

Latent-CycleGAN for Domain Shift [Group 7]

Heyuan Chi
4742393

Scientific Computing

Yutong Zeng
4732327

Data and Computer Science

Abstract

Generative Adversarial Networks (GANs) have shown remarkable success in unpaired image-to-image translation tasks, exemplified by the CycleGAN framework. However, GAN-based methods often encounter issues such as mode collapse and insufficient preservation of fine-grained details. On the other hand, diffusion models exhibit robust structural fidelity but are limited by high computational overhead and slower inference speeds. To address these challenges, we propose Latent-CycleGAN, a novel architecture that integrates a pre-trained Stable Diffusion model into the CycleGAN framework. By leveraging LoRA-based fine-tuning and introducing additional latent-space discriminators alongside cycle-consistency constraints, our approach effectively combines the strengths of GAN- and diffusion-based methodologies. Experimental results demonstrate that Latent-CycleGAN produces higher-fidelity translations with enhanced detail retention, improved style consistency, and greater training stability compared to conventional CycleGAN. This work provides a promising new direction for unpaired image translation and broader applications in unsupervised domain adaptation.

All source code and additional information are available: <https://github.com/HeyuanChi/latent-cycle-gan>

1. Introduction

Image-to-image (I2I) conversion is the task of transforming an image from one domain to another while preserving key content. A famous example is image style transfer, where the goal is to change the style of an image without altering its underlying content structure. Early work by Gatys et al. [3] demonstrated that deep convolutional features could be disentangled into separate representations for content and style, making it possible to generate images that combine the content of one source with the style of another. This finding established the groundwork for style transfer and, more broadly, for domain shift problems in computer vision, which are often formulated as special cases of I2I

transformation.

In recent years, many new methods employing deep generative models have been investigated for I2I translation, including variational autoencoders (VAEs), generative adversarial networks (GANs), and diffusion models. Among them, GAN-based approaches are notably successful, thanks to adversarial objectives that guide the generator to produce outputs resembling real images in the target domain. Furthermore, recent developments in diffusion models show promise for stable training and producing diverse outputs, though with higher computational overhead. These advancements have opened up new possibilities, such as zero-shot editing—using powerful pretrained models to perform style transformations without additional fine-tuning.

However, the main problem of diffusion models is that it is hard to ensure the consistence of between target and source image. In details, the structure of object in image could be lost and the edge of object might be distorted. In this case, the process could be considered as create a new image rather than domain shift.

In this paper, we propose a novel image-to-image translation framework, Latent-CycleGAN, which enhances the classic CycleGAN by integrating a pre-trained Stable Diffusion model and employing LoRA-based fine-tuning. Compared with using diffusion model to generate target image directly, our method introduces additional latent space discriminators and cycle-consistency constraints alongside traditional image-space adversarial training, thereby preserving structural fidelity and ensuring style consistency. Extensive experiments on datasets such as vangogh2photo, horse2zebra, and cat2dog demonstrate that our approach significantly improves visual quality and detail retention. This work offers a new perspective for unsupervised domain adaptation and style transfer, setting the stage for further advancements in the field.

2. Related Works

2.1. GAN-based Image Translation

Traditional I2I methods generally rely on paired data, which could be expensive or impractical to collect. Cy-

cleGAN [23] addressed this by introducing a cycle consistency loss, ensuring that an image translated from domain A to B and then back to A remains unchanged. Subsequent dual-generator architectures, such as DiscoGAN [11] and DualGAN [22], adopted similar cycle constraints. UNIT [15] further proposed a shared latent space for unsupervised translation. Although these GAN-based methods are effective, they often encounter issues like mode collapse and may not preserve fine-grained details when the domains differ significantly in semantics or style.

Still, GAN-based models are widely used and have shown impressive results in tasks such as photo coloring, image-to-image conversion, and style transformation. For instance, pix2pix [10] used paired training images, while CycleGAN [23] succeeded even without paired data by leveraging a cyclic loss. Enhanced techniques—like StarGAN [2] for multi-domain translation, BicycleGAN [25] for multi-modal outputs, and frameworks like MUNIT [9] or DRIT [14] for decomposing content and style—have made GAN-based methods increasingly flexible and controllable.

2.2. Diffusion-Based Approaches

Denoising Diffusion Probabilistic Models (DDPMs) [6] have emerged as a strong alternative to GANs, capable of generating high-fidelity images by reversing a progressive noising process. Improvements such as DDIM [20] and classifier-free guidance [7] further enhance generation speed and quality. More recently, latent diffusion models, as exemplified by Stable Diffusion [17], reduce computational cost by operating in a compressed latent space. These models can be conditioned on text or images to produce new outputs that preserve source content while adopting the desired style. Compared with GANs, diffusion methods tend to be more stable during training and yield diverse outputs, but often require longer sampling times.

2.3. Combining Diffusion and GANs

Recent efforts have explored hybrid approaches that integrate diffusion models with GANs to exploit the strengths of both. Methods like VQ-Diffusion [4] and DiffusionGAN [21] incorporate diffusion priors to regularize GAN outputs, improving training stability and structural consistency. Along this line, our work inserts a latent diffusion module into the CycleGAN framework, using a pretrained Stable Diffusion UNet as a noise predictor in the latent space. This additional constraint enforces structural integrity beyond pixel-level losses.

2.4. Low-Rank Adaptation

Large models such as Stable Diffusion contain billions of parameters, making full fine-tuning infeasible in many applications. Low-Rank Adaptation (LoRA) [8] addresses

this issue by inserting low-rank matrices into the attention layers of pretrained networks, enabling task-specific adaptation with minimal computational overhead. In our framework, we adopt LoRA to fine-tune the diffusion UNet modules, retaining the powerful generative prior of the original model while achieving effective, domain-specific customization.

