross-domain Correspondence

Utkarsh Ojha^{1,2} Yijun Li¹ Jingwan Lu¹ Alexei A. Efros^{1,3}
 Yong Jae Lee² Eli Shechtman¹ Richard Zhang¹
¹Adobe Research ²UC Davis ³UC Berkeley

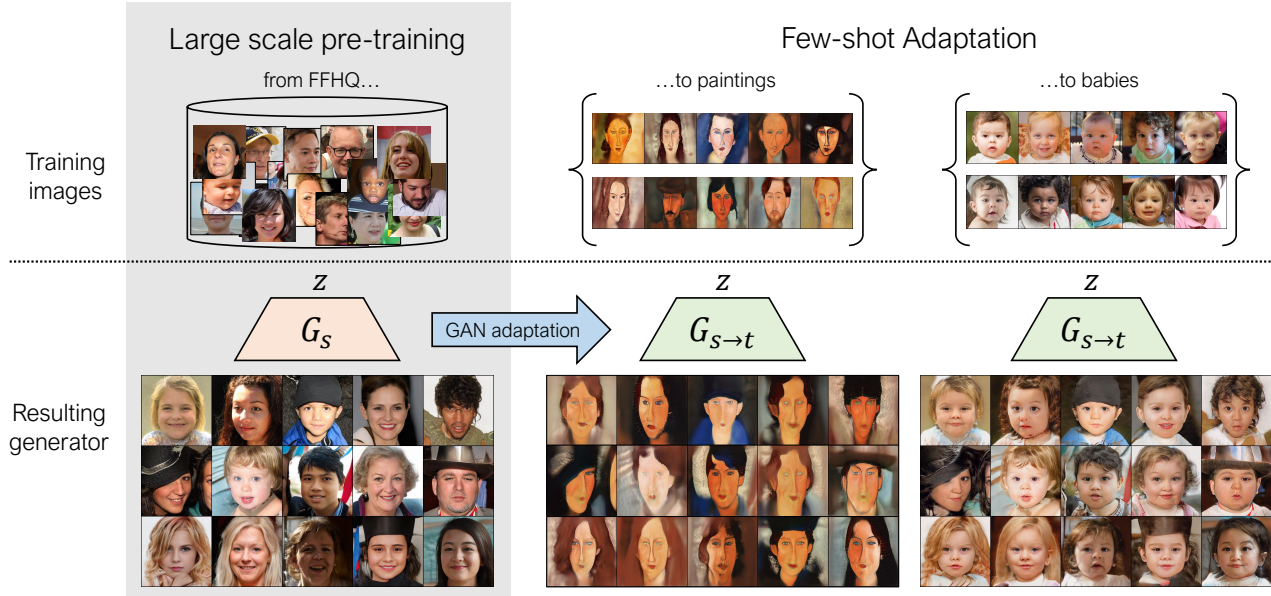


Figure 1: Given a model trained on a large source dataset (G_s), we propose to adapt it to arbitrary image domains, so that the resulting model ($G_{s \rightarrow t}$) captures these target distributions using extremely few training samples. In the process, our method discovers a one-to-one relation between the distributions, where noise vectors map to corresponding images in the source and target. Consequently, one can imagine how a natural face would look if Amedeo Modigliani had painted it, or how the face would look in its baby form. Please see our [webpage](#) for more results.

Abstract

Training generative models, such as GANs, on a target domain containing limited examples (e.g., 10) can easily result in overfitting. In this work, we seek to utilize a large source domain for pretraining and transfer the diversity information from source to target. We propose to preserve the relative similarities and differences between instances in the source via a novel cross-domain distance consistency loss. To further reduce overfitting, we present an anchor-based strategy to encourage different levels of realism over different regions in the latent space. With extensive results in both photorealistic and non-photorealistic domains, we demonstrate qualitatively and quantitatively that our few-shot model automatically discovers correspondences between source and target domains and generates

more diverse and realistic images than previous methods.

1. Introduction

Consider 10 portrait paintings by the incomparable Amedeo Modigliani [36], shown in Fig. 1 (middle). Given only these 10 paintings, would it be possible to train a model which can generate infinitely many paintings in the style of Modigliani? Unfortunately, contemporary generative models [11, 12, 13, 29, 3] require thousands of images to train properly, not 10. This problem is of practical importance, since many such domains of interest have a very limited collection of images (e.g., there are just 10 examples per artist in the Artistic-Faces dataset [36]).

Transfer learning serves as an alternative to training from scratch and has been explored in the context of generative

adversarial networks (GANs) to address the limited data regime. The key idea is to start with a source model, pre-trained on a large dataset, and adapt it to a target domain with limited data by either making only small changes to the network parameters to preserve as much information as possible [34, 23, 32, 20, 16], or by synthetically increasing the training data via data augmentation [39, 11]. Most of these methods, however, are designed for scenarios with more than 100 training images. When the number of available images is lowered to just a few [16], results often overfit to the training samples or are of poor quality.

In this work, we explore transferring a different kind of information from the source domain, namely *how images relate to each other*, to address the limited data setting. Intuitively, if the model can preserve the relative similarities and differences between instances in the source domain, then it has the chance to inherit the diversity in the source domain while adapting to the target domain. To capture this notion, we introduce a novel cross-domain distance consistency loss, which enforces similarity in the distribution of pairwise distances of generated samples before and after adaptation. Unlike domain adaptation approaches like image-to-image translation, here we are adapting models, not images.

Interesting properties emerge when enforcing this structure-level alignment between the two domains. Specifically, when the source and target domains are related (e.g., faces and caricatures), our approach automatically *discovers* a one-to-one correspondence between them and is able to more faithfully model the true target distribution in terms of both diversity and image realism, as shown in Fig. 1. When the two domains are unrelated (e.g., cars and caricatures), our approach is unable to model the target distribution but still discovers interesting part-level correspondences to generate diverse samples.

Since the few training samples only form a small subset of the target distribution we seek to approximate, we find it necessary to enforce realism in two different ways, to not inordinately penalize the diversity among the generated images. We apply an image-level adversarial loss on the synthesized images which should map to one of the real samples. For all other synthesized images, we only enforce a patch-level adversarial loss. In this way, only a small subset of our generated samples need to look like one of the few-shot training images, while the rest are only forced to capture their patch-level texture.

Contributions. Our main contribution is a novel GAN adaptation framework, which enforces cross-domain correspondence for few-shot image generation. Through extensive qualitative and quantitative results, we demonstrate that our model automatically discovers correspondences between related source and target domains to generate diverse and realistic images.

2. Related work

Few-shot learning. Representative few-shot classification approaches include learning a feature similarity function between the query and support examples [30, 26] and learning to learn how to adapt a base-learner to a new task [6, 22].

Few-shot image generation aims instead to hallucinate new and diverse examples while preventing overfitting to the few training images. Existing work mainly follows an adaptation pipeline, in which a base model is pre-trained on a large source domain and then adapted to a smaller target domain. They either embed a small number of new parameters into the source model [23, 32] or directly update the source model parameters, using different forms of regularization [20, 16]. Others employ data augmentation to reduce overfitting [39, 11] but are less effective under the extreme few-shot setting (e.g., 10 images). In contrast to prior work, we regularize the adaptation of the source model by transferring how images relate to each other in the source domain to the target domain, leading to plausible generation results, even with very few examples.

