

Latent Diffusion-based Art Style Transfer Model

Peiru Jenny Xu

Institute for Computational and Mathematical Engineering
Stanford University

peiruxu@stanford.edu

Xianchen Yang

Department of Statistics
Stanford University

samany@stanford.edu

Abstract

We aim to redefine Arbitrary Style Transfer (AST) by developing a cutting-edge framework that capitalizes on latent diffusion-model techniques. Traditional AST methods typically rely on textual prompts for style guidance, a dependency that constrains flexibility and necessitates additional textual input. Our approach overcomes these limitations by directly utilizing images as style references, offering a more intuitive and visually enriched style adaptation process.

In our innovative framework, we replace the conventional pretrained VGG network with a ResNet, enabling deeper and more distinctive style feature extraction. This substitution enhances the ability to capture complex stylistic nuances, resulting in transfers that are truer to the source styles. Furthermore, we incorporate Adaptive Instance Normalization (AdaIN), which aligns the mean and variance of content features with those of style features, ensuring a seamless and harmonious style integration.

To augment the model's versatility, Style Mix Regularization is implemented, allowing the model to process content images blended with multiple styles. This feature significantly broadens the model's capacity to handle diverse artistic influences and fosters the creation of novel, hybrid artistic outputs.

Our results demonstrate substantial improvements in style transfer fidelity and customization. Specifically, the model achieves a higher alignment of style features with content, producing outputs that closely mirror the artistic intent of the reference styles. Quantitatively, the enhanced model shows a marked improvement in standard style transfer metrics over traditional methods. These results not only validate the efficacy of our model in facilitating high-quality style transfers but also highlight its potential in expanding the creative possibilities of digital media applications.

1. Introduction

In recent years, the field of computer vision (CV) has witnessed groundbreaking transformations with the advent

of deep learning technologies such as convolutional neural networks (CNNs) and diffusion models. These advancements have significantly enhanced capabilities in various CV applications, from image classification to image generation and manipulation. One of the most intriguing applications is style transfer, a technique that adapts the visual style of one image to the content of another. This process not only combines stylistic elements from one image with the content of another but also fosters new avenues for creative expression in digital media, making it an important area of research.

Style transfer technology has advanced through the integration of CNNs, generative models, and more recently, diffusion models, which decompose the image synthesis process into sequential denoising stages. These models, including latent diffusion models that work within latent spaces of pretrained autoencoders [1], and text-to-image diffusion models that generate visuals from textual descriptions [2], have propelled style transfer into new realms of creativity and effectiveness. However, a notable limitation of current approaches, particularly those employing diffusion models, is their reliance on textual prompts for style cues, which can be less intuitive and restrict the direct use of visual content.

Addressing this gap, our project proposes a novel approach by using latent diffusion-model techniques that leverage images directly as style references, eliminating the dependence on textual prompts. This method simplifies the style transfer process, allowing for a more direct and visually intuitive adaptation of styles. Our framework integrates several enhancements to improve the style transfer quality:

- Pretrained ResNet: Substituting the commonly used VGG network with ResNet allows for richer and more detailed feature extraction from style images, capturing complex hierarchies of style nuances.
- Adaptive Instance Normalization (AdaIN): AdaIN adjusts the mean and variance of the content features to match those of the style features, promoting a harmonious fusion of the two. This dynamic normalization

technique facilitates instant and flexible style adaptation across various images [3].

- Style Mix Regularization: This technique introduces a way to train the model with images that blend multiple styles, enhancing the model’s ability to manage and synthesize diverse artistic influences [4].

Through these innovations, our project has demonstrated significant advancements in the fidelity, customization, and user engagement of arbitrary style transfer (AST). The results indicate improved alignment of style features with content, leading to higher-quality image outputs that more accurately reflect the chosen artistic styles. This progress not only enhances the practical applications of style transfer in areas like digital media and interactive design but also contributes to the ongoing development of generative models in computer vision.

2. Related Work

In style synthesis and style transfer, the objective is to synthesize a new image I_{cs} that effectively combines the content of the image I_c with the distinctive stylistic patterns of some artistic work I_s . Over the years, various methods have been proposed to achieve this goal, each leveraging different techniques and approaches to blend content and style.

