

Comparative Analysis of Image Style Transfer Using Diffusion and GAN Architectures

Chengke Zou, Yankai Mao, Zhongyuan Cao

1 Introduction

Generative AI is undeniably one of the most important and fastest-moving fields in computer science. Within this rapidly evolving landscape, neural style transfer represents a particularly well-defined and instructive problem. It combines the creativity of visual synthesis with the rigor of algorithmic design, offering a contained yet meaningful challenge for empirical exploration.

This project focuses on neural style transfer as a means to examine and compare key generative architectures. Specifically, it investigates two representative paradigms that have shaped the field in recent years: a Generative Adversarial Network (GAN) model and a modern Diffusion model. Both will be implemented and trained on the ArtBench dataset, which contains a diverse set of high-quality oil paintings. The resulting models will be evaluated through established quantitative metrics to assess their relative performance and to highlight their respective advantages and limitations in practical style transfer tasks.

2 Related Work

Image style transfer has long been a central topic in example-guided artistic image generation. Early approaches relied on handcrafted low-level features to match local patches between content and style images. More recently, pre-trained deep convolutional neural networks have been employed to capture feature distributions, enabling a more effective representation of complex style patterns.

2.1 Style Transfer with Diffusion model

Diffusion models are inspired by non-equilibrium thermodynamics (7). They define a Markov chain of diffusion steps that slowly add random noise to the data, and then learn to reverse the diffusion process to construct desired data samples from the noise. Building upon this principle, Denoising Diffusion Probabilistic Models (DDPM) (10) provided a tractable and scalable framework for training deep generative models using a simple Gaussian noise schedule and a reweighted variational objective, significantly improving image fidelity over prior likelihood-based methods.

Subsequent developments refined both training and sampling efficiency, such as Improved DDPM and Score-based Generative Models (SGMs). These advances enabled controllable and high-quality image generation, leading to large-scale diffusion systems such as GLIDE, Imagen, and Stable Diffusion, which further incorporate text conditioning via CLIP or transformer-based encoders.

For visual creativity, UnCLIP inverted the CLIP representation to guide diffusion sampling, allowing semantic control over generation consistent with textual and visual embeddings.

More recently, Inversion-Based Style Transfer with Diffusion Models (11) formalized the idea of reconstructing and editing specific diffusion trajectories to transfer styles between arbitrary images. By inverting the diffusion process to the latent noise space and re-synthesizing with altered conditions, these methods achieve fine-grained and faithful style transfer while preserving content structure.

2.2 Style Transfer with GANs

Research on Image-to-Image (I2I) translation, a task strongly shaped by Generative Adversarial Networks (GANs), started with supervised approaches. Pix2Pix (2) established an early baseline for paired translation, where each input image has a corresponding target.

The challenge of handling unpaired data, which is much more common in real-world settings, was later addressed by CycleGAN (3). Its introduction of the cycle-consistency loss made it possible to train models without aligned image pairs, as demonstrated in tasks such as photo-to-painting translation.

46 Follow-up research sought to improve both scalability and controllability. StarGAN v2 (4) proposed
47 a multi-domain framework with an instance-level style encoder, enabling diverse image synthesis
48 from reference examples.

49 More recently, studies have moved toward text-guided control, made possible by large-scale vision-
50 language models. StyleGAN-NADA (5), for instance, showed that a pre-trained StyleGAN can be
51 adapted to new domains (e.g., “a photo → a sketch”) using only text prompts from CLIP, removing
52 the need for explicit style references.

53 **3 Method and Algorithm**

54 **3.1 CycleGAN**

55 In first method, We adopt the classic GAN-based framework, CycleGAN(3), to address the style
56 transfer problem between real photos and artworks. Unlike traditional style transfer methods requiring
57 paired datasets, CycleGAN performs unpaired image-to-image translation by jointly learning two
58 mappings $G : \mathcal{X} \rightarrow \mathcal{Y}$ and $F : \mathcal{Y} \rightarrow \mathcal{X}$, along with discriminators $D_{\mathcal{X}}$ and $D_{\mathcal{Y}}$ that enforce realism in
59 each domain. The adversarial objective encourages generated samples to be indistinguishable from
60 real ones:

$$\mathcal{L}_{GAN}(G, D_{\mathcal{Y}}) = \mathbb{E}_{y \sim p_{\mathcal{Y}}(y)} [\log D_{\mathcal{Y}}(y)] + \mathbb{E}_{x \sim p_{\mathcal{X}}(x)} [\log(1 - D_{\mathcal{Y}}(G(x)))].$$

61 To preserve content, a cycle-consistency loss is introduced:

$$\mathcal{L}_{cyc}(G, F) = \mathbb{E}_x [\|F(G(x)) - x\|_1] + \mathbb{E}_y [\|G(F(y)) - y\|_1].$$

62 The full objective is $\mathcal{L} = \mathcal{L}_{GAN}(G, D_{\mathcal{Y}}) + \mathcal{L}_{GAN}(F, D_{\mathcal{X}}) + \lambda \mathcal{L}_{cyc}(G, F)$, where λ balances realism
63 and reconstruction.