Domain translation. Translating images from the source domain is an alternative approach for generating more target domain data. However, such methods [9, 40, 41] require a large amount of training data for both source and target domains and are not suitable for the few-shot scenario. Recent work [17, 33, 25] has begun to address this issue via learning to separate the content and style factor, but requires large amount of labeled data (class or style labels). In our case, we assume access to a large amount of unlabeled data in the source domain and focus on adapting the source model to the target domain for unconditional image generation.

Distance preservation. To alleviate mode collapse in GANs, DistanceGAN [2] proposes to preserve the distances between input pairs in the corresponding generated output pairs. A similar scheme has been employed for both unconditional [27, 18] and conditional [19, 35] generation tasks to increase diversity in the generations. In our work, we aim to inherit the learned diversity from the source model to the target model and achieve this via our novel cross-domain distance consistency loss.

3. Approach

We are given a source generator G_s , trained on a large source dataset \mathcal{D}_s , which maps noise vectors $z \sim p_z(z) \subset \mathcal{Z}$, drawn from a simple distribution in a low-dimensional space, into images x . We aim to learn an adapted generator $G_{s \rightarrow t}$ by initializing the weights to the source generator and fitting it to a small target dataset \mathcal{D}_t .

A naive translation can be obtained simply by using a

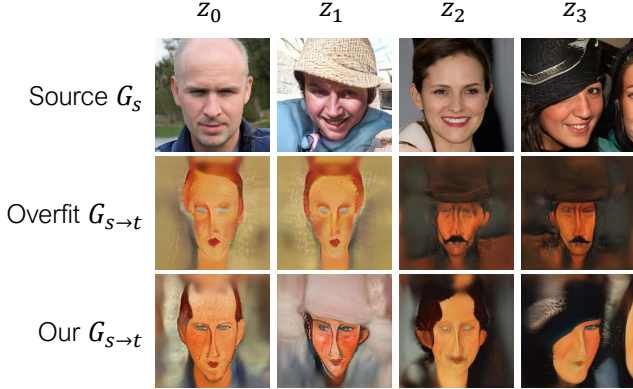


Figure 2: Adapting a source model (FFHQ) to 10 paintings using only the adversarial loss [34] results in overfitting, which leads to a loss in correspondence between source and target images. Our adaptation method preserves this property in a much better way, where the same noise maps to corresponding images in source/target.

GAN training procedure, with a learned discriminator D . With the non-saturating GAN objective, this corresponds to solving:

$$\begin{aligned} \mathcal{L}_{\text{adv}}(G, D) &= D(G(z)) - D(x) \\ G_{s \rightarrow t}^* &= \mathbb{E}_{z \sim p_z(z), x \sim \mathcal{D}_t} \arg \min_G \max_D \mathcal{L}_{\text{adv}}(G, D). \end{aligned} \quad (1)$$

Previous work [34] shows that this works well when the target dataset size exceeds 1000 training samples. However, in the extreme few-shot setting, this method overfits, as the discriminator can memorize the few examples and force the generator to reproduce them. This is shown in Fig. 2, where we see collapse after tuning the source model (top row) to the few-shot target dataset (middle row).

To prevent overfitting to generate diverse and realistic images (Fig. 2, bottom row), we propose a new cross-domain consistency loss (Sec. 3.1), which actively uses the original source generator to regularize the tuning process, and a “relaxed” discriminator (Sec. 3.2), which encourages different levels of realism over different regions in the latent space. Our approach is shown in Fig. 3.

3.1. Cross-domain distance consistency

A consequence of overfitting during adaptation is that relative distances in the source domain are not preserved. As seen in Fig. 2, the visual appearance between z_1 and z_2 collapses, disproportionately relative to z_1 and z_3 , which remain perceptually distinct. We hypothesize that enforcing preservation of relative pairwise distances, before and after adaptation, will help prevent collapse.

To this end, we sample a batch of $N + 1$ noise vectors $\{z_n\}_0^N$, and use their pairwise similarities in feature space to construct N -way probability distributions for each image. This is illustrated in Fig. 3 from the viewpoint of z_0 .

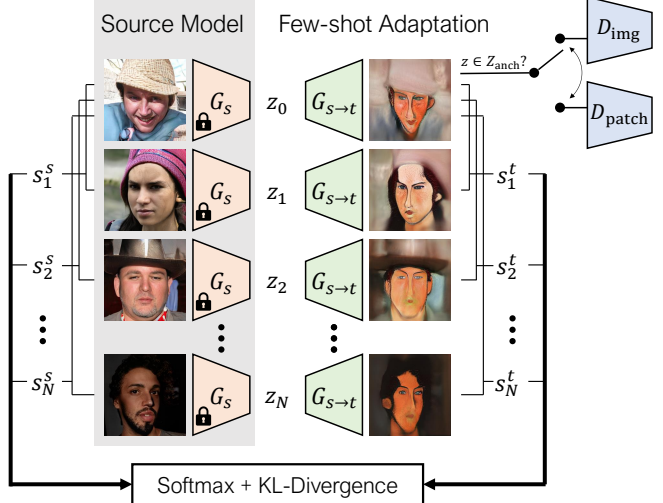


Figure 3: Our approach contains two key elements. (1) Cross-domain consistency loss $\mathcal{L}_{\text{dist}}$ aims to preserve the relative pairwise distances between source and target generations. In this case, the relative similarities between synthesized images from z_0 and other latent codes are encouraged to be similar. (2) Relaxed realism is implemented by using two discriminators, D_{img} for noise sampled from the anchor region (z_{anch}) and D_{patch} otherwise.

The probability distribution for the i^{th} noise vector, for the source and adapted generators is given by,

$$\begin{aligned} y_i^{s,l} &= \text{Softmax}(\{\text{sim}(G_s^l(z_i), G_s^l(z_j))\}_{\forall i \neq j}) \\ y_i^{s \rightarrow t,l} &= \text{Softmax}(\{\text{sim}(G_{s \rightarrow t}^l(z_i), G_{s \rightarrow t}^l(z_j))\}_{\forall i \neq j}), \end{aligned} \quad (2)$$

where sim denotes the cosine similarity between generator activations at the l^{th} layer. We are inspired by recent methods in contrastive learning [24, 4, 7], which converts similarities into probability distributions for unsupervised representation learning, as well as perceptual feature losses [10, 5, 28], which show that activations on multiple layers on discriminative networks help preserve similarity. We encourage the adapted model to have similar distributions to the source, across layers and image instances by using KL-divergence:

$$\mathcal{L}_{\text{dist}}(G_{s \rightarrow t}, G_s) = \mathbb{E}_{\{z_i \sim p_z(z)\}} \sum_{l,i} D_{KL}(y_i^{s \rightarrow t,l} || y_i^{s,l}). \quad (3)$$

While correspondence helps prevent collapse, we also modify the adversarial loss to further prevent overfitting.

