

AnimeDiffusion: Anime Face Line Drawing Colorization via Diffusion Models

Yu Cao[†], *Student Member, IEEE*, Xiangqiao Meng[†], P.Y. Mok, *Member, IEEE*,
Xueting Liu, Tong-Yee Lee, *Senior Member, IEEE* and Ping Li, *Member, IEEE*

Abstract—It is a time-consuming and tedious work for manually colorizing anime line drawing images, which is an essential stage in cartoon animation creation pipeline. Reference-based line drawing colorization is a challenging task that relies on the precise cross-domain long-range dependency modelling between the line drawing and reference image. Existing learning methods still utilize generative adversarial networks (GANs) as one key module of their model architecture. In this paper, we propose a novel method called AnimeDiffusion using diffusion models that performs anime face line drawing colorization automatically. To the best of our knowledge, this is the first diffusion model tailored for anime content creation. In order to solve the huge training consumption problem of diffusion models, we design a hybrid training strategy, first pre-training a diffusion model with classifier-free guidance and then fine-tuning it with image reconstruction guidance. We find that with a few iterations of fine-tuning, the model shows wonderful colorization performance, as illustrated in Fig. 1. For training AnimeDiffusion, we conduct an anime face line drawing colorization benchmark dataset, which contains 31696 training data and 579 testing data. We hope this dataset can fill the gap of no available high resolution anime face dataset for colorization method evaluation. Through multiple quantitative metrics evaluated on our dataset and a user study, we demonstrate AnimeDiffusion outperforms state-of-the-art GANs-based models for anime face line drawing colorization. We also collaborate with professional artists to test and apply our AnimeDiffusion for their creation work. We release our code on <https://github.com/xq-meng/AnimeDiffusion>.

Index Terms—Line drawing colorization, diffusion models, conditional generation



1 INTRODUCTION

LINE drawing colorization is an essential process in the animation industry, however, manually colorizing is time consuming, especially for the line drawings with complex structure content. So, it is necessary and valuable to design a kind of automatic line drawing colorization system. Line drawing colorization is challenging, because line drawings, different from grayscale images [1], [2], [3], [4], only contain structure content composing of a series of lines without any luminance or texture information. This question has greatly attracted attention of researchers in the field of Computer Graphics, therefore many approaches [5], [6], [7], [8] are proposed for manga and cartoon line drawing colorization during the past time.

Early work [9] utilized neural network to automatically colorize the cartoon images with random color, which is the first deep learning-based cartoon colorization method. Nevertheless, many interactions are usually needed to refine the colored results to satisfy what the user specified. In order to effectively control the color of the colored result, many user-hint based methods have been proposed successively, such as scribble colors [10], point colors [11], text-hint [12], and

language-based [13]. While these user-hint based methods are still not convenient and intuitive, especially for amateur users without aesthetic judgement. Reference-based colorization methods, such as [14], [15], [16], [17], [18], [19] provide a more convenient way. Users only need to prepare a line drawing and a corresponding reference color image, and the algorithm can automatically complete the colorizing process without other manual intervention.

Reference-base line drawing colorization can be formulated as a conditional image generation task. Since generative adversarial networks (GANs) has become the mainstream model for many generation tasks in the last decade, especially using images as generation conditions. Many previous work for line drawing colorization utilized GANs as one of the most important module of their model architecture design. However, this kind of approach mainly focuses on the improvement of feature aggregation module of two extracted deep features. For example, Lee et al. [14] proposed an attention based Spatial Correspondence Feature Transfer (SCFT) module. Li et al. [15] eliminated the gradient conflict among attention branches by using Stop-Gradient Attention (SGA) module. Cao et al. [19] designed an attention-aware model for generating high quality colored anime line drawing images. Otherwise, GANs-based models require to deployment of multiple losses, which also increases the instability of the training process.

Along with the diffusion probabilistic models [20] (“diffusion models” for brevity) has been proved to be an excellent model that is capable of generating high quality images, many algorithms based on diffusion models have been proposed in recent years. As a novel generation algo-

- [†] indicates equal contribution.
- Y. Cao, X. Meng, P.Y. Mok and P. Li are with The Hong Kong Polytechnic University, Hong Kong SAR, China. E-mail: {yu-daniel.cao, xiangqiao.meng}@connect.polyu.hk, {tracy.mok, p.li}@polyu.edu.hk.
- X. Liu is with Caritas Institute of Higher Education, Hong Kong SAR, China. E-mail: tliu@cihe.edu.hk.
- T.-Y. Lee is with National Cheng Kung University, Tainan, Taiwan. E-mail: tonylee@mail.ncku.edu.tw.

Manuscript received April 19, 2005; revised August 26, 2015.

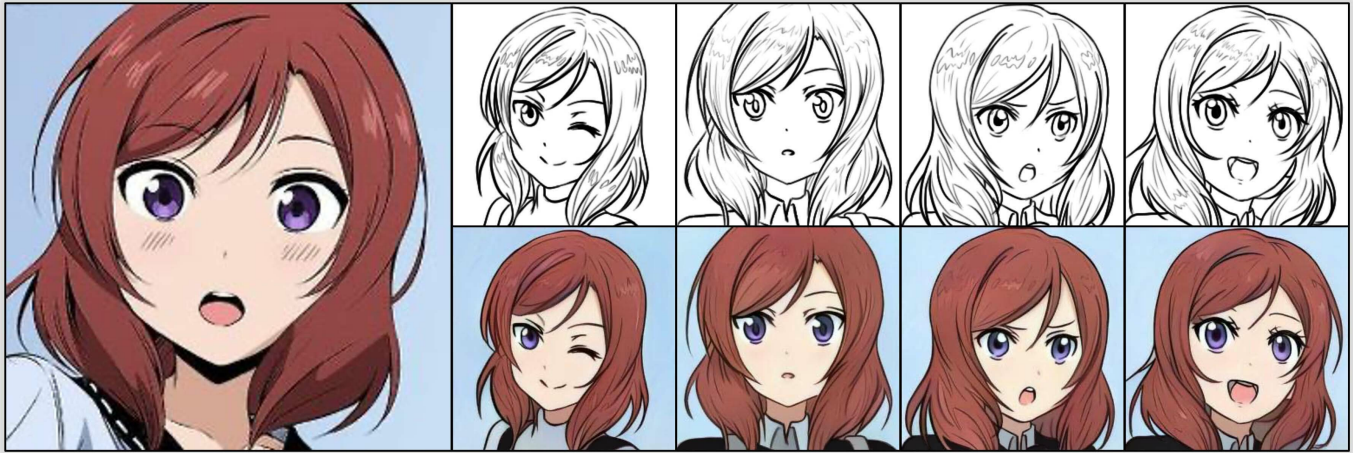


Fig. 1. We propose AnimeDiffusion that performs reference-based line drawing colorization. Give one reference color image (the left side) and four line drawings of the same character (on the top), AnimeDiffusion generates four colored results (on the bottom) with accurate color and semantic correspondence. In particular, we can generate precise eye color and surprising hair details. The anime character is Nishikino Maki of LoveLive and line drawings are drawn by Ms. Xiao Meng.