3. Method

In this section, we present our method in detail. Building upon the CycleGAN framework [24], we incorporate two pretrained Stable Diffusion models [18] (fine-tuned by LoRA [8]), and propose novel latent-space constraints to improve the translation quality.

3.1. CycleGAN (baseline)

CycleGAN [24] aims to learn a mapping between two image domains, denoted as A and B . It does not require to use paired training data. As shown in the middle of figure 1, the framework includes two generators and two discriminators:

- G_A : a generator translating images from domain A to domain B .
- G_B : a generator translating images from domain B to domain A .
- D_A : a discriminator distinguishing real B -domain images from translated (fake) images produced by G_A .
- D_B : a discriminator distinguishing real A -domain images from translated (fake) images produced by G_B .

3.1.1 Generator Losses

1. Image-Space Adversarial Loss

CycleGAN uses adversarial objectives in the image space. For generator G_A , which maps $x \in A$ to $\tilde{x} = G_A(x) \in B$, the corresponding discriminator D_A tries to classify real images $y \sim B$ as real and \tilde{x} as fake. In a least-squares GAN (LSGAN) formulation [16]:

$$\mathcal{L}_{\text{GAN}}(G_A, D_A) = \mathbb{E}_{x \sim A} [(D_A(G_A(x)) - 1)^2].$$

Similarly for G_B and D_B :

$$\mathcal{L}_{\text{GAN}}(G_B, D_B) = \mathbb{E}_{y \sim B} [(D_B(G_B(y)) - 1)^2].$$

2. Cycle-Consistency Loss

To ensure learned mappings are reversible, CycleGAN enforces

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G_A, G_B) = & \mathbb{E}_{x \sim A} [\|G_B(G_A(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim B} [\|G_A(G_B(y)) - y\|_1]. \end{aligned}$$

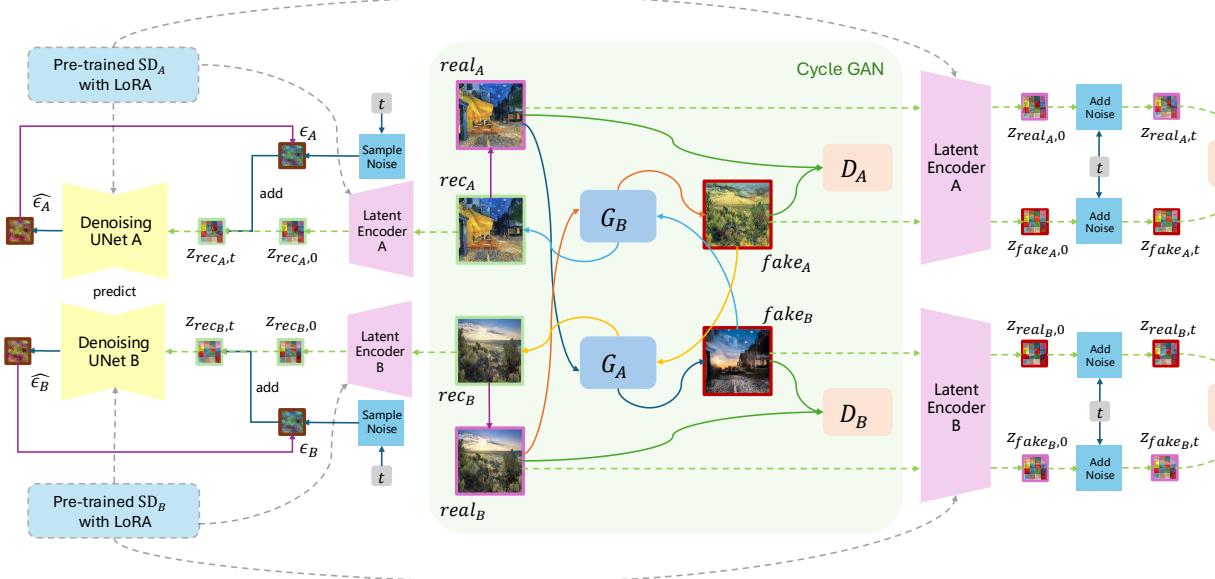


Figure 1: **Method Overview.** We start with the standard CycleGAN setup: two generators (G_A, G_B) translating between domains A and B , and two discriminators (D_A, D_B) for real-vs-fake image supervision. We then integrate a pretrained Stable Diffusion (SD) model [18] (optionally fine-tuned via LoRA [8]) and introduce two new latent-space discriminators (D_{LA}, D_{LB}). Each image is encoded by the SD VAE into latent space, partially noised at a small diffusion timestep, and fed to the latent discriminators. This combined image- and latent-level training enforces cycle consistency in both spaces, yielding more faithful and style-consistent image translations.

This penalizes large deviations between the reconstructed image and the original.

3. Identity Loss

An optional identity loss [24] can be used to encourage generators to preserve color or other low-level features:

$$\begin{aligned} \mathcal{L}_{id}(G_A, G_B) = & \mathbb{E}_{y \sim B} [\|G_A(y) - y\|_1] \\ & + \mathbb{E}_{x \sim A} [\|G_B(x) - x\|_1]. \end{aligned}$$

3.1.2 Discriminator Losses

Each discriminator is trained with a real vs. fake objective. For D_A ,

$$\mathcal{L}_{D_A} = \frac{1}{2} \left[\mathbb{E}_{y \sim B} [(D_A(y) - 1)^2] + \mathbb{E}_{x \sim A} [(D_A(G_A(x)))^2] \right],$$

and similarly for D_B ,

$$\mathcal{L}_{D_B} = \frac{1}{2} \left[\mathbb{E}_{x \sim A} [(D_B(x) - 1)^2] + \mathbb{E}_{y \sim B} [(D_B(G_B(y)))^2] \right].$$

3.2. Stable Diffusion and LoRA

In our work, we incorporate Stable Diffusion (SD) [18] as a powerful pretrained model for guidance in the latent space. Unlike conventional image-to-image translation

methods, SD encodes images into a latent representation and applies a denoising UNet to generate high-quality samples.