3.2. Relaxed realism with few examples

With a very small target data size, the definition of what constitutes a “realistic” sample becomes increasingly over-constrained, as the discriminator can simply memorize the few-shot target training set. We note that the few training

images only form a small subset of the desired distribution and extend this notion to the latent space. We define ‘‘anchor’’ regions, $\mathcal{Z}_{\text{anch}} \subset \mathcal{Z}$, which form a subset of the entire latent space. When sampled from these regions, we use a full image discriminator D_{img} . Outside of them, we enforce adversarial loss using a patch-level discriminator D_{patch} ,

$$\begin{aligned} \mathcal{L}'_{\text{adv}}(G, D_{\text{img}}, D_{\text{patch}}) = & \mathbb{E}_{x \sim \mathcal{D}_t} \left[\mathbb{E}_{z \sim \mathcal{Z}_{\text{anch}}} \mathcal{L}_{\text{adv}}(G, D_{\text{img}}) \right. \\ & \left. + \mathbb{E}_{z \sim p_z(z)} \mathcal{L}_{\text{adv}}(G, D_{\text{patch}}) \right]. \end{aligned} \quad (4)$$

To define the anchor space, we select k random points, corresponding to the number of training images, and save them. We sample from these fixed points, with a small added Gaussian noise ($\sigma = .05$). We use shared weights between the two discriminators by defining D_{patch} as a subset of the larger D_{img} network [9, 40]; using internal activations correspond to patches on the input. The size depends on the network architecture and layer. We read off a set of layers, with effective patch size ranging from 22×22 to 61×61 .

3.3. Final Objective

Our final objective consists of just these two terms: $\mathcal{L}'_{\text{adv}}$ for the appearance of the target and $\mathcal{L}_{\text{dist}}$, which directly leverages the source model to preserve structural diversity:

$$\begin{aligned} G_{s \rightarrow t}^* = \arg \min_G \max_{D_{\text{img}}, D_{\text{patch}}} & \mathcal{L}'_{\text{adv}}(G, D_{\text{img}}, D_{\text{patch}}) \\ & + \lambda \mathcal{L}_{\text{dist}}(G, G_s). \end{aligned} \quad (5)$$

The patch discriminator gives the generator some additional freedom on the structure of the image. The adapted generator is directly incentivized to borrow the domain structure from the source generator, due to the cross-domain consistency loss. As shown in the top and bottom rows in Fig. 2, the model indeed discovers *cross-domain correspondences* between source and target domains.

We use the StyleGANv2 architecture¹ [13], pre-trained on a large dataset (e.g. FFHQ [12]) as our source model. We use a batch size of 4. Empirically, we find that a high λ , from 10^3 to 10^4 , to work well. Additional training details can be found in the supplementary.

4. Experiments

We explore different source \rightarrow target adaptation settings to analyze the effectiveness of our approach in preserving part-level correspondences between images generated from G_s and $G_{s \rightarrow t}$. We also investigate what kinds of correspondences emerge when the source and target domains are unrelated.

Baselines: We compare to baselines, which similar to ours, adapt a pre-trained source model to a target domain with

limited data. (i) Transferring GANs (TGAN) [34]: fine-tunes a pre-trained source model to a target domain with the same objective used to train the source model; (ii) Batch Statistics Adaptation (BSA) [23]: only adapts the scale and shift parameters of the model’s intermediate layers; (iii) MineGAN [32]: for a given pre-trained source (e.g. MNIST 0-8) and target (e.g. 9) domain, it transforms the original latent space of source to a space more relevant for the target (e.g. mapping all 0-8 regions to 4, being more similar to 9); (iv) Freeze-D [20]: freezes the high resolution discriminator layers during adaptation; (v) Non-leaking data augmentations [11, 39]: uses adaptive data augmentations (TGAN + ADA) in a way that does not leak into the generated results; (vi) EWC [16]: extends the idea of Elastic Weight Consolidation [15] for adapting a source model to a target domain, by penalizing large changes to *important* weights (estimated via Fisher information) in the source model.

Datasets: We use source models trained on five different datasets: (i) Flickr-Faces-HQ (FFHQ) [12], (ii) LSUN Church, (iii) LSUN Cars, and (iv) LSUN Horses [37]. We explore adaptation to the following target domains: (i) face caricatures, (ii) face sketches [31], (iii) face paintings by Amedeo Modigliani [36], (iv) FFHQ-babies, (v) FFHQ-sunglasses, (vi) landscape drawings, (vii) haunted houses, (viii) Van Gogh’s house paintings, (ix) wrecked/abandoned cars. We operate on 256×256 resolution images for both the source and target domains. Adaptation is done on 10 images from the target domain, unless stated otherwise.

4.1. Quality and Diversity Evaluation

Qualitative comparison Fig. 4 shows results on two target domains using different methods, all of which start from the same source model pre-trained on FFHQ (Source). We observe that TGAN strongly overfits to the available training data, as was the case for Amedeo paintings (Fig. 2). Using adaptive data augmentations (TGAN + ADA) has little to no effect on sketches, and further degrades the quality for caricatures, where augmentations (e.g. 90° rotations) leak into the generated images. FreezeD, MineGAN and EWC perform better than TGAN by generating slightly more diverse images. However, the diversity is only introduced through minor modifications among the few captured modes. For example, (i) the caricature results for EWC show multiple generations capturing the same person with different expression/hairlines; (ii) the results on sketches using MineGAN depict a person with similar attributes in multiple generations. Our method better captures the distribution of caricatures and sketches, and generates diverse and realistic images containing objects which do not appear in the training images (e.g. hats, in sketches). This is because our method is flexible enough to not penalize the generated images which do not adhere to the 10 training samples.