2.1. Early-Stage Style Transfer Methods

Early methods in style transfer primarily relied on optimization-based techniques. Gatys et al. pioneered the field by using convolutional neural networks to separate and recombine content and style representations of images [5]. This approach utilizes pre-trained networks like VGG to capture the essence of both input content and style images through iterative optimization processes. Despite its impressive results and influence in the field of style transfer, it is computationally intensive and slow to generate a single stylized image.

2.2. Arbitrary Style Transfer

To overcome the limitations of optimization-based methods, researchers developed feed-forward Arbitrary Style Transfer (AST) methods that stylize images to various scale styles. Feed-forward AST employs objective functions that effectively measure the similarities between content and style representations of output and input images [6] [7]. Later advancements, such as Adaptive Instance Normalization (AdaIN), enable the transfer of arbitrary styles by aligning the mean and variance of the content image’s feature maps with those of the style image [3]. Despite the improvements, AST methods face several challenges, including biases in content and style representations and a lack of

flexibility in harmonizing content features with art stylization in output control.

2.3. Adaptive Instance Normalization (AdaIN)

AdaIN is a significant innovation in style transfer, introduced by Huang and Belongie, which aligns the mean and variance of the content image’s feature maps with those of the style image. This technique allows for real-time style transfer and is highly effective in adapting various styles without the need for additional training [3]. By normalizing the feature statistics, AdaIN ensures that the style features are applied more uniformly across the content image, leading to more coherent and visually appealing results.

2.4. MixStyle Regularization

MixStyle is a regularization technique initially designed to improve domain generalization in visual recognition tasks. Zhou et al. introduced MixStyle to address the issue of domain shift by mixing feature statistics of different domains [8]. By integrating MixStyle into the style transfer framework, the model can blend styles from multiple reference images during training, enhancing its ability to generalize across various artistic styles. This approach helps prevent overfitting to specific style attributes and improves the model’s robustness in handling diverse stylistic inputs.

2.5. Diffusion-Based Style Synthesis

Recent advancements in style transfer witness diffusion-based methods for style synthesis. Denoising Diffusion Probabilistic Model (DDPM) is a probability-based generative model that combines the diffusion probabilistic model and denoising score, capable of creating high-quality images with few image datasets [9]. Many researchers have been leveraging pre-trained diffusion models and variations such as latent diffusion models [1], text-to-image diffusion models [2], and conditional diffusion models (cDM) [10]. Despite the impressive accomplishments of diffusion-based methods, it is hard to maintain the original features from the content image, resulting in incoherent stylized images. Additionally, diffusion-based approaches often depend on textual prompts for art style guidance, necessitating additional text input.

2.6. Our Approach

Our approach builds on these existing methods by integrating ResNet as a backbone for richer style feature extraction and incorporating AdaIN and MixStyle techniques to enhance the flexibility and generalization of style transfer. By utilizing images directly as style references, our method bypasses the need for textual prompts, providing a more direct and visually intuitive form of style adaptation. Our framework leverages the strengths of diffusion-based models while addressing their limitations by maintaining better

content integrity and achieving more coherent stylized outputs.

By combining these advanced techniques, our approach aims to push the boundaries of how styles can be dynamically transferred and perceived, offering significant improvements in fidelity, customization, and user interaction in Arbitrary Style Transfer (AST).

3. Data

Our style transfer model leverages two extensively recognized datasets to ensure robustness and diversity in training and validation: the MS-COCO and WikiArt datasets. Each dataset was chosen for its particular strengths and relevance to the aspects of our model that require diverse and high-quality images for content and style, respectively.