The core of Stable Diffusion models consists of:

- A Variational Autoencoder (VAE): encodes an input image into a latent code and decodes a latent code back into the image space.
- A Denoising UNet: trained via a diffusion process to remove noise from latent representations at various timesteps.

When an image x is provided, the VAE encoder produces a latent code $z = \text{VAE.encode}(x)$. During sampling or inference, the UNet iteratively denoises a noisy latent until it converges on a clean sample, which the VAE decoder then maps back to image space.

To adapt the pretrained SD model to specific styles or domains, we can apply LoRA [8], a parameter-efficient fine-tuning technique. LoRA injects low-rank adaptations into the weight matrices of the UNet, reducing the number of additional trainable parameters while preserving most of the original SD capabilities. This fine-tuning step allows SD to capture domain-specific cues without overfitting or forgetting the general knowledge from the original training.

In practice, we freeze all of the original Stable Diffusion parameters—namely, the main UNet convolutional weights,

normalization layers, and the VAE encoder/decoder—and only insert LoRA adapters [8] (with rank $r = 8$) into the cross-attention blocks. Because we rely on text prompts for domain-specific conditioning, cross-attention is pivotal in how textual context is integrated into the latent denoising process. By fine-tuning just these cross-attention layers, we can effectively adapt the model to new styles without sacrificing the broader, pretrained representations that maintain overall stability and quality. Restricting the learnable parameters to these low-rank LoRA matrices also keeps memory overhead low, reduces the risk of overfitting, and aligns with the original LoRA rationale, enabling faster, more stable training.

3.3. Latent-CycleGAN: Our Proposed Method

We now introduce **Latent-CycleGAN** as shown in figure 1, which augments the standard CycleGAN training with latent-space discriminators and latent cycle-consistency constraints. This leverages Stable Diffusion (optionally fine-tuned with LoRA) to better guide the translation process at the latent level.

In addition to the image-space discriminators D_A and D_B , we introduce new latent discriminators:

- D_{LA} : operates in the SD latent space for domain- B images, evaluating real vs. fake latents.
- D_{LB} : operates in the SD latent space for domain- A images, evaluating real vs. fake latents.

These discriminators aim to enforce that the translated samples (once encoded by the SD VAE) remain realistic in the diffusion latent space.

Unlike conventional CycleGAN, which focuses solely on pixel-level fidelity, our approach applies adversarial and cycle-consistency constraints within the latent space learned by Stable Diffusion. This latent space often encodes more abstract, high-level features about structure and content. By enforcing that translated images produce realistic latent codes, we reduce mode collapse and encourage the model to preserve key semantic details. Similarly, the latent cycle-consistency loss ensures that once we re-encode and diffuse translated samples, we can still recover the original noise distributions. This can help stabilize training and yields higher-fidelity transformations, especially for domains with large style gaps.

3.3.1 Latent Generator Losses

1. Latent Adversarial Loss

For $x \in A$, let $\tilde{y} = G_A(x)$. We encode both real y (from B) and fake \tilde{y} with the SD VAE, obtaining $z_{\text{real}} = \text{VAE.encode}(y)$ and $z_{\text{fake}} = \text{VAE.encode}(\tilde{y})$. To add mild stochasticity, we choose a small diffusion

timestep t and inject noise according to the forward diffusion equation:

$$z_{\text{real}}^{(t)} = \sqrt{\alpha_t} z_{\text{real}} + \sqrt{1 - \alpha_t} \epsilon, \epsilon \sim \mathcal{N}(0, \mathbf{I}),$$

and analogously for $z_{\text{fake}}^{(t)}$. The latent discriminator D_{LA} is trained with:

$$\mathcal{L}_{\text{GAN-latent}}(G_A, D_{LA}) = \mathbb{E}[(D_{LA}(z_{\text{fake}}^{(t)}) - 1)^2].$$

Similarly for G_B and D_{LB} .

To introduce moderate but increasing noise during training, we randomly select a diffusion timestep $t \in [0, t_{\max}/4 \times \text{training_iteration_ratio}]$ where `training_iteration_ratio` is the fraction of training completed (`current_iteration / total_iterations`). At the start of training the ratio close to 0, so t stays near 0. This keeps the latent representation almost noise-free, helping the model learn basic structural and semantic cues. As training progresses the ratio increases, allowing t to grow—up to $t_{\max}/4$ at the end of training. By capping t at $t_{\max}/4$, we avoid excessive noise that could destabilize training, yet still provide enough variation to make the latent-space discriminators robust. Consequently, the model experiences a gentle shift from near noise-free conditions to mildly corrupted latents, striking a balance between preserving key details and learning to handle moderate uncertainty.