rithm, it has greatly promoted the progress of AI-Generated Content (AIGC) technology. Inspired by this, we propose the first diffusion model called AnimeDiffusion tailored for anime face line drawing colorization. Since the diffusion models usually have the problem of high computing consumption, we design a hybrid training strategy that consists of classifier-free guidance pre-training stage and image reconstruction guidance fine-tuning stage. During the AnimeDiffusion training procedure, we train a U-Net which regards the line drawings and reference images as conditional denoising input. In order to make the model learn semantic correspondence ability, the reference image is a geometry distorted version of original reference image by applying Thin-Plate Splines (TPS) transformation. The original reference image is added to a Gaussian noise and concatenated with line drawing and reference image together. The U-Net uses the concatenated images as input and predicts the noise that added onto the original reference feature map. This pre-training process mainly makes AnimeDiffusion learn the denoising ability. In the fine-tuning stage of AnimeDiffusion, we calculate the MSE loss between the reconstructed image and the original reference image and update the parameters of AnimeDiffusion when performing reverse sampling task. It is worth to note that our fine-tuning is different from the one applied in other approaches. Existing methods mainly fine-tune the pre-trained image-to-image [21], [22] or text-to-image [23], [24] diffusion models for various kinds of downstream tasks. However, our fine-tuning allows us to train a diffusion model from sketch to perform colorization more efficiently. Experimental results demonstrate that AnimeDiffusion generates better results than the state-of-the-art GANs-based line drawing colorization models both qualitatively and quantitatively. In order to train AnimeDiffusion and fill the gap of no available high resolution anime face dataset, we conduct a novel benchmark dataset for academic research purpose. All original images are chosen from [25].

Our main contributions can be summarized as follows:

- We propose the first AnimeDiffusion model tailored for anime face line drawing colorization. Experiments demonstrate that AnimeDiffusion notably outperforms the GANs-based counterparts and achieves the state-of-the-art anime face line drawing colorization results.
- We design a hybrid training strategy for AnimeDiffusion in order to tackle the problem of high computing consumption of diffusion models. The proposed strategy can accelerate the network convergence and improve colorization performance.
- We conduct a new anime face line drawing colorization benchmark dataset, which contains 31696 training data and 579 testing data. Our dataset aims to fill the gap of no available high resolution (256×256) anime face dataset for training and evaluation.

2 RELATED WORK

2.1 Line Drawing Colorization

Since line drawing contains only structure information with sparse line sets, existing colorization methods for grayscale images cannot directly be used. Many colorization methods tailored for line drawings have been developed. Traditional line drawing colorization approaches [5], [7] are commonly optimized-based which allow users to use brushes to inject desired color into specific regions. With the advancement of deep learning technology and for the better control of color, many user-hint colorization methods spring up. The color hints are usually concatenated with line drawing and encoded as the input for neural network in many deep learning based methods. Ci et al. [10] proposed a conditional GAN model to colorize the anime line drawing using color scribbles, which can generate colored results with accurate shading. Zhang et al. [11] developed a color points hint two-stage colorization method, which divided the complex colorization task into two simpler and goal-clearer subtasks. Kim et al. [12] utilized their SECat module to generate illustrations with quality details using text tags

as their hints. Zou et al. [13] for the first time presented a language-based system for interactive colorization of scene sketches. However, the complexity of such user-hints methods will become more labor-intensive as the number of line drawings increase, many interactions are usually needed to refine the colored results and these methods are not user-friendly for amateur users without aesthetic judgement, especially for preparing appropriate color hints. Therefore, many reference based colorization methods have been proposed, and they are very suitable for colorizing line drawing sets or videos of anime characters, which need to keep the same characters with consistent colors during each frame. Sato et al. [26] segmented the target and reference image into different regions, then represented regions as nodes of a graph structure and colorized the monochrome target image by matching the graphs of the target and reference images. Furusawa et al. [6] proposed the first semi-automatic system to colorize an entire manga with color features extracted from the input reference image. Chen et al. [8] proposed an active learning based framework to match local regions between line arts and reference color image, followed by mixed-integer quadratic programming (MIQP) which considers the spatial contexts to further refine matching results. Shi et al. [27] proposed a new line art video colorization method using 3D convolutional module to refine the temporal consistency of the colored result. Dou et al. [28] is the first work that utilizes the HSV color space for anime sketch colorization. Maejima et al. [29] proposed colorization method for anime character using few-shot learning. Sun et al. [17] trained a dual conditional GAN to colorize contours in different styles which helps designers create icons. Li et al. [30] presented an icon colorization system that is composed of an encoder-decoder network and a conditional normalizing flow. Our AnimeDiffusion is a novel reference-based colorization tailored for anime face line drawing colorization. Compared with previous GANs-based methods, AnimeDiffusion can generate better results both in visual quality and quantitative metrics.

2.2 Semantic Correspondence

Semantic correspondence [31] is one of the fundamental problems in computer vision which goal is to establish dense correspondences across images containing the targets of the same category or with similar semantic information. In computer graphics, it is also very important for exemplar-based image colorization task, and there usually exists the same semantic information between the target image and exemplar image in practical colorization usage. For line art colorization, this can also be viewed as cross-domain correspondence since the texture difference between line drawing and reference color image. Zhang et al. [32] proposed an exemplar-based image translation system based on cross-domain correspondence learning. Lee et al. [14] proposed an attention-based module to spatially match and aggregate the sketch feature and reference color image feature. Li et al. [33] designed a model to colorize grayscale natural image, even if the exemplar image has no similar semantic information of the target grayscale image. He et al. [34] presented a novel progressive color transfer model, which jointly optimizes dense semantic correspondences in the

deep feature domain and the local color transfer in the image domain. Zhang et al. [35] proposed the first end-to-end exemplar-based video colorization algorithm, which unified the semantic correspondence and colorization into a single network. Lu et al. [36] proposed a unified semantic color transfer system from reference image to the target grayscale image. Most of these semantic correspondence learning methods are designed for grayscale image or video colorization. However, we design AnimeDiffusion to perform anime face line drawing colorization, which can generate results with clear and accurate semantic colors.