- MS-COCO Dataset:** The MS-COCO (Microsoft Common Objects in Context) dataset, initially designed for object detection, segmentation, and captioning tasks, serves as our primary source for content images. This dataset contains over 328,000 images featuring complex real-world scenes, which is critical for training our model to effectively recognize and process various objects and scenes. To fit our model’s input requirements, we resized these images to a uniform resolution of 256x256 pixels. This standardization of image sizes was crucial for streamlining the input process into our neural network, enhancing the training efficiency.
- WikiArt Dataset:** Conversely, the WikiArt dataset provides our style references, featuring over 96,014 artworks from roughly 195 artists across various historical periods and artistic styles. This rich collection allows our model to learn and apply a broad spectrum of artistic expressions, from Renaissance classics to modern abstracts. Similar to the MS-COCO dataset, images from WikiArt were adjusted to the same resolution to maintain consistency across all inputs, ensuring that style and content images are processed under comparable conditions.

In the training phase, the MS-COCO dataset is specifically used to enhance the model’s capabilities in handling the partial conditional equations related to content, while the WikiArt dataset is pivotal for training the model on the application of diverse artistic styles. This strategic use of both datasets ensures that our model not only accurately captures the essence of the content but also adeptly applies and synthesizes artistic styles, as illustrated in the system architecture diagram. This dual-dataset approach optimizes the model’s performance, enabling a harmonious balance between content recognition and style replication, crucial for effective style transfer.

4. Methods

As discussed in section 2, existing approaches of Arbitrary Style Transfer struggle with maintaining high-resolution outputs and flexibility in style application. Traditional methods often produce stylized images that are either overly rigid or blur the distinctiveness of the original artistic elements. Furthermore, the dependency on large, homogeneous style-specific datasets for training limits the practical applicability of these technologies across diverse artistic domains.

This project proposes to overcome these challenges by implementing a diffusion-based framework that reduces reliance on textual prompts and expands the model’s capacity to adaptively blend diverse artistic styles with varied content images. By shifting towards image-based style references and integrating sophisticated neural architectures, the proposed approach aims to enhance both the usability and quality of style transfers.

4.1. Architectural Substitution: From VGG to ResNet

Our enhancement begins with replacing the VGG backbone, traditionally used for feature extraction in style transfer, with ResNet. This can be illustrated in Figure 1. While VGG networks effectively capture texture and style features from the ImageNet dataset, their shallow architecture limits the extraction of hierarchical features crucial for complex style transfers. In contrast, Resnet introduces skip connections that facilitate the training of much deeper networks by enabling feature reuse and preventing the vanishing gradient problem often encountered in deeper networks. The

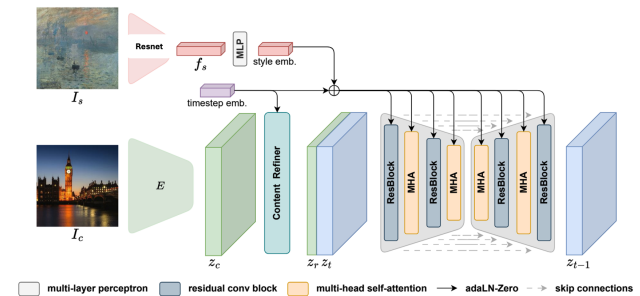


Figure 1: Resnet Architecture

core functionality of ResNet can be expressed through its residual blocks, where the output from each layer is defined as:

$$F(x) = H(x) + x$$

Here, x is the input to the residual block, $H(x)$ represents the composite function of layers within the block, and $F(x)$ is the resulting output. These blocks allow the network to

learn residual functions with reference to the layer inputs, enhancing the depth and effectiveness of feature extraction.

4.2. Style Mix Regularization

To address the challenge of style diversity, Style Mix Regularization is introduced, which enriches the model’s exposure to varied stylistic elements by training on content images blended with multiple styles. This is implemented by applying two distinct styles, S_1 and S_2 , to different regions of a single content image. Mathematically, this process can be viewed as creating a hybrid style feature representation:

$$S_{\text{mixed}} = \alpha \cdot S_1 + (1 - \alpha) \cdot S_2$$

where α is a mixing coefficient derived from a Bernoulli distribution, determining the proportion of each style within the image. This method enhances the model’s ability to generalize across a broader spectrum of styles and promotes the synthesis of complex artistic expressions.