2. Latent Cycle-Consistency Loss

To maintain consistency through the two-stage translation, we apply the SD forward diffusion to the reconstructed images and then measure alignment in noise space. For example, let $\hat{x}_A = G_B(G_A(x))$. The VAE encodes \hat{x}_A into $z_{\hat{x}_A}$, which is then diffused to timestep t . The diffusion UNet predicts the noise ϵ_{pred} , which we compare to the actual noise ϵ added:

$$\mathcal{L}_{\text{cyc-lat-A}} = \mathbb{E}[\|\hat{\epsilon}(z_{\hat{x}_A}) - \epsilon_A\|_2^2],$$

with an analogous term for \hat{x}_B ,

$$\mathcal{L}_{\text{cyc-lat-B}} = \mathbb{E}[\|\hat{\epsilon}(z_{\hat{y}_B}) - \epsilon_B\|_2^2].$$

3.3.2 Latent Discriminator Losses

The discriminators D_{LA} and D_{LB} are optimized similarly to standard LSGAN:

$$\mathcal{L}_{D_{LA}} = \frac{1}{2} [\mathbb{E}((D_{LA}(z_{\text{real}_A}^{(t)}) - 1)^2) + \mathbb{E}((D_{LA}(z_{\text{fake}_A}^{(t)}))^2)],$$

$$\mathcal{L}_{D_{LB}} = \frac{1}{2} [\mathbb{E}((D_{LB}(z_{\text{real}_B}^{(t)}) - 1)^2) + \mathbb{E}((D_{LB}(z_{\text{fake}_B}^{(t)}))^2)].$$

3.3.3 Full Objective

Finally, we combine the image-space and latent-space losses with weighting factors. Denote λ_{GAN} , λ_{cyc} , λ_{id} , α_{GAN} , α_{cyc} . The total generator objective is:

$$\begin{aligned}\mathcal{L}_G = & \lambda_{\text{GAN}} (\mathcal{L}_{\text{GAN}}(G_A, D_A) + \mathcal{L}_{\text{GAN}}(G_B, D_B)) \\ & + \lambda_{\text{cyc}} \mathcal{L}_{\text{cyc}}(G_A, G_B) + \lambda_{\text{id}} \mathcal{L}_{\text{id}}(G_A, G_B) \\ & + \alpha_{\text{GAN}} (\mathcal{L}_{\text{GAN-lat}}(G_A, D_{LA}) + \mathcal{L}_{\text{GAN-lat}}(G_B, D_{LB})) \\ & + \alpha_{\text{cyc}} (\mathcal{L}_{\text{cyc-lat-A}}(G_A, G_B) + \mathcal{L}_{\text{cyc-lat-B}}(G_B, G_A)).\end{aligned}$$

Each discriminator has a matching real-vs-fake loss:

$$\mathcal{L}_D = \{\mathcal{L}_{D_A}, \mathcal{L}_{D_B}, \mathcal{L}_{D_{LA}}, \mathcal{L}_{D_{LB}}\}.$$

3.3.4 Implementation Details

We implement our approach in PyTorch, building upon the original CycleGAN codebase [24]. The generators (G_A , G_B) and image-space discriminators (D_A , D_B) use ResNet or U-Net backbones. The latent discriminators (D_{LA} , D_{LB}) are adapted from PatchGAN [10] to operate on four-channel latent inputs from Stable Diffusion’s VAE. The SD model can remain frozen or be partially updated using LoRA [8] to capture domain-specific styles. We randomly sample a diffusion timestep t from a small range (early in the diffusion schedule) so that some semantic information remains while still challenging the latent discriminators with mild noise.

4. Experiments

This section details our experimental setup and evaluation. We begin with an overview of the datasets and evaluation metrics, followed by implementation and training specifics. We then present results from the LoRA fine-tuning baseline, and proceed to a comprehensive analysis of our proposed Latent-CycleGAN framework under varying hyperparameter configurations. In the end, we conclude with additional experiments, a full set of quantitative metrics, and qualitative comparisons.

4.1. Datasets

We primarily evaluate the proposed Latent-CycleGAN on three datasets. The vangogh2photo and horse2zebra datasets are sourced from CycleGAN [24], while cat2dog is adopted from DRIT [14]. The amounts of the classes for training are 6,287 (photo), and 400 (vangogh) in vangogh2photo; and 1,067 (horse), 1,334 (zebra) in horse2zebra. And for testing are : 751 (photo), and 400 (vangogh) in vangogh2photo; 120 (horse), 140 (zebra) in horse2zebra. The amounts for training are 871 (cat) and 1,364 (dog). And for testing are 100 (cat) and 100 (dog). All images are resized to 256×256 for training.

4.2. Evaluation Metrics

To evaluate the performance of our method, we employ Kernel Inception Distance (KID) as the quantitative measure. KID introduced by Bińkowski et al.[1], is an unbiased alternative to the widely used FID [5]. FID often produces biased estimates, especially with smaller datasets, since it relies on Gaussian assumptions and covariance matrix inversion. In contrast, KID adopts a polynomial-kernel MMD, making it distribution-free, more robust in low-data scenarios, and numerically simpler. Empirical studies show that KID correlates with FID on large samples but outperforms it with limited data[13][19]. Thus, KID is recommended when evaluating generative models with small datasets.

4.3. Training Details

All experiments were conducted on a server equipped with One NVIDIA A100-SXM4-80GB GPU and 120 GB of system RAM (AutoDL Cloud Server). We implemented our models in PyTorch using Python 3.8. The Stable Diffusion VAE and UNet components were loaded from CompVis Stable Diffusion v1.5[18] and remained partially frozen unless otherwise stated.

We used the Adam optimizer [12] ($\beta_1 = 0.5$, $\beta_2 = 0.999$) for all networks and applied standard data augmentation (random cropping, flipping) for both the A and B domains. Batch size was set to 16 with an initial learning rate of 1e-4 (decayed after some epochs to zero).

4.4. LoRA Fine-Tuning

To explore the benefit of stylistic adaptation within the diffusion model, we incorporated LoRA adapters [8] into the Stable Diffusion UNet. Specifically, we inserted low-rank adapters ($r = 8$) into each cross-attention layer to capture domain-specific style cues for images. We fine-tuned LoRA parameters on each dataset for 5000 epochs with a reduced learning rate of 1e-5, while freezing most of the pretrained diffusion weights.