2.3 Diffusion Models

Diffusion models such as denoising diffusion probabilistic models (DDPM) [37] have achieved great success in image generation tasks. It is shown that image generation models based on diffusion models have better performance in terms of training stability and generation quality [38], [39]. Denoising diffusion implicit models (DDIM) [40] accelerates the sampling procedure and enables a determined generation process with given Gaussian noise. In addition to generating high quality images based on random noise, diffusion models also show good performance in conditional image-to-image translation tasks. An image-to-image diffusion models (Palette) [41] offers a versatile and general framework for image manipulation. Stochastic Differential Editing (SDEdit) [21] is a guided image editing and synthesis method, which synthesizes realistic images by iteratively denoising through a stochastic differential equation (SDE). Combined with the contrastive language-image pre-training (CLIP) [42] model, it can also be used for multimodal generation tasks. A text-guided diffusion models (Diffusion-CLIP) [22] shows the flexibility to add text guidance conditions. Latent diffusion models (LDMs) [23] can be trained on limited computational resources using powerful pre-trained autoencoders in the latent space. Compared with GANs-based models, image generation based on diffusion models can easily add a variety of guidance, such as texts, strokes, and reference images. Since existing diffusion models are designed for natural image generation with random noise or text prompt. Some diffusion models can perform natural image colorization based on the prior knowledge of color in real world. They cannot be directly used for our task. Very recently, Zhang et al. [24] proposed ControlNet which can generate diversity colored cartoon images according to the sketch input and text prompt. Since it used diffusion models pre-trained on natural image dataset, it sometimes produced some distorted results compared to input sketches. Our AnimeDiffusion is the first one to perform reference-based anime face line drawing colorization with accurate color information using diffusion models.

3 OVERVIEW

Given a reference image, we aim to colorize the line drawing with clear geometry structure and accurate semantic color. The core problem to solve is how to inject color from the corresponding position of reference image into line drawing. Previous GANs-based methods usually design two encoders to extract line drawing features and reference

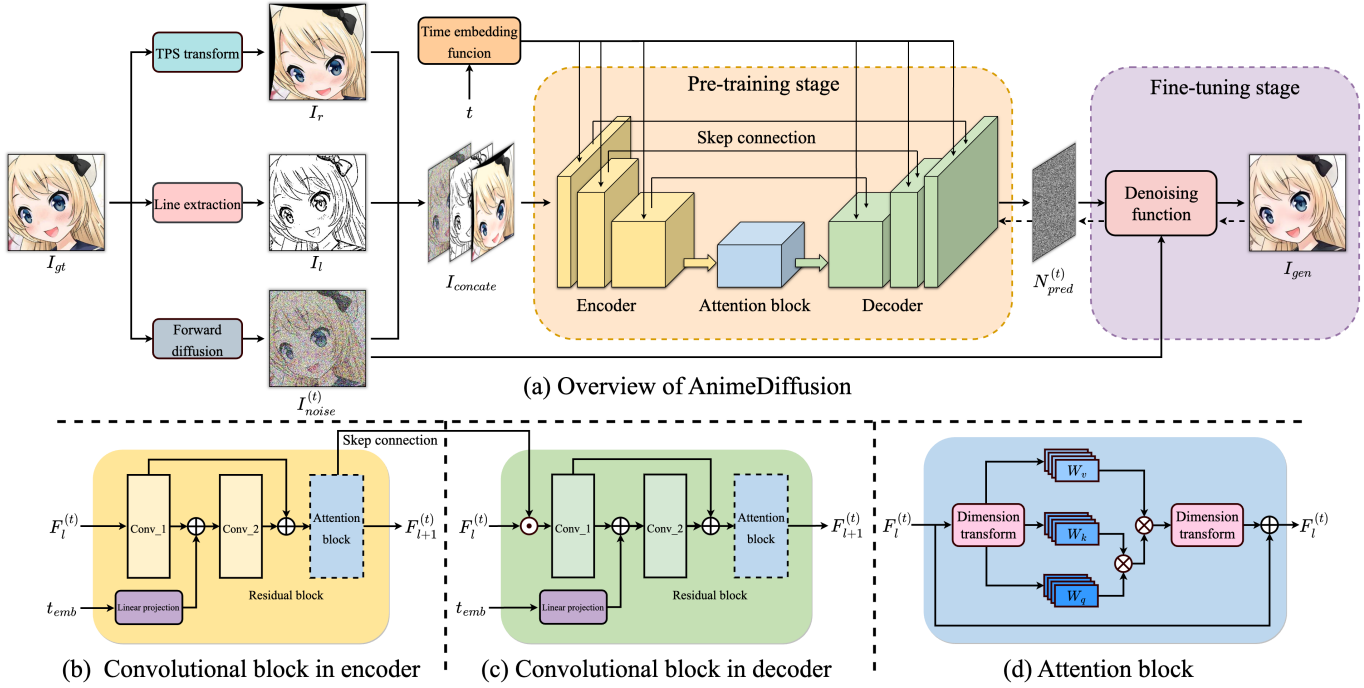


Fig. 2. The flowchart of AnimeDiffusion.

image features respectively, and then use feature aggregation block to integrate two cross-domain features in the latent space. This operation can make model learn semantic correspondence ability to generate colored results, and one discriminator is needed to distinguish the generated results from real colored images which makes the colored results more realistic. However, these GANs-based methods suffer from two problems, one is that line drawing features and reference images features are integrated in the latent space using complicated feature aggregation module. This part of network design is usually stacked with various novel neural network modules and it is not intuitive and explainable for algorithm designers, such as [14], [15], [18]. The other one is that the weighted sum of multiple loss is employed [43], [44] combined with adversarial loss, this makes it difficult to train the network and the training process is more unstable. Additional normalization techniques [45], [46] or improvements of the GAN model itself [47], [48], [49] are needed to solve this problem.