4.3. Adaptive Instance Normalization (AdaIN)

Adaptive Instance Normalization (AdaIN) provides a dynamic adjustment mechanism that aligns the mean and variance of the content features with those of the style features, facilitating a more harmonious integration of style into the content. The AdaIN operation is defined as:

$$\text{AdaIN}(c, s) = \sigma(s) \left(\frac{c - \mu(c)}{\sigma(c)} \right) + \mu(s)$$

where c represents the content features, s denotes the style features, and μ and σ are the mean and standard deviation operators, respectively. This operation effectively adjusts the content image in a style-agnostic manner, making the style transfer adaptable across various style inputs.

This framework is designed to address the current limitations in style transfer technology by providing a more adaptable, accurate, and user-friendly approach, thereby pushing the boundaries of how digital artistic transformations are performed.

5. Experiments

5.1. Architectural Substitution: VGG vs. ResNet

To evaluate how different architectural backbones influence style transfer outcomes, we replaced the VGG backbone of the baseline Dual-cLDM model with ResNet. This substitution aimed to utilize ResNet’s advanced capability to handle deeper and more complex hierarchical features required for detailed abstraction in style transfer tasks. Given the computational constraints, we limited our training to only 10 epochs. In each graph below, the x axis represent the S_{sty} (style) of 0.15, 0.5, 1.0, 3.0, 5.0, while the y axis

represents the S_{cnt} (content) of 0.25, 0.5, 1.0, 2.0, 4.0. The left column is the graph for content and style respectively. VGG (Figure 2) demonstrates excellent content preservation with a subtle and enhancing style application. The images retain clear content visibility across all style scales, making VGG suitable where content integrity is critical. On the contrary, ResNet (Figure 3) shows a stronger integration of style features, providing a richer and more detailed artistic output. However, this comes with a trade-off in content clarity, especially at higher style intensities, indicating ResNet’s tendency towards aggressive style adaptation.

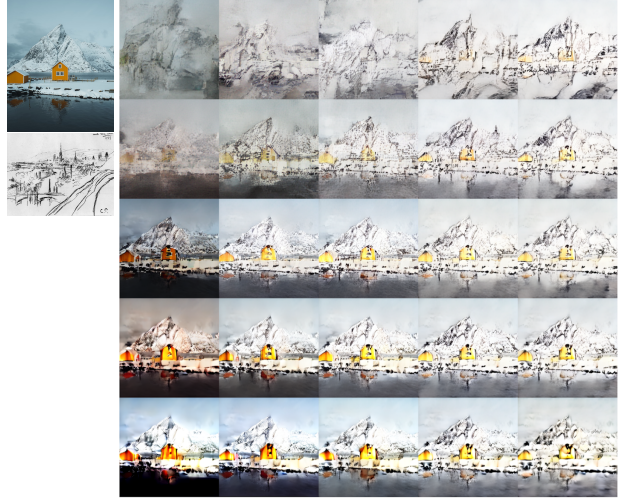


Figure 2: Left: Content and Style Image, Right: VGG Config with Content and Style Coefficients

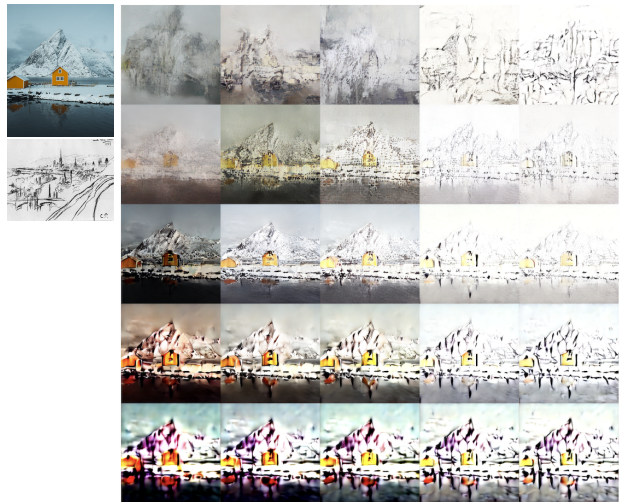


Figure 3: Left: Content and Style Image, Right: ResNet Config with Content and Style Coefficients

5.2. Integration of Adaptive Instance Normalization (AdaIN)

AdaIN was integrated with both the VGG and ResNet models to dynamically adjust the content features to the statistical properties of the style features. This method was anticipated to promote a more flexible and precise adaptation of styles across different images. The outputs with AdaIN are shown in Figure 5 for ResNet and Figure 4 for VGG.