Figure 2 compares outputs on the vangogh2photo dataset, generated by both the original Stable Diffusion model and our LoRA-adapted SD model. We observe that the original SD typically produces overly uniform Van Gogh-style images, lacking sufficient generative diversity. Furthermore, its photo outputs often diverge substantially from the training data in terms of structure and realism. In contrast, the LoRA-enhanced diffusion model effectively captures domain-specific style cues—such as the characteristic brushstrokes and color palettes of Van Gogh—while maintaining closer adherence to the photo domain’s distribution. These LoRA-enhanced diffusion modules are then integrated into our final pipeline to serve as the latent encoder/decoder backbone.

By fine-tuning the original SD models through the LoRA method, we obtain models with stronger generalization in the corresponding field. In this way, the LoRA-adapted SD can be used to guide CycleGAN in later stages, facilitating more faithful domain-to-domain translations.

4.5. Latent-CycleGAN

First, we trained the models on vangogh2photo dataset. We tested different hyperparameter configurations for our Latent-CycleGAN framework. Then we did additional experiments on other datasets.

4.5.1 Qualitative Observations

We recall that our method introduces two additional weighting factors: α_{GAN} for latent-space adversarial loss and α_{cyc} for latent cycle-consistency.

In our first experiment, we completely remove the original losses by setting

$$\lambda_{\text{GAN}} = 0, \lambda_{\text{cyc}} = 0, \lambda_{\text{id}} = 0,$$

and retaining only the latent-space adversarial and cycle-consistency terms. Specifically, we use

$$\alpha_{\text{GAN}} = 1, \alpha_{\text{cyc}} = 0.1.$$

We refer to this configuration as Set A. Unlike the standard CycleGAN, our goal is to investigate whether latent-space losses alone are sufficient to train a successful model.

To this end, we initialize the model using a pretrained CycleGAN [24] trained for 200 epochs (approximately 100K iterations), and continue training for an additional 50 epochs (around 25K iterations) using the parameters defined above. Qualitative results are presented in Figure 3.

When only the latent adversarial and cycle-consistency losses are used, the generated images exhibit noticeable artifacts such as holes, grid-like distortions, and structural inconsistencies. We hypothesize that this degradation arises from the absence of pixel-level supervision: the image-space discriminators in standard CycleGAN see actual images and can directly constrain local details and texture fidelity. By removing them, the model relies solely on how well the latent-space discriminators recognize encoded features with added noise. This indirect supervision is not sufficient to preserve fine structure, leading to more pronounced spatial corruption in the outputs.

Motivated by these artifacts, we next investigate how restoring the original image-space losses, together with the latent-space losses, affects the results. In particular, we consider two additional sets of parameters (which we call Set B and Set C) to assess the combined influence of latent and image-space constraints.

- Set B: $\alpha_{\text{GAN}} = 1, \alpha_{\text{cyc}} = 0.1; \lambda_{\text{GAN}} = 1, \lambda_{\text{cyc}} = 10, \lambda_{\text{id}} = 10.$

- Set C: $\alpha_{\text{GAN}} = 0.1, \alpha_{\text{cyc}} = 0.01; \lambda_{\text{GAN}} = 1, \lambda_{\text{cyc}} = 10, \lambda_{\text{id}} = 10.$

Both Set B and Set C incorporate the original CycleGAN losses ($\lambda_{\text{GAN}}, \lambda_{\text{cyc}}, \lambda_{\text{id}}$) and the latent-space losses ($\alpha_{\text{GAN}}, \alpha_{\text{cyc}}$). They differ in the relative weights assigned to the latent components. Again, we initialize from the same pretrained CycleGAN checkpoint and train for 50 epochs under these new parameter settings.

From the results of translating real photos to Van Gogh style in Figure 3, sets B and C show only marginal improvements over the baseline. Among these, set B exhibits more accurate colors, making it slightly better than the baseline. When translating Van Gogh paintings to real photos, both sets B and C outperform the baseline, with set B once again producing more realistic color rendering. In contrast, set C, which completely removes the original loss, performs quite poorly. It suggests that the latent loss can provide assisted guidance, but it cannot fully replace the original loss.

These results confirm that latent losses alone (Set A) are insufficient to produce high-quality images. However, when used in conjunction with the original image-space constraints (Sets B and C), they can further enhance style-transfer fidelity and color consistency.

Figure 4 compares the training losses. In the top plot, the standard CycleGAN converges steadily after 200 epochs, whereas in the bottom plot (Latent-CycleGAN), we see that Set A (only latent-space losses) fails to stabilize, remaining at a higher loss level over 50 epochs. This verifies our qualitative observation that removing all image-space losses leads to insufficient supervision. Meanwhile, SetB and Set C show more stable convergence under combined image-space and latent-space constraints, further corroborating our conclusion that latent losses alone are not enough, but can serve as an effective complement to the original CycleGAN objectives.

4.5.2 Additional Experiments

To further validate the robustness of our approach, we also conducted experiments on other datasets (horse2zebra and cat2dog). All hyperparameters follow the same setting as in the set C. Figure 5 illustrates some translated examples.

4.5.3 Quantitative Comparison

We computed the aforementioned quantitative metric KID on each dataset and compared our Latent-CycleGAN against the standard CycleGAN models (baseline).

Table 1 presents the comprehensive metrics on different datasets. Our models achieved a lower KID except on each dataset for vangogh2photo. From these experiments, we observe that incorporating LoRA fine-tuning and latent-space adversarial/cycle-consistency consistently improves

both perceptual quality and style accuracy across diverse domains.