However, diffusion models can fundamentally solve the above two problems. We formulate our training procedure as conditional noise prediction task, therefore, multiple condition images including line drawing and reference are directly concatenated in the pixel space. It is easier to extend diffusion models for other conditional generation tasks without the need to design special feature extractors. In addition, the loss function of diffusion models is simple and closely related to the training task. Combined with our designed hybrid training strategy, our model finally has excellent coloring ability. Without introducing additional discriminator network, the quality of the colored images is extremely close to that of real colored images.

4 ANIMEDIFFUSION

4.1 Model Architecture

As illustrated in Fig. 2(a), assuming I_{gt} is an original colored image, and a line drawing I_l is extracted using XDoG [50] extractor. The more detailed information about our data preparation will be introduced in section 5.1. Since there are usually large spatial structure discrepancy between the line drawing and reference image, in order to make AnimeDiffusion learn the accurate semantic correspondence ability during the training process, we apply TPS transformation [51] to convert I_{gt} to a geometry distorted version I_r . The forward diffusion goes from I_{gt} to $I_{noise}^{(t)}$ with random t step in the range of T . Then $I_{noise}^{(t)}$, I_l and I_r are concatenated together to comprise I_{concat} with 7 channels. We build an U-Net to predict the noise with 3 channels that added onto the I_{gt} . We propose a novel conditional noise prediction proxy task for the pre-training stage by introducing I_l and I_r as additional conditional inputs. The information of t is embedded using the time embedding function and is transmitted into all convolutional blocks of both encoder and decoder in the U-Net. As is shown in Fig. 2(b) and (c), the convolutional blocks in encoder and decoder have the same structure containing a residual block followed by an attention block. It is worth to note that we use dotted box to represent attention block is a selective usage. At the shallow layers of the encoder, we do not use attention block due to the large dimension of feature maps. The encoder and decoder equipped with multi-head self-attention make our model efficiently capture global and local features in different convolutional layers. There is a linear projection module to map the embedded time information t_{emb} to the one with the same size as the feature map after the first convolution operation. We use the common add operation

to encode the time information into the convolutional block. The attention block is not only used as a sub-block in the encoder and decoder of U-Net, but also as an independent block in the bottleneck of U-Net. The detailed structure of attention block is illustrated in Fig. 2(d). The use of attention block can make model learn long-range features and multi-scale features which is essential for our colorization task. Then we use denoising function to transfer the predicted noised $N_{pred}^{(t)}$ combined with $I_{noise}^{(t)}$ to the generated colored image I_{gen} .

4.2 Line Extraction

Due to the lack of a large amount of hand-drawn line data, it is quite time-consuming and laborious to expand the data volume by hand-drawn method. We use XDoG [50] line style as the intermediate representation of line drawings during the training and inference stage. Any input line drawings created by artist will be automatically converted to XDoG style to fit the model. For this reason, we build a line extraction module integrated in AnimeDiffusion. During the training stage, reference color images are used as input to extract the line drawings, and during the inference stage, hand-drawn line drawings are used as input to transform the line draft style into XDoG style.

For given image x , the line extractor is described in the form [50]

$$S_{\sigma,k,p}(x) = (1 + p) \cdot G_{\sigma}(x) - p \cdot G_{k\sigma}(x) \quad (1)$$

where G_{σ} and $G_{k\sigma}$ is Gaussian convolution operation, σ is the variance of Gaussian convolution kernel, k is scaling ratio of the variance between two convolution, and p is used to control the edge emphasis lines.

We need a line extractor in which line extraction results are as close as possible to the effect of a painter’s hand-drawn line. The variance σ of the Gaussian convolution has a significant effect on the line thickness of the line drawing, and we choose σ to be 0.3 to get a reasonable line width. We hired a professional artist to draw line art for some of the images in our dataset. For a color image I_c , its line drawing extraction result is $S_{k,p}(I_c)$, and the hand-drawn image by the painter is represented as $H(I_c)$. The objective of parameters for the line extractor is

$$\arg \min_{k,p} \sum_i \|S_{k,p}(I_i) - H(I_i)\|_2^2 \quad (2)$$

4.3 Training Strategy

We design a hybrid training strategy to train AnimeDiffusion, which consists of classifier-free guidance pre-training stage and an image reconstruction guidance fine-tuning stage. This training strategy separates the denoising task from the image reconstruction task, makes the network learning a specific task at each stage, and is more beneficial to network training and weight updating. Each training step will be introduced in detail in section 4.3.1 and 4.3.2.

4.3.1 Pre-training Stage

During the classifier-free guidance pre-training stage, AnimeDiffusion mainly learns denoising ability. As shown in Fig. 2, the original image I_{gt} goes through a forward diffusion process which is a Markov chain since it adds Gaussian noise to I_{gt} and obtains noisy image I_{noise}^t for time step t iteratively. Each step of the forward process is a Gaussian transition.

$$q(I_{noise}^{(t)}|I_{noise}^{(t-1)}) = \mathcal{N}(I_{noise}^{(t)}; \sqrt{1 - \beta_t}I_{noise}^{(t-1)}, \beta_t\mathbf{I}) \quad (3)$$

where β_t is variance schedule at time step t . The forward process of the diffusion model represents the addition of noise from step 0 to t . For the cumulative t steps of noise addition $I_{noise}^{(1:T)}$, the marginal distribution is

$$q(I_{noise}^{(1:T)}|I_{gt}) = \prod_{t=1}^T q(I_{noise}^{(t)}|I_{noise}^{(t-1)}) \quad (4)$$

Under the condition of equation 3, the marginal distribution of each forward step is a standard Gaussian distribution

$$q(I_{noise}^{(t)}|I_{gt}) = \mathcal{N}(I_{noise}^{(t)}; \sqrt{\bar{\alpha}_t}I_{gt}, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (5)$$

where $\bar{\alpha}_t = \prod_{i=1}^t (1 - \beta_i)$. After t times iteration, the result latent variable $I_{noise}^{(t)}$ can be simplified as

$$I_{noise}^{(t)} = \sqrt{\bar{\alpha}_t}I_{gt} + \sqrt{1 - \bar{\alpha}_t}\epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (6)$$

The training objective of the model is to predict the noise $N_{pred} = \epsilon_{\theta}(I_l, I_r, I_{noise}^{(t)}, t)$ with given noised data point $I_{noise}^{(t)}$, time step t and condition I_l and I_r , and optimizing the objective

$$\mathbb{E}_{I_{gt} \sim q(I_{gt}), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), I_l, I_r, t} \|\epsilon - \epsilon_{\theta}(I_l, I_r, I_{noise}^{(t)}, t)\|_p^p \quad (7)$$

For a simplified description, we subsequently use ϵ_{θ} represent $\epsilon_{\theta}(I_l, I_r, I_{noise}^{(t)}, t)$. Palette [41] indicates that the L2 norm can capture the output distribution more faithfully, and we adopt $p = 2$ in our pre-training stage.