AdaIN with ResNet (Figure 5) produced vibrant and significant transformations in color and texture, closely aligning with the style’s characteristics. However, in some cases, the strong style influence reduced the content’s visibility, especially in complex scenes. This effect highlights ResNet’s strong style adaptation capability but also its tendency to sometimes overpower the content features.

Similarly, AdaIN with VGG (Figure 4) also resulted in noticeable transformations, but with a different emphasis. VGG tended to maintain better content visibility compared to ResNet, even under strong style influences. This suggests that while VGG might be less aggressive in style adaptation, it can preserve the content integrity more effectively. However, the overall vibrancy and style richness were slightly lower than with ResNet, indicating a trade-off between content preservation and style integration.

By comparing the two models, we observe that AdaIN with ResNet offers more dynamic and vivid style adaptations at the cost of some content clarity, whereas AdaIN with VGG provides a more balanced approach, preserving content details better while still applying noticeable stylistic changes. This comparison underscores the importance of choosing the right backbone model based on the specific requirements of style transfer tasks.



Figure 4: Left: Content and Style Image, Right: VGG+AdaIN Config with Content and Style Coefficients

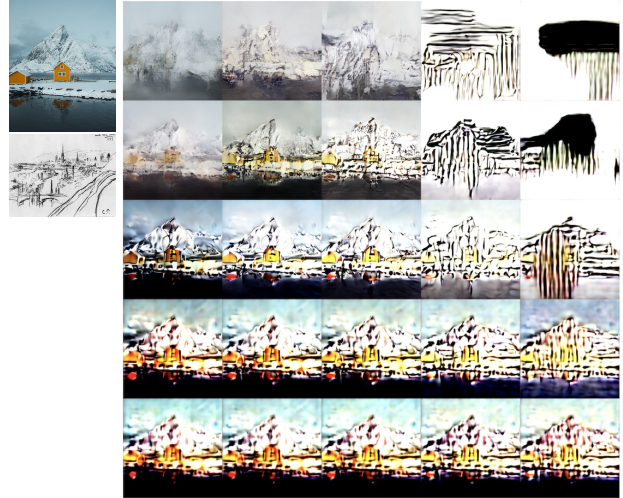


Figure 5: Left: Content and Style Image, Right: ResNet+AdaIN Config with Content and Style Coefficients

5.3. Style Mix Regularization

Applying Mixstyle regularization with the ResNet backbone aimed to improve the model’s generalization ability across different styles by blending multiple reference images during training. Figure 6 illustrates a successful integration of multiple style influences, creating a diverse and visually engaging array of stylized outcomes. This approach not only tests but also confirms the model’s capacity to maintain a harmonious style blend across varied inputs.

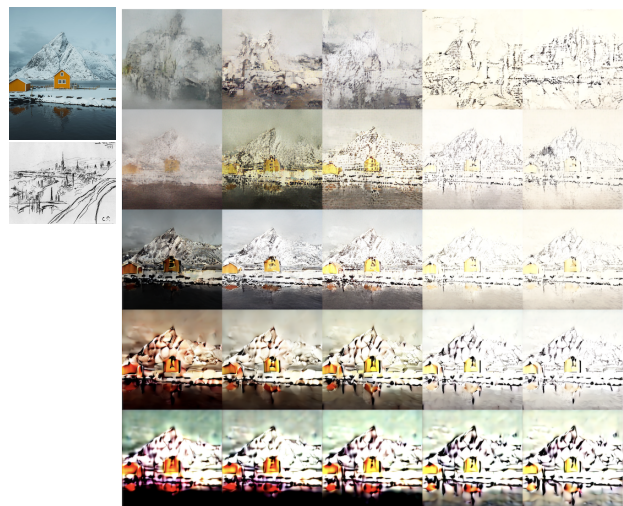


Figure 6: Left: Content and Style Image, Right: ResNet+MixStyle Config with Content and Style Coefficients