Table 1: Comparison of Kernel Inception Distances on different datasets. Lower is better.

| Dataset | Baseline | Latent-CycleGAN |
|---------------|------------------|------------------|
| vangogh2photo | 0.03 ± 0.004 | 0.03 ± 0.004 |
| photo2vangogh | 0.03 ± 0.005 | 0.03 ± 0.005 |
| horse2zebra | 0.09 ± 0.006 | 0.08 ± 0.005 |
| zebra2horse | 0.13 ± 0.004 | 0.13 ± 0.004 |
| cat2dog | 0.23 ± 0.005 | 0.17 ± 0.005 |
| dog2cat | 0.16 ± 0.003 | 0.13 ± 0.003 |

Based on these results, we conclude that LoRA-based fine-tuning and latent-space adversarial training together enhance both the fidelity and the style-consistency of unpaired image translation across multiple datasets. This highlights the general applicability of our proposed framework.

5. Ablation Study

To investigate the contribution of each component in our Latent-CycleGAN framework, we perform three ablation experiments on the vangogh2photo dataset. We systematically remove or replace individual elements and then evaluate the models using the KID metric. Table 2 summarizes the quantitative results: 1. Remove LoRA. We replace the LoRA-fine-tuned Stable Diffusion model with the original Stable Diffusion, using the same latent losses; 2. Remove Latent Adversarial Loss. Here, we disable the latent adversarial loss, retaining only the latent cycle-consistency term; 3. Remove Latent Cycle-Consistency. Conversely, we remove only the latent cycle-consistency loss while keeping the latent adversarial term.

Table 2 shows that each ablation degrades the translation performance compared with our full Latent-CycleGAN. And the qualitative results are shown in Figure 6.

Table 2: Ablation results (Kernel Inception Distances) on vangogh2photo. Lower KID indicates better performance.

| Model | vangogh2photo | photo2vangogh |
|------------------------------------|-------------------|-------------------|
| Latent-CycleGAN | 0.028 ± 0.004 | 0.029 ± 0.005 |
| w/o LORA | 0.037 ± 0.004 | 0.031 ± 0.005 |
| w/o $\mathcal{L}_{\text{GAN-lat}}$ | 0.028 ± 0.004 | 0.029 ± 0.005 |
| w/o $\mathcal{L}_{\text{cyc-lat}}$ | 0.037 ± 0.004 | 0.032 ± 0.005 |

From the table and the figure, we can see that removing LoRA or removing the latent cycle-consistency loss both lead to noticeably higher KID scores in the vangogh2photo direction, indicating that these two components are particularly important for improving style-to-real translation. On

the other hand, removing the latent adversarial loss yields results that are almost identical to the full model, suggesting that for this dataset, cycle consistency and LoRA fine-tuning contribute more critically to enhancing translation quality.

6. Conclusion

In this paper, we presented **Latent-CycleGAN**, a framework that augments CycleGAN by incorporating additional latent-space constraints from a pretrained Stable Diffusion model with LoRA. Our experiments show that:

Strengths

1. Integrating latent-space constraints helps retain structural details, particularly for challenging cross-domain translations.
2. LoRA fine-tuning offers a practical way to adapt large diffusion models to specific style domains without excessive computational overhead.
3. Compared with vanilla CycleGAN, qualitative results often reveal better texture consistency and reduced mode collapse, even when KID gains are modest.

Limitations

1. On tasks with smaller domain gaps (e.g. horse2zebra), improvements are not always pronounced, indicating latent-space constraints may be less critical when CycleGAN alone suffices.
2. When the target domain is dominated by scenic imagery, the model often “blends” facial or bodily features into landscape-like textures. As a result, human portraits may be forced into an outdoor or environmental context, leading to unnatural merges and distorted facial details (shown in figure 7).

Future Work To address these limitations, we plan to:

1. Incorporate face or body detection/segmentation priors when translating human portraits into domains dominated by landscape imagery. This may prevent facial features from being forcefully reinterpreted as background and help maintain clearer structural details.
2. Explore more advanced or localized latent adversarial losses to better capture fine-grained features without over-penalizing irrelevant regions.
3. Scale up to higher-resolution datasets and fine-tune the balance between image- and latent-space constraints to improve both detail fidelity and training efficiency.

By uniting the simplicity and unpaired learning advantage of CycleGAN with the powerful latent-space representations of Stable Diffusion, Latent-CycleGAN presents a promising direction for advanced image-to-image translation and broader unsupervised domain adaptation tasks.

References

- [1] Mikołaj Bińkowski, Dougal J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations (ICLR)*, 2018. [5](#)
- [2] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8789–8797, 2018. [2](#)
- [3] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016. [1](#)
- [4] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10696–10706, 2022. [2](#)
- [5] 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 *Advances in Neural Information Processing Systems (NIPS)*, pages 6626–6637, 2017. [5](#)
- [6] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [2](#)
- [7] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. [2](#)
- [8] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhu Li, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models, 2021. *arXiv preprint arXiv:2106.09685*. [2, 3, 4, 5](#)
- [9] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *European Conference on Computer Vision (ECCV)*, pages 172–189, 2018. [2](#)
- [10] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5967–5976, 2017. [2, 5](#)
- [11] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. In *International conference on machine learning*, pages 1857–1865. Pmlr, 2017. [2](#)
- [12] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [5](#)
- [13] Karol Kurach, Mario Lucic, Xiaohua Zhai, Marcin Michalski, and Sylvain Gelly. A large-scale study on regularization and normalization in GANs. In *International Conference on Machine Learning (ICML)*, pages 3581–3590, 2019. [5](#)
- [14] Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Diverse image-to-image translation via disentangled representations. In *European Conference on Computer Vision (ECCV)*, pages 36–52, 2018. [2, 5](#)
- [15] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 700–708, 2017. [2](#)
- [16] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2813–2821, 2017. [2](#)
- [17] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. [2](#)
- [18] 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 (CVPR)*, pages 10684–10695, 2022. [2, 3, 5](#)
- [19] Mehdi S. M. Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5228–5237, 2018. [5](#)
- [20] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. [2](#)
- [21] Zhendong Wang, Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Diffusion-gan: Training gans with diffusion. *arXiv preprint arXiv:2206.02262*, 2022. [2](#)
- [22] Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *Proceedings of the IEEE international conference on computer vision*, pages 2849–2857, 2017. [2](#)
- [23] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2223–2232, 2017. [2](#)
- [24] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2223–2232, 2017. [2, 3, 5, 6](#)
- [25] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 465–476, 2017. [2](#)