¹<https://github.com/rosinality/stylegan2-pytorch>

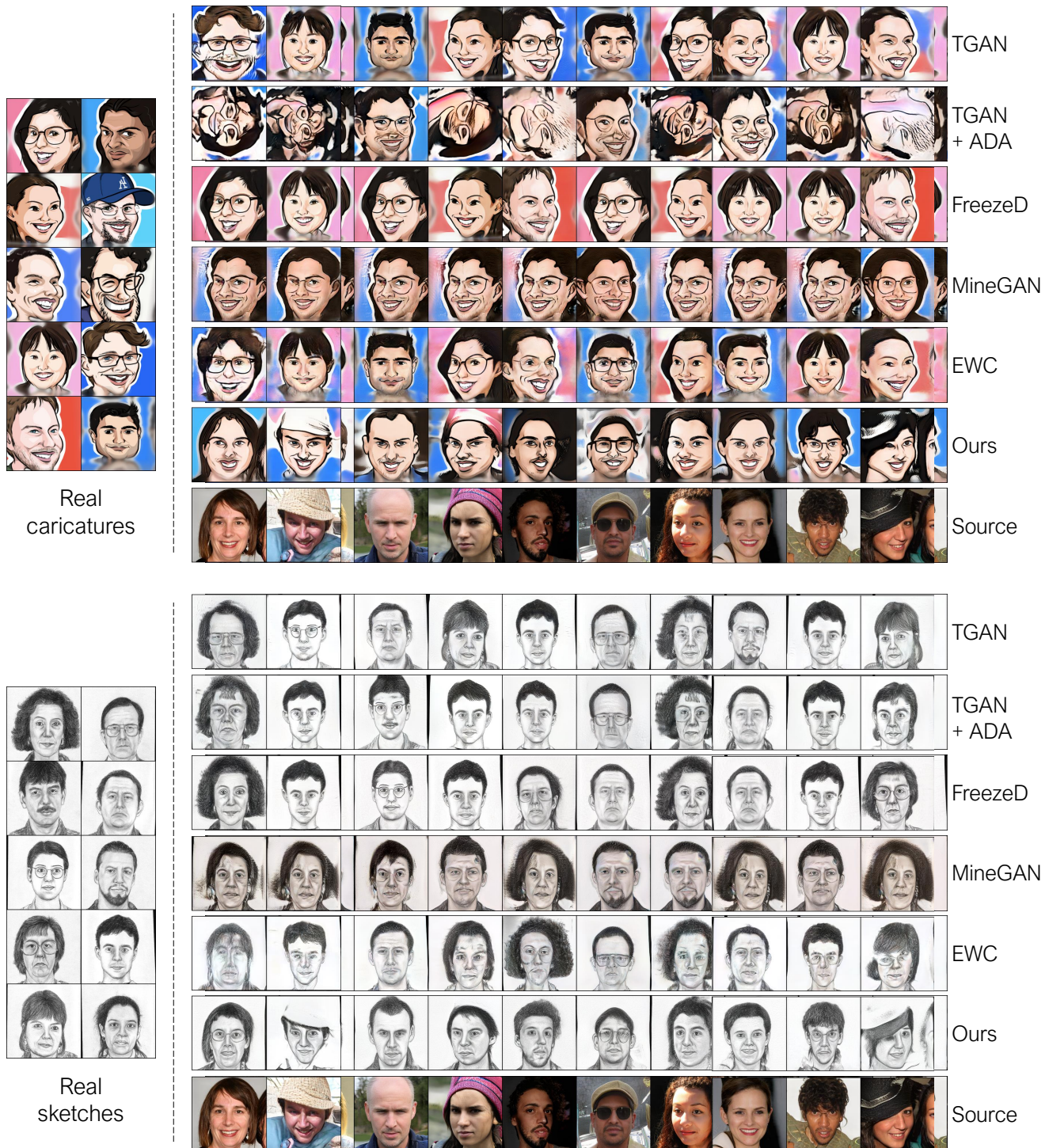


Figure 4: 10-shot image synthesis results for different methods, which start from the same source model (bottom). Keeping the noise vectors same (across columns), we observe that the baselines either overfit, or only capture a few modes in the target domain. Our method generates higher quality and more diverse results which better correspond to the source domain images generated from the same noise.

Quantitative comparison The original Sketches, FFHQ-babies, and FFHQ-sunglasses datasets roughly contain 300,

2500, and 2700 images, respectively. To simulate a few-shot setting, we randomly sampled 10 images from each

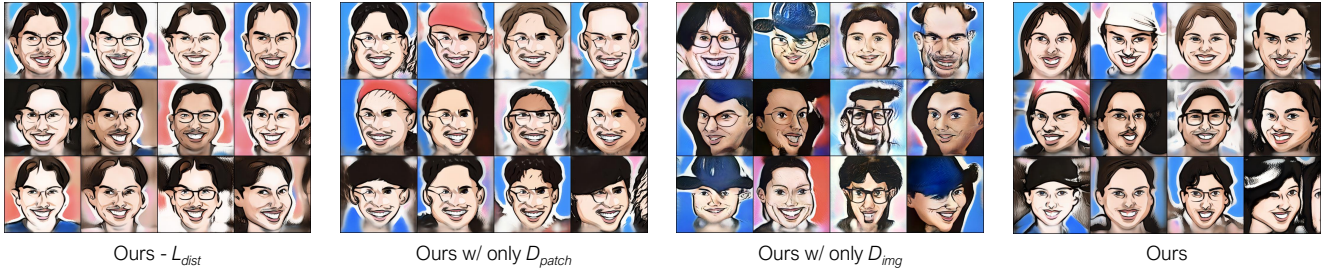


Figure 5: Effect of different components of our method. Absence of \mathcal{L}_{dist} makes some properties of the images (e.g. hairline) look very similar. Application of only one of D_{img} and D_{patch} degrades the image quality by distorting the face structure.



Figure 6: Visualizing the emerging correspondence in different adaptation settings. The generations from the domain of Van Gogh houses resemble the building structure from the source Church images. The generated wrecked/abandoned cars preserve the source car’s body and pose. Generations from FFHQ \rightarrow Sunglasses learn to add sunglasses to people’s faces.

dataset to train our model. For evaluation purposes, however, we can use the entire dataset to measure how well our generated images model the true distribution.

Table 1 shows the FID [8] scores. Our method significantly outperforms all baselines for Babies and Sunglasses. For domains with limited data, however, the FID score would not reflect the overfitting problem.

Ideally, we wish to assess the number of visually distinct images an algorithm can generate. In the worst case, the algorithm will simply overfit to the original k training images. To capture this, we first generate 1000 images and assign them to one of the k training images, by using lowest LPIPS distance [38]. We then compute the average pairwise LPIPS distance within members of the same cluster and then average over the k clusters. A method that reproduces the original images exactly will have a score of zero by this metric. Table 2 summarizes the distances for different baselines over three target domains. We see that

	Babies	Sunglasses	Sketches
TGAN [34]	104.79 \pm 0.03	55.61 \pm 0.04	53.41 \pm 0.02
TGAN+ADA [11]	102.58 \pm 0.12	53.64 \pm 0.08	66.99 \pm 0.01
BSA [23]	140.34 \pm 0.01	76.12 \pm 0.01	69.32 \pm 0.02
FreezeD [20]	110.92 \pm 0.02	51.29 \pm 0.05	46.54 \pm 0.01
MineGAN [32]	98.23 \pm 0.03	68.91 \pm 0.03	64.34 \pm 0.02
EWC [16]	87.41 \pm 0.02	59.73 \pm 0.04	71.25 \pm 0.01
Ours	74.39 \pm 0.03	42.13 \pm 0.04	45.67 \pm 0.02

Table 1: FID scores (\downarrow) for domains with abundant data. Standard deviations are computed across 5 runs.

our method consistently achieves higher average LPIPS distances, indicating more distinct images being generated. We also visualize the cluster centers and their members, to see if they are semantically meaningful. See supp. for details.

What role do different components of our method play? We use caricature as the target domain, and first study the effect of our framework with and without \mathcal{L}_{dist} ; see Fig. 5.



Figure 7: Visualizing the effect of unrelated source domains. In most cases, accurate modeling of a target distribution is not feasible when starting from an unrelated source domain. However, our method still discovers a correspondence on a part-level between the two domains, where different parts of the source (car’s tires) correspond to parts of the target (caricature’s eyes).

	Caricatures	Amedeo’s paintings	Sketches
TGAN [34]	0.39 ± 0.06	0.41 ± 0.03	0.39 ± 0.03
TGAN+ADA[11]	0.50 ± 0.05	0.51 ± 0.04	0.41 ± 0.05
BSA [23]	0.35 ± 0.01	0.39 ± 0.04	0.35 ± 0.01
FreezeD [20]	0.37 ± 0.01	0.40 ± 0.03	0.39 ± 0.03
MineGAN [32]	0.39 ± 0.07	0.42 ± 0.03	0.40 ± 0.05
EWC [16]	0.47 ± 0.03	0.52 ± 0.03	0.42 ± 0.03
Ours	0.53 ± 0.01	0.60 ± 0.01	0.45 ± 0.02

Table 2: Intra-cluster pairwise LPIPS distance (\uparrow). Standard deviation is computed across the k clusters (k = no. of training samples).