We used a model of U-Net with attention blocks to predict the noise ϵ added in equation 6. The U-Net needs to accept line drawing I_l , reference image I_r , and noisy image I_{noise}^t as inputs. Considering the high spatial consistency between line drawing and color images, we concatenate the above three in the channel dimension. And the subsequent experimental in section 5 results demonstrate that without a complex feature fusion mechanism, the model we proposed can achieve the semantic correspondence between the line drawings and the reference maps. Our previous work [19] for the first time emphasized the importance of semantic correspondence for reference-based line drawing colorization. The method in this paper uses clever algorithm design to achieve more amazing results in accurate semantic correspondence, especially in the region of anime character face.

4.3.2 Fine-tuning Stage

After training $\epsilon_\theta(I_l, I_r, I_{noise}^{(t)}, t)$, diffusion models inference through the learned reverse process. Since the result distribution of forward process $p(I_{noise}^{(T)})$ approximates a standard Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$, the sampling process starts from pure Gaussian noise, followed by T rounds of denoising. On the one hand, training diffusion models often requires a large batch size and long iteration rounds with large computing consumption. On the other hand, we want to strike a balance between the diversity and accuracy of generated results. Based on the pre-trained model already having some denoising ability, we introduce the image reconstruction guidance fine-tuning stage to improve the generation ability of AnimeDiffusion. Since the great diversity of the generated results of the diffusion model, as the training iterations grow, the quality of generated images is affected less by the input noise and more by the guidance condition. Therefore we input the noise generated according to the reference image instead of random noise in the fine-tuning and inference process of the model. According to equation 6, I_{gt} is estimated as

$$\tilde{I}_{gt} = \frac{1}{\sqrt{\bar{\alpha}_t}} (I_{noise}^{(t)} - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(I_l, I_r, I_{noise}^{(t)}, t)) \quad (8)$$

The mean value of reverse process $p_\theta(I_{noise}^{(t-1)} | I_{noise}^{(t)}, I_l, I_r)$ is parameterize as

$$\tilde{\mu}_\theta(I_{noise}^{(t)}, t) = \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \cdot \tilde{I}_{gt} + \frac{(1 - \bar{\alpha}_{t-1}) \sqrt{\bar{\alpha}_t}}{1 - \bar{\alpha}_t} \cdot I_{noise}^{(t)} \quad (9)$$

With the estimation of $\mu_\theta(I_{noise}^{(t-1)}, t)$, each iteration of reverse process is

$$I_{noise}^{(t-1)} = \tilde{\mu}_\theta(I_{noise}^{(t)}, t) + \sigma_t^2 \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (10)$$

where σ_t is the sampling variance with $\sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$.

The time consumption and space consumption for fine-tuning by directly adopting the reverse process of DDPM is vast, so we use DDIM [40] as our denoising function, which is an alternative non-Markov chain denoising process with different sampling or reverse process

$$I_{noise}^{(t'-1)} = \sqrt{\bar{\alpha}_{t'-1}} \tilde{I}_{gt} + \sqrt{1 - \bar{\alpha}_{t'-1} - \eta \sigma_{t'}^2} \epsilon_\theta + \eta \sigma_{t'} \epsilon \quad (11)$$

DDIM obtains a sub-sequence of time $[0, T)$, where t' is sampled time sequence [40]. η is a hyper-parameter that controls whether noise is added during the reverse process. If η is set to 0, the process of image generation is deterministic.

Since both the forward and reverse processes of DDPM are random, the colorization results are different for the same sample. DDIM provides a deterministic reverse sampling strategy, but due to the different initial noise, there is no guarantee that the image reconstruction can be completed with the original image as the reference image. To fully utilize the image synthesis performance of the diffusion model for image processing purposes, we borrow the deterministic forward process from DiffusionCLIP [22] in our fine-tuning stage. According to equation 8 and 11,

DDIM is considered as an Euler method to solve an ordinary differential equation (ODE)

$$d \frac{I_{noise}^{(t)}}{\sqrt{\bar{\alpha}_t}} = \epsilon_\theta \cdot d \sqrt{\frac{1 - \bar{\alpha}_t}{\bar{\alpha}_t}} \quad (12)$$

The above ODE holds in a finite number of steps, and to obtain an accelerated forward process, we use the same sampling time series t' as the reverse process of DDIM. The recursive relation from $I_{noise}^{(t')}$ to $I_{noise}^{(t'+1)}$ is simplified as

$$I_{noise}^{(t'+1)} = \sqrt{\bar{\alpha}_{t'+1}} \tilde{I}_{gt}(I_{noise}^{(t')}) + \sqrt{1 - \bar{\alpha}_{t'+1}} \epsilon_\theta \quad (13)$$

$\tilde{I}_{gt}(I_{noise}^{(t')})$ represents the ground truth estimation function of equation 8. Accordingly, to obtain the deterministic reverse process, we apply η as 0 in the reverse process of DDIM.