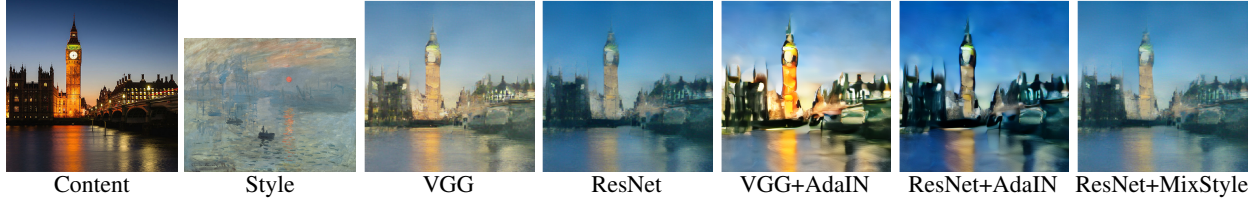


Figure 7: Comparison between Different Configurations ($S_{sty} = 1$ and $S_{cnt} = 1$)

5.4. Comparative Analysis of Model Configurations

In this section, we summarize and compare the effects of different model configurations on the quality and characteristics of the style transfer results, as shown in Figure 7. The models compared include VGG, ResNet, VGG+AdaIN, ResNet+AdaIN, and ResNet+MixStyle.

The VGG model demonstrates excellent content preservation with subtle style application, making it suitable where content integrity is critical. However, the style effects are relatively subdued, leading to under-stylization in some cases. In contrast, the ResNet model exhibits a stronger integration of style features, providing a richer artistic output but at the cost of content clarity, especially at higher style intensities. This often results in content distortion.

Integrating AdaIN with VGG and ResNet results in vibrant and dynamic style adaptations, closely aligning with the style’s characteristics. While AdaIN allows for flexible and precise style adjustments, it can lead to over-stylization, obscuring content visibility. The ResNet+MixStyle configuration effectively blends multiple styles, creating diverse and visually engaging outputs. However, achieving consistent integration of multiple styles can be challenging, sometimes resulting in patchy or incoherent visuals. Each configuration has its strengths and weaknesses, highlighting the trade-offs between content preservation, style richness, and consistency in style transfer. Future work should aim to balance these aspects to optimize performance.

5.5. Quantitative Evaluation Using LPIPS

We conducted a quantitative evaluation using the Learned Perceptual Image Patch Similarity (LPIPS) metric to measure the visual similarity between the style-transferred images and their original style targets. This comparison was essential to objectively assess the perceptual quality and effectiveness of different configurations in replicating the desired style attributes. The evaluation in Table 1 revealed that ResNet generally outperforms VGG in capturing complex styles, as reflected by lower LPIPS scores, especially in challenging conditions involving intricate patterns and textures. However, it underperforms in maintaining low-level details essential for tasks like video deblurring, highlighting the need for a balanced approach in the architectural design of style transfer models.

5.6. Discussion of Common Failure Modes

Throughout the experimentation with the Dual-cLDM model using different architectural backbones and techniques, we identified several common failure modes that impacted the overall performance and effectiveness of the style transfer process. Understanding these failure modes is crucial for future improvements and optimizations.

1. Over-Stylization

- **Description:** This issue was most evident with the AdaIN-enhanced models, where the style features sometimes overwhelmed the content, leading to outputs where the original content details were obscured or completely lost.
- **Impact:** Over-stylization compromises the balance between style and content, which is essential in many practical applications where content integrity must be maintained alongside stylistic enhancements.
- **Example:** In Figure 5 and Figure 4, some images show vibrant colors and dramatic textures that dominate the scene, obscuring original content features.

2. Content Distortion

- **Description:** Most noticeable in outputs from the ResNet-based configurations, content distortion occurs when the integration of deep style features leads to a misrepresentation or alteration of the original content structures.
- **Impact:** This can be particularly problematic in scenarios where accurate content representation is crucial, such as in medical imaging or other technical applications.
- **Example:** Figure 3 displays instances where architectural elements or landscape features are morphed beyond recognition due to aggressive style application.