64 In our project, we plan to use ArtBench(1) as the painting dataset (domain \mathcal{Y}) and ImageNet as the
65 real-world photo dataset (domain \mathcal{X}) to train and evaluate the model.

66 **3.2 Inversion-Based Style Transfer with Diffusion Models**

67 The diffusion model we selected here is the Inversion-based Style Transfer model (11), which is built
68 upon Stable Diffusion models (12). InST learns a textual embedding directly from a single reference
69 image to generate new artistic images in the same style without fine-tuning the diffusion model. Uses
70 the CLIP text encoder to represent a new, learnable concept token. The embedding for this concept,
71 denoted \hat{v} is optimized so that the generated image matches the target artistic image.

$$\hat{v} = \arg \min_v \mathbb{E}_{x, \epsilon, t} [\|\epsilon - \epsilon_{\theta}(z_t, t, v(y))\|_2^2]$$

72 To improve efficiency and generalization, the model introduces multi-layer cross attention: Image
73 embeddings from the CLIP encoder are projected through attention layers to extract key image
74 information. The cross-attention mechanism updates query, key, and value at each layer as:

$$Q_i = W_Q^{(i)} \cdot Q_{i-1}, \quad K_i = W_K^{(i)} \cdot \tau_{\theta}(y), \quad V_i = W_V^{(i)} \cdot \tau_{\theta}(y)$$

75

$$v_{i+1} = \text{softmax} \left(\frac{Q_i K_i^T}{\sqrt{d}} \right) V_i$$

76 Finally, a stochastic inversion step introduces controlled noise to preserve content and enable stylistic
77 diversity, computed as:

$$\hat{\epsilon}_t = (z_{t-1} - \mu_T(z_t, t))\sigma_t$$

78 **3.3 Evaluation**

79 We evaluate results using two complementary metrics. The Structural Similarity Index (SSIM)(8)
80 measures how well the generated image preserves the content and structure of the original, while
81 a CLIP-based similarity score(9) assesses semantic and stylistic alignment between the output and
82 target style description. Higher SSIM and CLIP similarity indicate better content retention and style
83 consistency.

84 **References**

- 85 [1] Peiyuan Liao, Xiuyu Li, Xihui Liu, and Kurt Keutzer. *The ArtBench Dataset: Benchmarking*
86 *Generative Models with Artworks*. In *Proceedings of the European Conference on Computer*
87 *Vision (ECCV) Workshops*, 2022.
- 88 [2] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. *Image-to-Image Translation*
89 *with Conditional Adversarial Networks*. In *Proceedings of the IEEE Conference on Computer*
90 *Vision and Pattern Recognition (CVPR)*, pages 1125–1134, 2017.
- 91 [3] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. *Unpaired Image-to-Image Trans-*
92 *lation using Cycle-Consistent Adversarial Networks*. In *Proceedings of the IEEE International*
93 *Conference on Computer Vision (ICCV)*, pages 2223–2232, 2017.
- 94 [4] Yunjey Choi, Youngik Uh, Jaejun Yoo, and Jung-Woo Ha. *StarGAN v2: Diverse Image Synthesis*
95 *for Multiple Domains*. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
96 *Pattern Recognition (CVPR)*, pages 8187–8197, 2020.
- 97 [5] Rinon Gal, Or Patashnik, Haggai Maron, Amit H. Zada, and Daniel Cohen-Or. *StyleGAN-*
98 *NADA: CLIP-Guided Domain Adaptation of Pretrained StyleGANs*. In *ACM SIGGRAPH 2022*
99 *Conference Proceedings*, pages 1–10, 2022.
- 100 [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. *ImageNet: A Large-Scale*
101 *Hierarchical Image Database*. In *Proceedings of the IEEE Conference on Computer Vision and*
102 *Pattern Recognition (CVPR)*, pages 248–255, 2009.
- 103 [7] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. *Deep*
104 *Unsupervised Learning Using Nonequilibrium Thermodynamics*. In *Proceedings of the 32nd*
105 *International Conference on Machine Learning (ICML)*, 2015.
- 106 [8] Jie An, Siyu Huang, Yibing Song, Dejing Dou, Wei Liu, and Jiebo Luo. *ArtFlow: Unbiased*
107 *Image Style Transfer via Reversible Neural Flows*. In *Proceedings of the IEEE/CVF Conference*
108 *on Computer Vision and Pattern Recognition (CVPR)*, pages 862–871, 2021.
- 109 [9] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
110 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
111 Sutskever. *Learning Transferable Visual Models from Natural Language Supervision*. In
112 *Proceedings of the International Conference on Machine Learning (ICML)*, pages 8748–8763.
113 PMLR, 2021.
- 114 [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. *Denoising Diffusion Probabilistic Models*. In
115 *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 116 [11] Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and
117 Changsheng Xu. *Inversion-Based Style Transfer With Diffusion Models*. In *Proceedings of the*
118 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 119 [12] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. *High-*
120 *Resolution Image Synthesis with Latent Diffusion Models*. In *Proceedings of the IEEE/CVF*
121 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.