After several iterations, the denoising function finally generates the colored image I_{gen} . We calculate the mean square error (MSE) between I_{gen} and I_{gt} , in order to constrain the generated image as close to the original image as possible in pixel level. It is worth noted that the dotted arrow in Fig. 2 represents reverse gradient propagation to update the parameters of U-Net during the fine-tuning stage.

$$L_{rec} = \mathbb{E}_{I_{gt} \sim q(I_{gt}), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \| I_{gen} - I_{gt} \|_p^p \quad (14)$$

Algorithm 1 Anime Diffusion Fine-tuning

Input: ϵ_θ pretrained noise prediction model,
 \mathcal{I} training anime face images,
 S number of sampling steps

Output: $\tilde{\epsilon}_\theta$ fine-tuned noise prediction model

- 1: Initialize list of forward noisy images $\hat{\mathcal{I}}$;
- 2: Sampling a increasing sub-sequence $\mathcal{T}' = \{t'_0, \dots, T'\}$ of length S ;
- 3: **for** I_{gt} **in** \mathcal{I} **do**
- 4: **for** t' **in** \mathcal{T}' **do**
- 5: Calculate $\tilde{\epsilon} \leftarrow \epsilon_\theta(I_l, I_r, I_{noise}^{(t')}, t')$;
- 6: Predict the ground truth $I_{gt}(I_{noise}^{(t')}, t')$;
- 7: Forward step $I_{noise}^{(t'+1)} \leftarrow \sqrt{\bar{\alpha}_{t'+1}} \tilde{I}_{gt} + \sqrt{1 - \bar{\alpha}_{t'+1}} \tilde{\epsilon}$;
- 8: **end for**
- 9: Update $\hat{\mathcal{I}}$;
- 10: **end for**
- 11: **for** $I_{noise}^{T'}$ **in** $\hat{\mathcal{I}}$ **do**
- 12: **for** t' **in** reverse(\mathcal{T}') **do**
- 13: Calculate $\tilde{\epsilon} \leftarrow \epsilon_\theta(I_l, I_r, I_{noise}^{(t')}, t')$;
- 14: Predict the ground truth $I_{gt}(I_{noise}^{(t')}, t')$;
- 15: Reverse step $I_{noise}^{(t'-1)} \leftarrow \sqrt{\bar{\alpha}_{t'-1}} I_{gt} + \sqrt{1 - \bar{\alpha}_{t'-1}} \tilde{\epsilon}$;
- 16: **end for**
- 17: Update reconstruction loss L_{rec} ;
- 18: Gradient step $\nabla_{\tilde{\epsilon}_\theta} L_{rec}$
- 19: **end for**

We perform a small number of steps of fine-tuning based on the pre-trained model, as shown in algorithm 1. In summary, we perform our proposed hybrid training strategy to reduce computation consumption of training and inference significantly in comparison with the traditional manner of diffusion method. The training objective of the

pre-training stage is to obtain the solution of equation 12, the derivative of the path from the initial distribution to the target distribution at time step t . When applying equation 11 for inference, the distribution of the sub-sequences of time step $\{t'\}$ will affect the colorization results due to the deviation between the predicted derivatives ϵ_θ and the true path from time step t' to $t' - 1$. Therefore, we fix the sub-sequence t' in the fine-tuning stage and select a smaller number of sampling steps to make our model save inference time. Fine-tuning the pre-trained model so that the noise prediction at time step t' is closer to the direction pointing to the next time step $t' - 1$, rather than accurately predicting the derivative.

4.4 User Interface

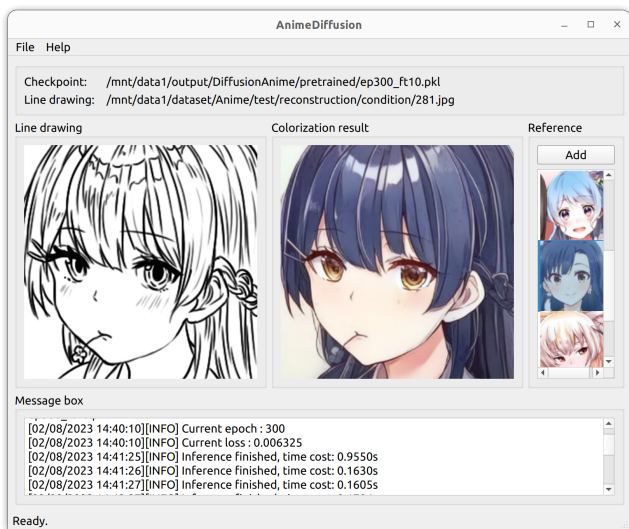


Fig. 3. User interface of AnimeDiffusion for coloring anime face line drawings.

As is shown in Fig. 3, a user interface is developed for users to perform line drawing colorization by our AnimeDiffusion. Users only need to provide the line drawing and the reference image as two inputs of AnimeDiffusion, then it can one-key automatically complete the coloring process to generate colored results without additional human intervention. The elapsed time of the colorizing operation is printed at the bottom of the interface. Owing to the adoption of the DDIM acceleration sampling strategy and the fine-tuning stage, our method generated high-quality colorization results with only a small number of sampling steps. The colorization time consumption for a line drawing can be controlled within 0.2s on a machine equipped with RTX 4090, excluding the initialization period. Our end-to-end AnimeDiffusion model can be directly integrated into the practical colorization pipeline in the animation creation industry. In contrast to other diffusion methods [24], [41], our method can accurately edit face line drawings according to reference images. Especially as shown in Fig.1, the artist added a lot of detail lines to the hair, our AnimeDiffusion can also add detailed textures instead of just flat color to the hair after coloring.

5 EXPERIMENTS

5.1 Dataset

In this paper, we focus on the anime face line drawing colorization task. In order to train AnimeDiffusion, a large dataset of anime face images is necessary, and the image resolution should not be too low in order to adequately express the color and detail information of the face. However, there is no dataset that meets our needs and can be directly used for training. So we build a benchmark dataset for anime face line drawing colorization. All anime character images are collected from Danbooru2020 [25], which is a large-scale anime image database with 4.2m+ images. According to our task requirements, we only cut out the face part. After simple manual alignment and denoising operation, a total of 31696 training data and 579 testing data are produced. Due to limitations in GPU memory and model computing efficiency, all images are resized to 256×256 resolution. To simulate the manual line drawing style by artists and generate paired line drawing images, we use XDoG [50] to extract line drawings from colored anime images and set the parameters of XDoG algorithm with $\phi = 1 \times 10^9$ to keep a step transition at the border of lines in line drawings. We randomly set σ to be 0.3/0.4/0.5 to get different levels of line thickness, which generalizes AnimeDiffusion on various line widths to avoid overfitting. And we set $p = 9, k = 4.5, \epsilon = 0.01$ in XDoG. As mentioned before, in the practical colorization scenario, there are large space discrepancies between the target line drawing and reference image. In order to make AnimeDiffusion learn accurate semantic correspondence ability during training and avoid learning trivial solution by directly using pixel aligned training data, we randomly set the parameters of TPS transformation on colored reference images when loading training data, i.e. each image of one batch data will have different geometry distortions. To some extent, this is a data augmentation trick.