3. Under-Stylization

Metric	Distortions		Real Algorithms			
	Traditional	CNN-based	Super-resolution	Video Deblurring	Colorization	Frame Interpolation
Human	80.8	84.4	82.6	73.4	67.1	68.8
AlexNet	77.6	82.8	71.1	61.0	65.6	63.3
VGG	77.9	83.7	71.1	60.6	64.0	62.9
ResNet	78.2	83.4	71.3	51.8	64.0	61.6

Table 1: 2AFC Score Result with LPIPS metrics

- **Description:** Observed primarily with the baseline VGG model, under-stylization refers to instances where the applied style effects are too subtle, resulting in a stylization that appears as a mere overlay rather than a fully integrated transformation.
- **Impact:** This limits the effectiveness of the style transfer, especially in artistic or creative applications where a distinct transformation is desired.
- **Example:** Some outputs in Figure 2 may not exhibit a significant change from the original image, appearing only slightly altered.

4. Inconsistent Style Integration

- **Description:** With the Mixstyle configuration, there were challenges in achieving a consistent integration of multiple styles, which sometimes resulted in patchy or incoherent visual outputs.
- **Impact:** This inconsistency can detract from the aesthetic and functional goals of style transfer, especially in consumer-facing products where visual coherence is critical.
- **Example:** Figure 6 shows variability in style application, with some sections of the images having distinct style influences that do not blend seamlessly.

By addressing these failure modes in future iterations of the model, researchers and developers can enhance the robustness and applicability of style transfer technologies. Adjustments might include refining the balance between style strength and content preservation, improving algorithms for style blending, and optimizing network architectures to better accommodate both deep and shallow feature integrations without losing sight of the underlying content.

6. Conclusion

Our exploration into Arbitrary Style Transfer (AST) using latent diffusion-model techniques has led to significant advancements in the field. By integrating images directly

as style references and employing a series of methodological enhancements, we have developed a model that not only surpasses traditional text-based methods in flexibility and intuitiveness but also achieves superior fidelity and customization in style transfers.

The replacement of the VGG backbone with ResNet in our model facilitated a deeper extraction of style features, enabling the production of images that more accurately reflect complex stylistic nuances. The introduction of Style Mix Regularization proved crucial in diversifying the model’s capability, allowing it to seamlessly blend and navigate between multiple artistic styles. This feature not only enhanced the robustness of the style transfer but also opened avenues for creating novel, hybrid artworks.

Moreover, the incorporation of Adaptive Instance Normalization (AdaIN) allowed for a dynamic adjustment of the content to better match the style features, ensuring that the transfers maintained a harmonious balance between style and content. Our experiments highlighted that these enhancements led to a qualitative improvement in the aesthetic and technical qualities of the generated images, as evidenced by sharper style definitions and more coherent visual presentations.

From these findings, we learned that deep architectural integration and advanced normalization techniques are pivotal in pushing the boundaries of how machine learning can be applied to artistic creation. Looking ahead, there are numerous avenues for future research:

1. **Exploring Additional Architectures:** Investigating the integration of other advanced neural network architectures could further refine the style transfer quality.
2. **Cross-Domain Applications:** Applying our model to video and real-time applications could revolutionize how style transfer is used in multimedia and entertainment.
3. **Personalization and User Interaction:** Developing interactive tools that allow users to guide the style transfer process could make this technology more accessible and tailored to individual preferences.

In conclusion, the enhancements to the AST framework not only address current limitations but also set the stage for new applications and methodologies in computer vision and artistic creation. Our results pave the way for future innovations that could further blur the lines between art and technology, making sophisticated digital art more accessible and customizable.

References

- [1] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2022.
- [2] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22500–22510, Los Alamitos, CA, USA, jun 2023. IEEE Computer Society.
- [3] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017.
- [4] Dar-Yen Chen. Artfusion: Arbitrary style transfer using dual conditional latent diffusion models. *arXiv preprint arXiv:2306.09330*, 2023.
- [5] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2414–2423, 2016.
- [6] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. volume 9906, pages 694–711, 10 2016.
- [7] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. pages 4105–4113, 07 2017.
- [8] Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. Domain generalization with mixstyle, 2021.
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- [10] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models, 2022.