A. Additional Figures



(a) Generated Van Gogh's paints by SD-v1.5



(b) Generated Van Gogh's paints by LoRA-SD.

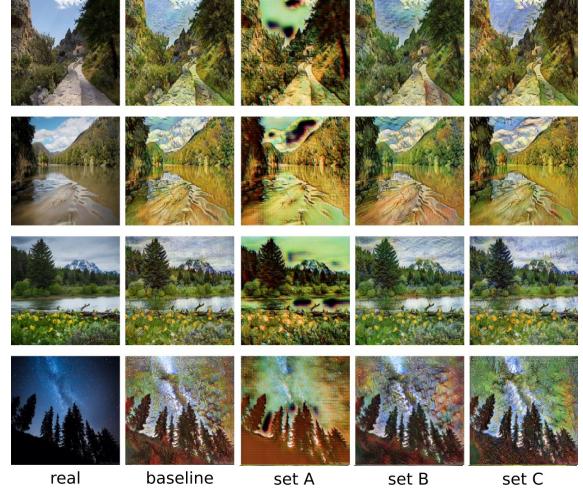


(c) Generated real photos by SD-v1.5.



(d) Generated real photos by LoRA-SD.

Figure 2: Comparison of outputs on the vangogh2photo dataset using original Stable Diffusion v1.5 vs. LoRA-SD. Rows (a) and (b) highlight Van Gogh-style images, while rows (c) and (d) show real-photo-style images. We observe that LoRA-SD models capture domain-specific style cues more effectively, improving stylistic fidelity and overall image quality.



(a) Example of translation results from real photos to Van Gogh style.



(b) Example of translation results from Van Gogh's Paintings to real style.

Figure 3: Visualization results of image translation based on the Latent-CycleGAN framework.

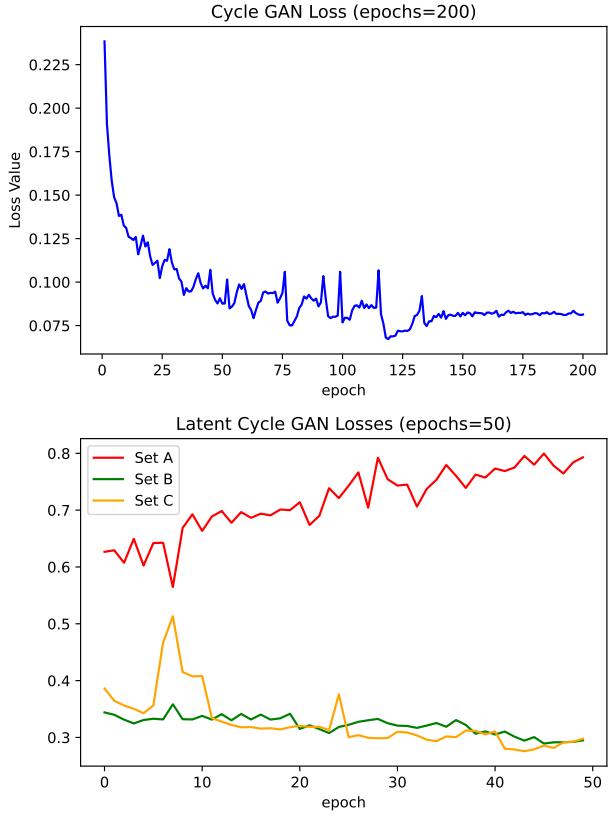


Figure 4: The training loss curve of a standard CycleGAN model, converging steadily after 200 epochs (around 100K iterations). The training losses of our Latent-CycleGAN variations (Set A, Set B, and Set C) over 50 epochs (about 25K iterations). Observe that Set A (red) diverges or fluctuates at a higher loss level, suggesting insufficient pixel-level supervision, whereas Sets B (green) and C (orange) exhibit relatively lower, more stable losses, reflecting the benefit of combining original image-space and latent-space constraints.