5.2 Implementation Details

We implement our AnimeDiffusion model based on the PyTorch framework, and it is trained on 1 NVIDIA A100 GPU. All input image size is fixed at 256×256 . For the diffusion hyper-parameters setting, we use a linear noise schedule of $(1e^{-6}, 1e^{-2})$ with 1000 time steps. We pre-train the model with a batch size of 32 for 300 epochs, and we don't find overfitting, and we fine-tune the model with a batch size of 4 for 1 epoch. On our devices, the pre-training stage takes 40 hours and the fine-tuning stage takes 110 minutes. We apply the Adam optimizer with a learning rate of $1e^{-5}$ for both of the above stages. Besides, we have no other hyper-parameters to adjust, like learning rate decay or warm-up schedule.

5.3 Qualitative Evaluation

We compare our AnimeDiffusion with another three state-of-the-art GANs-based methods. Lee et al. [14] proposed the self-augmented supervised training strategy and designed a model with an attention based Spatial Correspondence Feature Transfer (SCFT) module. We regard it as the baseline for line drawing colorization task. Li et al. [15] designed

a Stop-Gradient Attention (SGA) module to eliminate the gradient conflict among attention branches. Cao et al. [19] proposed an attention-aware improved method based on [14], which focus on the anime line drawing colorization task. Since the variety of anime characters' faces and combine with the actual needs of anime character creation, we set two anime cases separately including anime face with homochromatic pupils and anime face with heterochromatic pupils. The latter is a very challenging case, which needs model require a high precision extraction of local features and semantic correspondence. To the best of our knowledge, we are the first learning based work can generate results with accurate color in pupils according to the reference image with no extra eyes segmentation label [8] or pupil position estimation network [52].

For homochromatic pupils case, we show detailed comparison results in Fig. 4. Yellow region shows that AnimeDiffusion can recognize the ear semantic information from line drawing and inject the right color the same as face. Green region shows that AnimeDiffusion can maintain the light-reflecting effect compared with the original color image (Fig. 4(c)). Blue region indicates that AnimeDiffusion can accurately transfer the color information from reference image (Fig. 4(a)) into the line drawing (Fig. 4(b)) in the eyes part. However, the other three methods have flaws in the areas we have marked in three colors. Since the diffusion models are based on maximum likelihood estimation method that can estimate the probability density more accurately than GANs-based method. So our results are much clearer than the other three methods. The performance of feature aggregation modules in other three models is not the best, so the coloring effect is defective. Combined with our proposed training strategy, denoising and reconstruction tasks are separated during the training procedure, making the network training more stable. Therefore, our model acquires better detailed features capture ability.

More comparison results are illustrated in Fig. 5, given line drawings and reference images, AnimeDiffusion generates colored results with accurate color and good semantic correspondence. The image is clear without noise, and the color texture is smooth and soft. Especially in the eyes, the color is very precise, and the sense of light in the eyes is kept very good and full of charm. By contrast, Lee et al. [14] generates results with color bleeding and wrong semantic correspondence, the color texture is rough and detail information is unclear. Li et al. [15] generates results with inaccurate color. Cao et al. [19] generates results with sharpen image quality and wrong color in eyes. In general, the quality of images produced by GANs-based methods is not very stable and random flaw sometimes occurs. Our results are of high quality with beautiful color and rich details. Compared with the other three GANs-based methods, the image texture quality has been significantly improved.

For heterochromatic pupils case which is a very challenging task of anime face line drawing colorization. As is illustrated in Fig. 6, AnimeDiffusion can generate results with accurate color in different pupils according to the reference image. Although Lee et al. [14] produces some results with heterochromatic pupils, but the overall quality of results is not high, there are still color bleeding and inaccurate semantic correspondence. Li et al. [15] generates

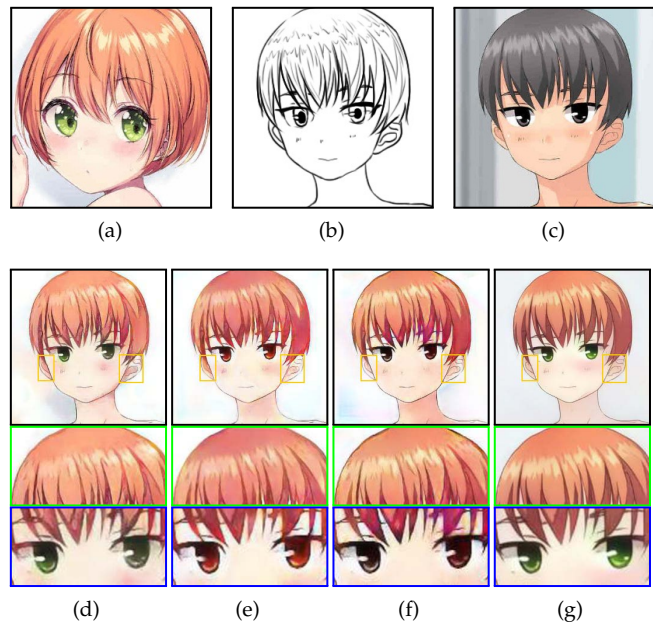


Fig. 4. Detailed comparison of colorization results. (a) reference image, (b) line drawing, (c) original color image, (d) Lee et al. [14], (e) Li et al. [15], (f) Cao et al. [19], (g) AnimeDiffusion

results with distorted color globally and inaccurate color in pupils. Cao et al. [19] fail to handle the heterochromatic pupils case, but the image sharpness and semantic information are still accurate. As heterochromatic pupil is a fine-grained feature in the image space, and GAN is not accurate in data distribution modeling, the three SOTA GANs-based methods [14], [15], [19] cannot handle it well. In contrast, our diffusion-based solution takes advantage of its precise data modeling property, combined with the use of multi-scale feature self-attention modules. Therefore, pupils can be colorized accurately according to the reference images without introducing additional processing modules.