We see that leaving out $\mathcal{L}_{\text{dist}}$ reduces diversity among the generations, all of which have very similar head structure and hair style. We next study the different ways we enforce realism. What happens if we keep $\mathcal{L}_{\text{dist}}$, but use image-level adversarial loss through D_{img} on *all* generations? ‘Ours w/ only D_{img} ’ results reveal the problem of mode collapse at the part level (same blue hat appears in multiple generations) and the phenomenon where some results are only slight modifications of the same mode (same girl with and without the blue hat). Could we then only use D_{patch} to enforce patch-level realism on *all* generations? ‘Ours w/ only D_{patch} ’ shows the results, where we observe more diversity, but poorer quality compared to ‘Ours w/ only D_{img} ’. This is because the discriminator never gets to see a *complete* caricature image, and consequently does not learn the part-level relations which makes a caricature look realistic. ‘Ours’ combines all the ideas, resulting in generations which are diverse and realistic at both the part and image level.

4.2. Analyzing source ↔ target correspondence

As mentioned before, our method captures, to a large extent, the correspondence between the images from source and target domains, and preserves it during the adaptation process. In this section, we study this property in detail in different source to target adaptation settings.



Figure 8: Embedding an unseen caricature image (left), and visualizing its reconstruction from models adapted from different source to the caricature domain.

Related source/target domains Our method builds around the idea of discovering correspondence between a source and a target domain. But how well does that property actually hold? We first consider the FFHQ → Caricatures/sketches setting in Fig. 4, and analyze the baselines’ results. As seen previously in Sec. 4.1, the generated images mostly overfit to the training samples and are unable to borrow anything other than rough pose from the corresponding source. Using our method, on the other hand, we observe that the content of the natural source faces, as well as any of the accessories (e.g. hats/sunglasses), are preserved in both the sketch and caricature domain, depicting a much cleaner correspondence. We further study this for other source → target settings: (i) Church → Van Gogh houses, (ii) Cars → Abandoned cars, (iii) FFHQ → Sunglasses, as shown in Fig 6. When the source and target domains have similar semantics, the results generated from the same noise vectors in respective domains have clear correspondence.

Unrelated source/target domains We adapt four source models (FFHQ, Church, Cars, and Horses) to two target domains (Caricatures and Haunted houses) and present the results in Fig. 7. For FFHQ → Caricatures and Church → Haunted, the generated results from the adapted model mimic the target domain appearance and reflect correspondences with the source. For all the remaining scenarios, the adaptations do not capture the target distribution accurately. However, some part-level correspondences still

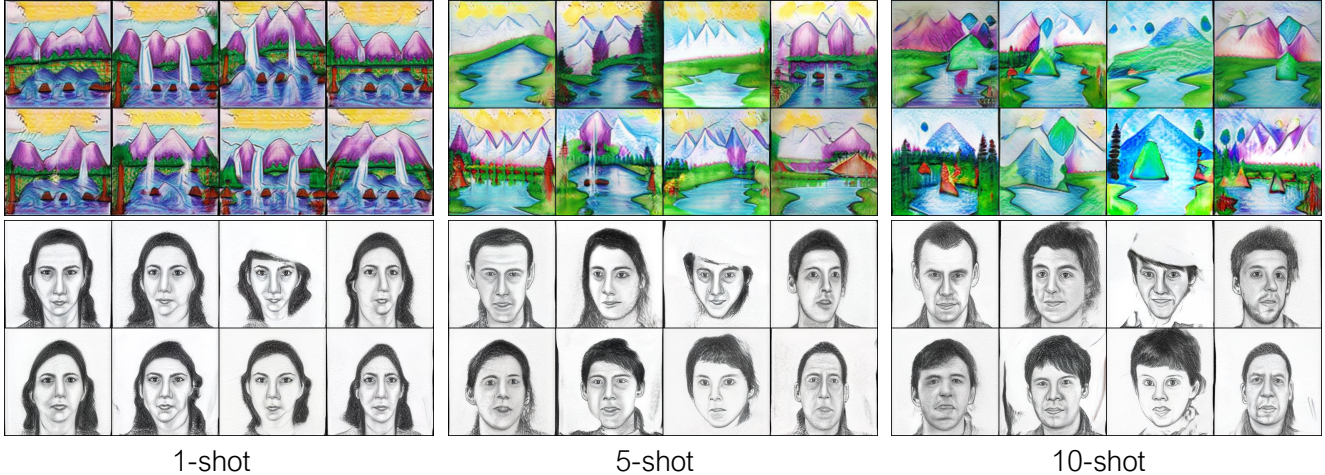


Figure 9: **Effect of training data size:** Even with just a single image, our method can capture it in different modes through the generations. Increasing the number of training samples results in more diversity in the generated sketches and landscapes.

emerge in these settings. E.g., (i) Church \rightarrow Caricatures: windows/doors of the Church roughly map to eyes of the caricature; (ii) Cars \rightarrow Caricatures: wheels/bumpers of the cars adapt to represent eyes/mouth of caricatures respectively; (iii) Cars \rightarrow Haunted: moon in the haunted houses (see supp. for the 10 images used) maps to the headlights, lighting them up; (iv) Horses \rightarrow Haunted: bottom of the horse legs adapt to doors of haunted houses.

Quantitative analysis of source/target relevance We translate four source models to four target domains: caricatures, haunted houses, landscape drawings, and abandoned cars. We then test how well a translated model can embed an *unseen* image from the respective target domain, E.g., after translating the four source models to the caricature domain, we use a new caricature and embed it into the four models using Image2StyleGAN [1]. The results are shown in Fig. 8. In Table 3, for each domain, we report the average similarity scores between five unseen inputs and their reconstructions. We see that FFHQ, Church, and Cars best reconstruct images from caricatures, haunted houses, and abandoned cars respectively which aligns with our intuition.

4.3. Effect of target dataset size

So far, all results are generated with 10 training images per target domain. We now explore how the dataset size affects the quality and diversity of the generated images. We consider two adaptation setups, Church \rightarrow Landscape drawings and FFHQ \rightarrow Sketches, and present results in 1-shot, 5-shot and 10-shot settings; see Fig. 9. The real images used in these settings can be found in the supplementary.

In 1-shot, our method introduces small variations to the single target image, for example the lady appears in different poses in the generated sketches, and the mountains and waterfall have different structures. The diversity of the re-

	Caricature	Haunted house	Wrecked car
FFHQ	0.158 ± 0.045	0.645 ± 0.024	0.643 ± 0.012
Church	0.294 ± 0.077	0.599 ± 0.028	0.621 ± 0.032
Cars	0.233 ± 0.106	0.635 ± 0.031	0.606 ± 0.057
Horses	0.299 ± 0.083	0.631 ± 0.028	0.619 ± 0.038

Table 3: Relevance of source and target domains, measured via LPIPS (\downarrow) between an unseen image and its reconstruction. Better the reconstruction \rightarrow more similar domains.

sults increases with 5 training images. The sketches now reflect distinct identities. Further increasing the number of training samples (10-shot) introduces more details for the sketch domain, and generates more diverse landscapes.

5. Conclusion and Limitations

We proposed to adapt a pretrained GAN learned on a large source domain to a small target domain by discovering cross-domain correspondences. While our method generates compelling results, it is not without limitations. Cars \rightarrow Abandoned cars in Fig. 6 depicts an example, where the color of the red car changes to orange in its abandoned form, likely because of the existence of an orange car (and no red one) in the 10 training images. FFHQ \rightarrow Sunglasses depicts another example, where a blonde hair turns dark with sunglasses. These show that there is a need for discovering better correspondence between the source and target domains, which will lead to more diverse generations. Nevertheless, we believe this work takes an important step towards creating more data-efficient generative models, demonstrating that existing source models can be leveraged in an effective way to model new distributions with less data.