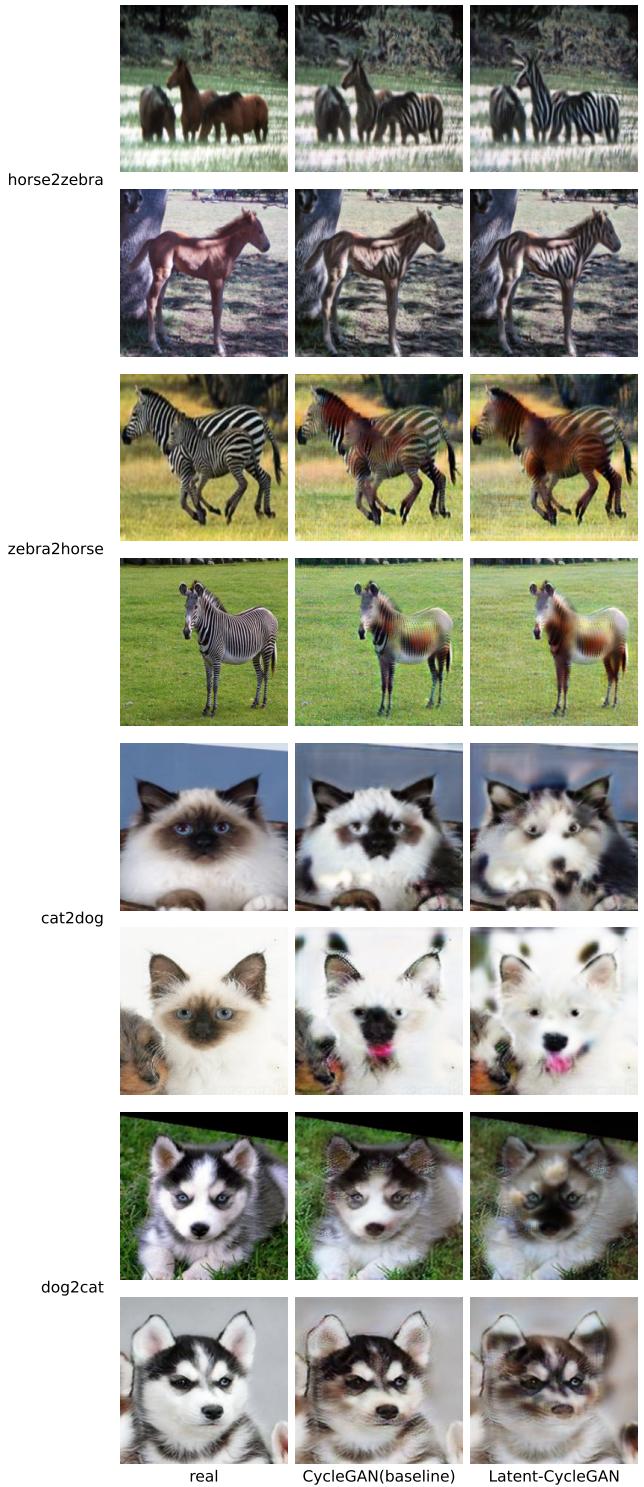
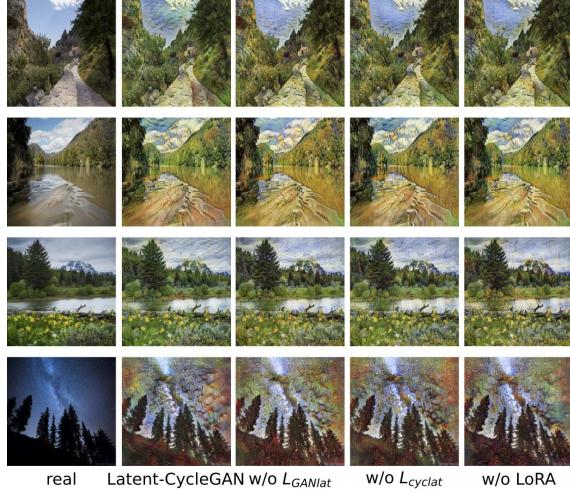


Figure 5: Visualization results of image translation on datasets horse2zebra and cat2dog based on the Latent-CycleGAN framework.



(a) Example of translation results from real photos to Van Gogh style.



(b) Example of translation results from Van Gogh's Paintings to real style.

Figure 6: Visualization results of image translation based on the Latent-CycleGAN framework.

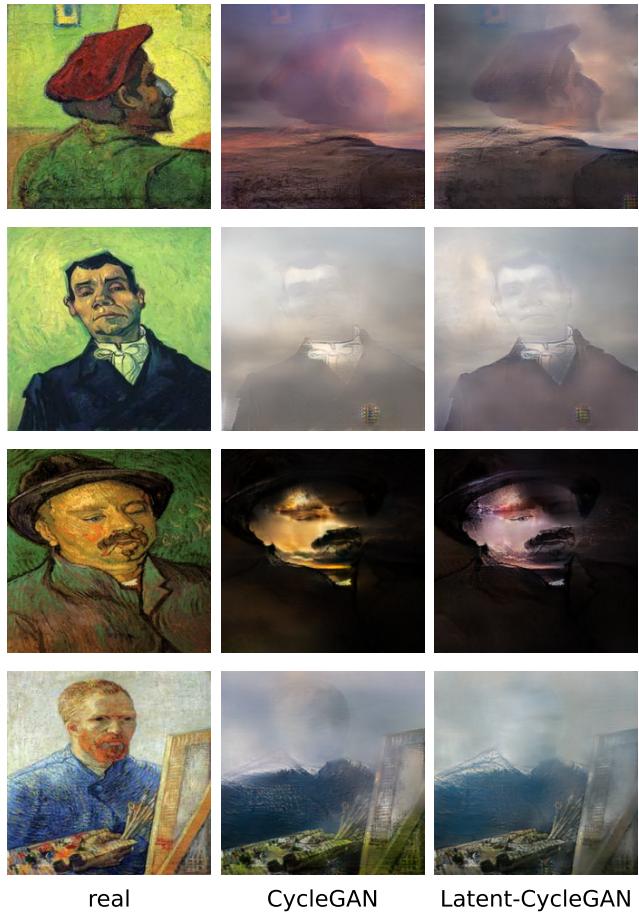


Figure 7: An illustrative failure case where a Van Gogh style portrait is translated into a domain primarily composed of landscape scenes. As a result, the person's facial features and body structure become “blended” into background-like textures, demonstrating how the model may unintentionally force human portraits to conform to a scenic context.