Acknowledgements: We thank Daichi Ito for the beautiful caricatures. Part of the work was supported through NSF CAREER IIS-1751206 and Adobe Data Science Research Award.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Int. Conf. Comput. Vis.*, 2019.
- [2] Sagie Benaïm and Lior Wolf. One-sided unsupervised domain mapping. In *Adv. Neural Inform. Process. Syst.*, 2017.
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- [5] Alexey Dosovitskiy and Thomas Brox. Inverting visual representations with convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4829–4837, 2016.
- [6] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- [7] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
- [8] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Adv. Neural Inform. Process. Syst.*, 2017.
- [9] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.
- [10] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.
- [11] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Adv. Neural Inform. Process. Syst.*, 2020.
- [12] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [13] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. *arXiv preprint arXiv:1912.04958*, 2019.
- [14] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Int. Conf. Learn. Represent.*, 2015.
- [15] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [16] Yijun Li, Richard Zhang, Jingwan Lu, and Eli Shechtman. Few-shot image generation with elastic weight consolidation. In *Advances in Neural Information Processing Systems*, 2020.
- [17] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *Int. Conf. Comput. Vis.*, 2019.
- [18] Shaohui Liu, Xiao Zhang, Jianqiao Wang, and Jianbo Shi. Normalized diversification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [19] Qi Mao, Hsin-Ying Lee, Hung-Yu Tseng, Siwei Ma, and Ming-Hsuan Yang. Mode seeking generative adversarial networks for diverse image synthesis. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [20] Sangwook Mo, Minsu Cho, and Jinwoo Shin. Freeze discriminator: A simple baseline for fine-tuning gans. *arXiv preprint arXiv:2002.10964*, 2020.
- [21] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. In *International Conference on Machine Learning*, 2020.
- [22] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- [23] Atsuhiko Noguchi and Tatsuya Harada. Image generation from small datasets via batch statistics adaptation. In *Int. Conf. Comput. Vis.*, 2019.
- [24] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [25] Kuniaki Saito, Kate Saenko, and Ming-Yu Liu. Cocommit: Few-shot unsupervised image translation with a content conditioned style encoder. *arXiv preprint arXiv:2007.07431*, 2020.
- [26] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Adv. Neural Inform. Process. Syst.*, 2017.
- [27] Ngoc-Trung Tran, Tuan-Anh Bui, and Ngai-Man Cheung. Dist-gan: An improved gan using distance constraints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [28] Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, and Victor S Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *ICML*, volume 1, page 4, 2016.
- [29] Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. In *Neural Information Processing Systems (NeurIPS)*, 2020.
- [30] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Adv. Neural Inform. Process. Syst.*, 2016.
- [31] Xiaogang Wang and Xiaoou Tang. Face photo-sketch synthesis and recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(11):1055–1067, 2009.
- [32] 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 *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [33] Yaxing Wang, Salman Khan, Abel Gonzalez-Garcia, Joost van de Weijer, and Fahad Shahbaz Khan. Semi-supervised learning for few-shot image-to-image translation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [34] Yaxing Wang, Chenshen Wu, Luis Herranz, Joost van de Weijer, Abel Gonzalez-Garcia, and Bogdan Raducanu. Transferring gans: generating images from limited data. In *Eur. Conf. Comput. Vis.*, 2018.
- [35] Dingdong Yang, Seunghoon Hong, Yunseok Jang, Tianchen Zhao, and Honglak Lee. Diversity-sensitive conditional generative adversarial networks. *arXiv preprint arXiv:1901.09024*, 2019.
- [36] Jordan Yaniv, Yael Newman, and Ariel Shamir. The face of art: landmark detection and geometric style in portraits. *ACM Transactions on Graphics (TOG)*, 38(4):1–15, 2019.
- [37] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.
- [38] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.
- [39] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable augmentation for data-efficient gan training. *arXiv preprint arXiv:2006.10738*, 2020.
- [40] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Int. Conf. Comput. Vis.*, 2017.
- [41] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *Adv. Neural Inform. Process. Syst.*, 2017.

Supplementary

This document provides additional information about the proposed method. First, we continue the discussion on training and architecture details of our method and the baselines. We use the official pre-trained models for Church, Cars and Horses [13], and a 256 x 256 resolution pre-trained model for FFHQ.² The generator’s (G) and the discriminator’s (D_{img}) architecture are the same as StyleGAN2. For D_{pch} , we consider the first l layers of D_{img} , and convert the corresponding feature to a $N \times N$ output, where each member’s receptive field corresponds to a patch in the input image. We use Adam optimizer [14], and use the rest of the hyperparameters (e.g. learning rate) from [13]. The ideal training duration for adapting to different target domains is as follows: for domains whose appearance matches very closely with the source, i.e. babies/sunglasses with FFHQ as source, we get good results within 1000 iterations. For other face-based target domains (e.g. Modigliani’s paintings), we train our models for 5000 iterations to get decent results. For more complex domains (e.g. landscape drawings), we observe our best results at around 10000 iterations. Note that in eq. 4 (main paper), we have used $z \sim p_z(z)$ instead of $z \sim p_z(z) - Z_{anch}$ for the second term. This is because if an image from Z_{anch} is (globally) realistic, its patches will be realistic as well. Also note that the other way around is not true.

Details about MineGAN We use the publicly available code of MineGAN³ to produce its results (Fig. 4, Tab 1/2). However, since there is an involvement of an additional *miner* network during the adaptation process, we tried experimenting with smaller networks (than the default) for the extreme few-shot setting (10 training images). The results obtained were similar to the default setting – FID: 96.72, 68.67 for babies and sunglasses domain respectively, suggesting that reducing the complexity of the miner network alone is not sufficient to achieve better results.

FFHQ \rightarrow face domains Fig. 10 shows the real images used for different target domains in our experiments (apart from those presented in the main paper). Next, we show the results of translating a source model trained on natural faces (FFHQ) to different kinds of target domains. We observe the diversity in the generated images, which come as a result of preserving correspondence between the source and target distribution. Fig. 12 shows more examples for the idea discussed in Fig. 8 of the main paper. These four caricature images are unseen during the adaptation of a source model X (X is FFHQ/Church/Cars/Horses) to the caricature domain. We again observe that adapting FFHQ (natural faces) to the

caricature domain best embeds and reconstructs unseen images, indicating that caricature as a domain is most related to FFHQ than any other source domain. Fig. 13 presents an extension of Fig. 9 of main paper, where we study the 1-shot, 5-shot and 10-shot setting for two baselines, and compare it with our method. We notice that both FreezeD and EWC, overfit to the target sample in 1-shot setting, generating virtually identical sketches/scenes. This trend of overfitting for these baselines continues in 5/10-shot settings as well, where generations collapse to small variations around a few modes. Our method, on the other hand, takes the benefit of increasing training data size, by learning to generate more and more diverse samples, different from the images used for training.

Visualizing the clusters Fig. 14 visualizes the clustering-based diversity assessment introduced in Sec. 4.1 of the main paper. We group the generated images from a method into k clusters, with the k training images serving as the cluster center. After this, we study the resulting clusters, where we visualize how similar is (i) the closest member to the center (measured via LPIPS), (ii) the farthest member to the center. The intuition is that a method whose generated images overfit to the training data will result in clusters where the closest member is very similar to the corresponding cluster center. Each column deals with one cluster, where the cluster centers (real images used for training) are shown in the middle. The top half of the figure visualizes the closest members for different methods, whereas the bottom half visualizes the farthest ones. When no images get assigned to a cluster, the concept of closest/farthest members doesn’t apply, and we depict this with a red cross. Summarily, we observe that the closest members from TGAN/EWC are much more similar to the corresponding center than our method, whose even closest members are visually distinct. This observation also helps explain the better performance of our method compared to others in Table 2 (main paper).

Hand gestures experiment We find the property of emerging correspondences within seemingly unrelated source/target domains interesting, and hence for creativity purposes, take a further step to explore the idea. We collect images of arbitrary hand gestures being performed over a plain surface, and train a *source* model from scratch using that dataset. Next, we adapt it to various domains such as landscapes, fire, maps. During inference, we observe different aspects of the target domains a pair of hands can control (e.g. structure of river/islands). Please see our teaser video, which shows the correspondence results in this case, as well as better explains the benefits of our method in previously discussed scenarios (e.g. FFHQ \rightarrow caricatures, Church \rightarrow Van Gogh houses).

²We use <https://github.com/rosinality/stylegan2-pytorch>

³<https://github.com/yaxingwang/MineGAN/tree/master/styleGAN>

	Density	Coverage
TGAN [34]	0.379	0.250
TGAN+ADA [11]	0.434	0.285
FreezeD [20]	0.418	0.217
MineGAN [32]	0.803	0.125
EWC [16]	0.301	0.325
Ours	0.690	0.467

Table 4: Density (\uparrow) and Coverage (\uparrow) scores for FFHQ babies.

Precision and recall metrics A limitation of FID [8] is that it packs two aspects of the generated images, sample quality and diversity, into one score. This makes it difficult to disentangle and study the two properties separately. To overcome this, density and coverage metrics were proposed to evaluate the generative models [21]. In some feature space (e.g. CNN embeddings), density measures how many real-sample neighbourhood regions contain a fake sample. Coverage, in the same space, measures the ratio of real samples whose neighbourhood contains at least one fake sample. In both the definitions, *neighbourhood* is defined as a spherical region around a real sample, with its radius given by the distance from the next nearest real sample. A high score for both the metrics is preferred. Density is unbounded, whereas coverage is bounded by 1. We present evaluation of the baselines here using these metrics on FFHQ babies dataset in Table 4. We observe that MineGAN achieves a superior *density* score, i.e. quality of the generated image, but suffers in the *coverage* aspect. This is again an indication of mode collapse to a small number of high quality samples. Our method achieves a better balance between the quality as well as diversity of the generated samples. Note that this result is in alignment with the one presented in Table 2 (main paper), which studies diversity among the generated samples in a different way.

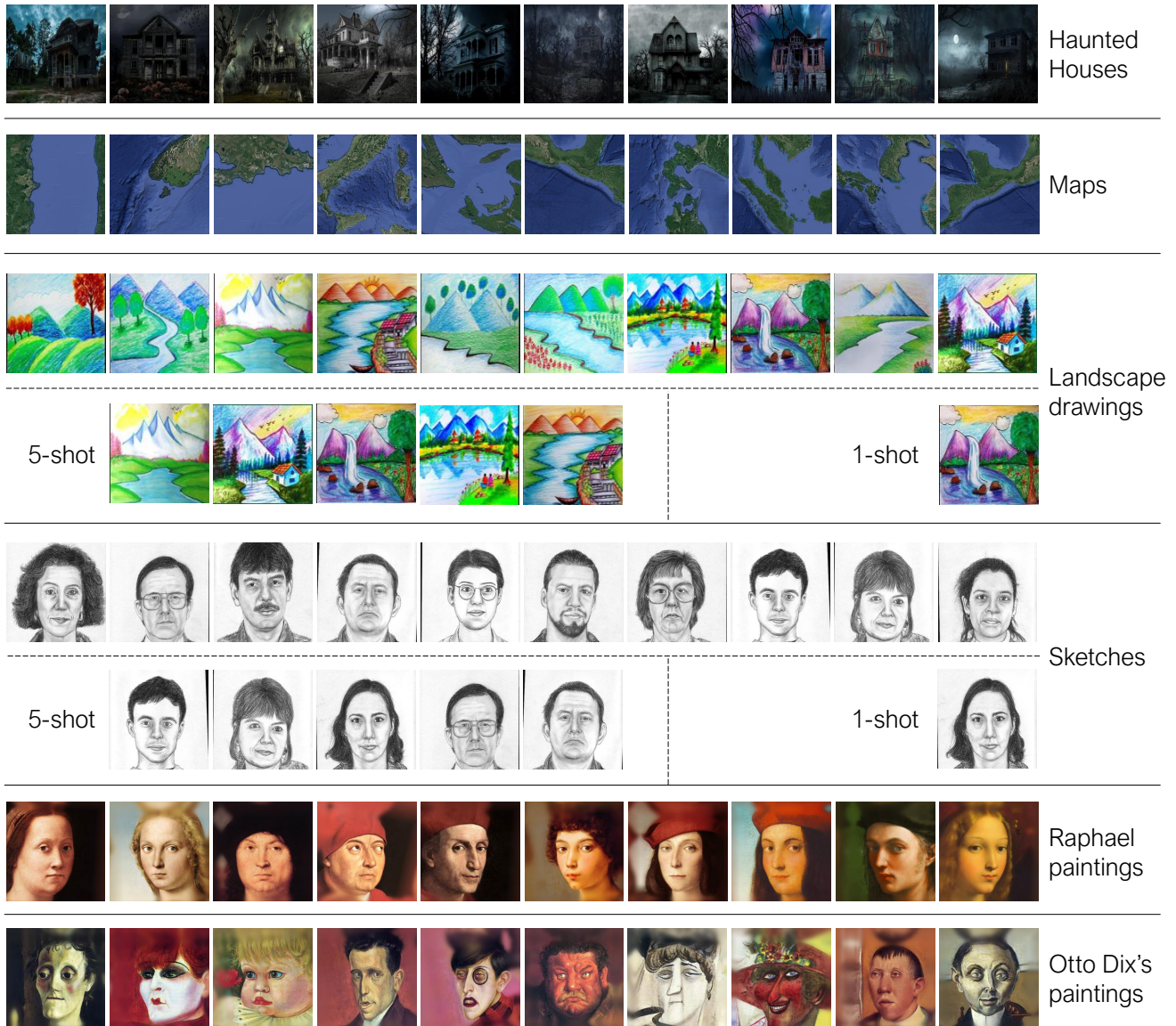


Figure 10: Real images used for translating a source model to different target domains. For landscape drawings and sketches, we have shown the images used in 1-shot and 5-shot scenario, for the experiment presented in Fig. 9 of the main paper, and Fig. 13 of this document.

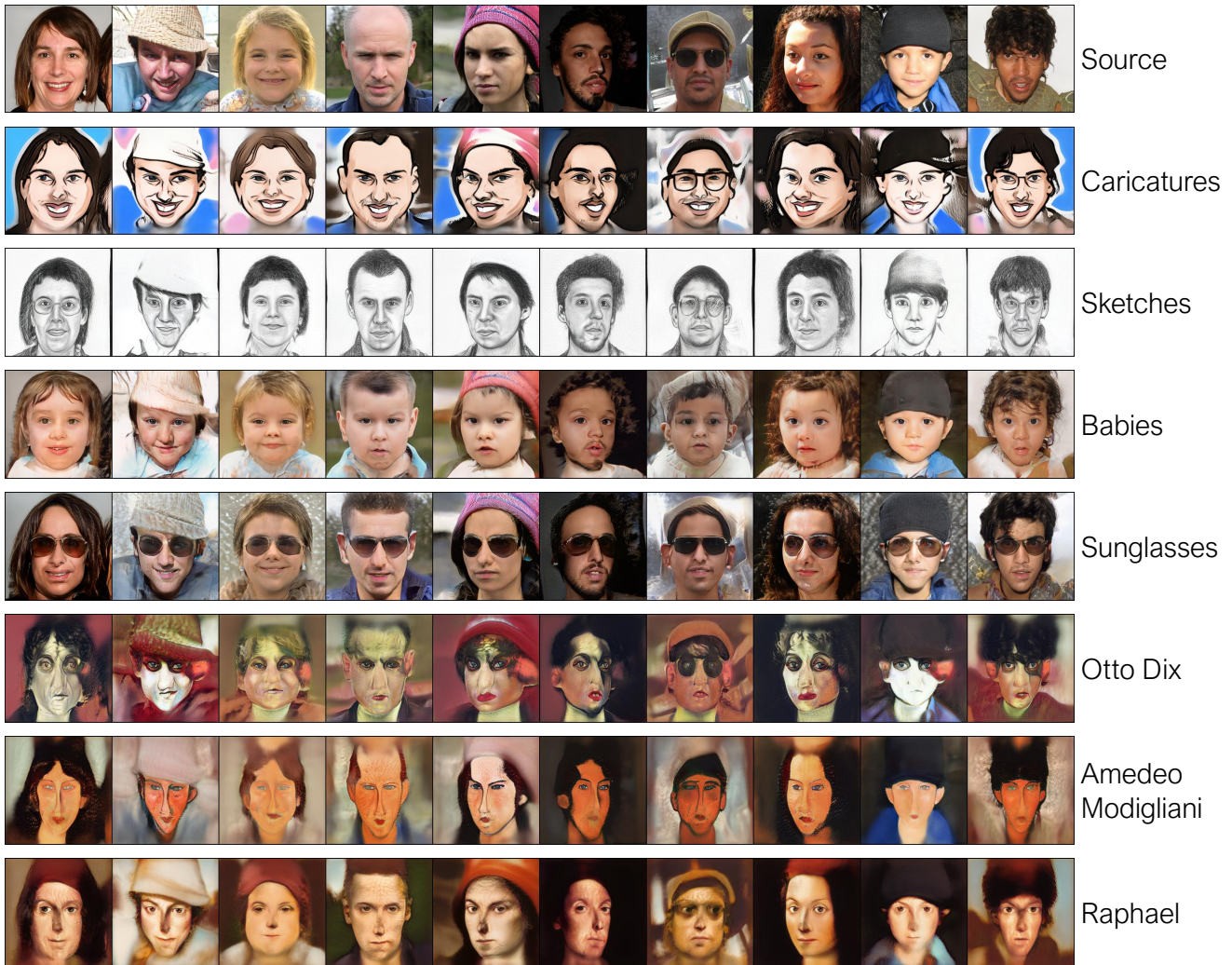


Figure 11: Translating the FFHQ source model to different target domains. The noise vector is kept same across the columns, so that we can study the relation between the corresponding source and target image.

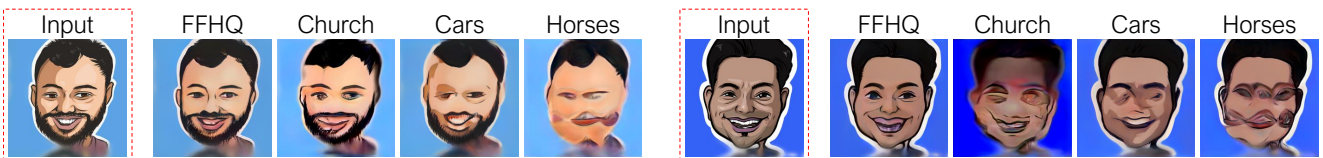


Figure 12: Embedding unseen caricature images into models adapted from different source to the same target domain (caricature). We observe that $\text{FFHQ} \rightarrow \text{caricatures}$ best captures the caricature properties, resulting in best reconstructions.

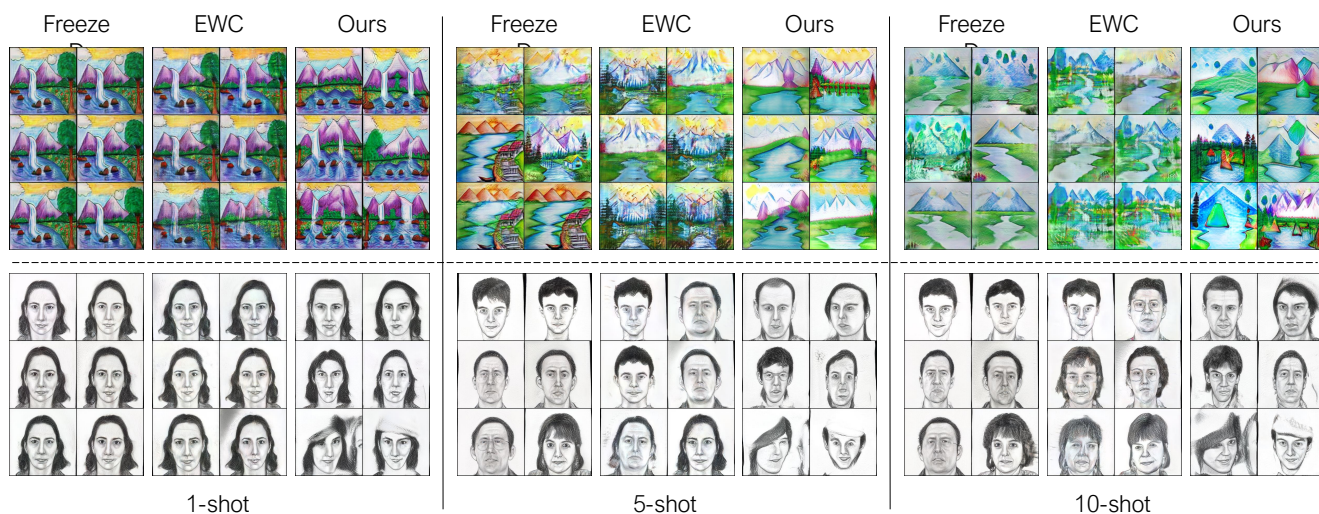


Figure 13: Comparison of our method compared to EWC [16] and FreezeD [20] in 1-shot, 5-shot and 10-shot setting.



Figure 14: Visualizing the clusters formed using the technique described in Sec. 4.1 of the main paper. The closest members produced by TGAN/EWC are much more similar to the corresponding cluster center than our method, indicating that the generations using the proposed method possess more diversity.