5.4 Quantitative Evaluation

5.4.1 Evaluation Metrics

We mainly use three evaluation metrics for quantitative comparison AnimeDiffusion with other methods. The popular Fréchet Inception Distance (FID) is used to assess the generation ability of algorithms in perceptual level. Besides measuring the perceptual credibility, we also adopt Peak Signal-to-Noise Ratio (PSNR) and Multi-Scale Structural Similarity Index Measure (MS-SSIM) to evaluate the image reconstruction ability of algorithms in pixel level. We design two kinds of colorization tasks respectively including self-reference reconstruction and random-reference colorization to analyze the colorization performance of AnimeDiffusion.

5.4.2 Self-reference Reconstruction

For self-reference colorization, the line drawing and reference image are paired, ideally the colorized output should be exactly the same as the reference image. We directly use our paired testing data for conducting self-reference colorization. In fact, during the training phase, the main task

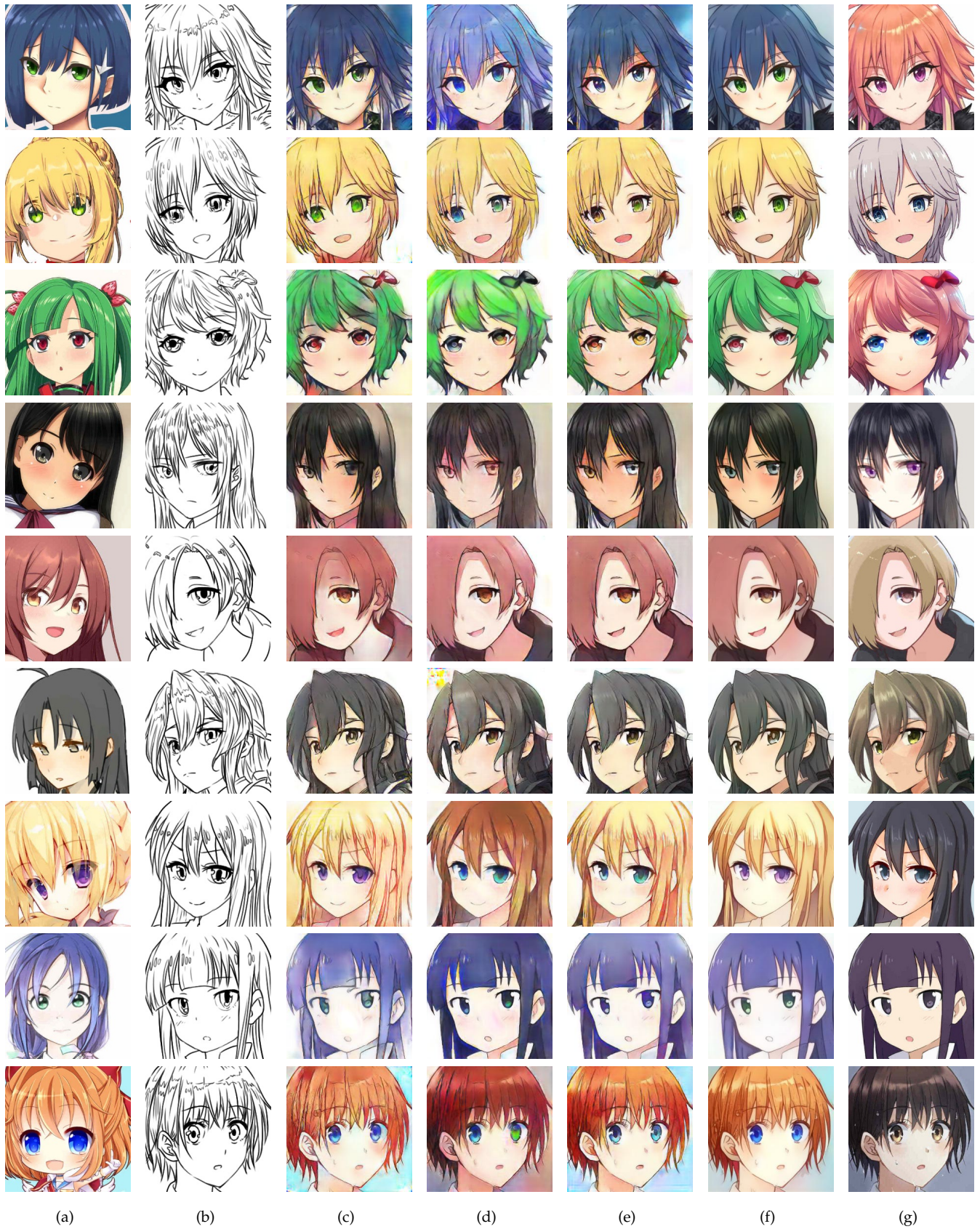


Fig. 5. Qualitative comparison for anime face with homochromatic pupils. (a) reference images, (b) line drawings, (c) Lee et al. [14], (d) Li et al. [15], (e) Cao et al. [19], (f) AnimeDiffusion, and (g) original color images.



Fig. 6. Qualitative comparison for anime face with heterochromatic pupils. (a) reference images, (b) line drawings, (c) Lee et al. [14], (d) Li et al. [15], (e) Cao et al. [19], (f) AnimeDiffusion, and (g) original color images.

of AnimeDiffusion is to do the image self-reconstruction, through this proxy task, the network can learn the colorization ability. For fairness, we train AnimeDiffusion and other three GANs-based methods sufficiently to compute PSNR and MS-SSIM. As is shown in Table 1, AnimeDiffusion acquires the best image reconstruction performance.

5.4.3 Random-reference Colorization

For the random-reference colorization, it is more like the common practical usage when using reference-based colorization method. We shuffle all the reference images in our testing data, then use unpaired line drawings and reference images to perform random-reference colorization. Using the total 579 generated images and 579 reference images to compute FID score. A smaller FID indicates that the distribution of the colored images is closer to the reference images and indicates that the model with wonderful generation ability. As is shown in Table 1, AnimeDiffusion shows better generation ability than other three GANs-based methods. One thing needs to note is that although Cao et al. [19] shows little poor image reconstruction performance than Lee et al. [14], but shows better generation ability than Lee et al. [14] and Li et al. [15], that is because Lee et al. [14] just learns a trivial solution. This point is also discovered and discussed in Li et al. [15].

TABLE 1
Quantitative Comparison between AnimeDiffusion and Other Three SOTA GANs-based Methods

Method	PSNR↑	MS-SSIM↑	FID↓
Lee et al. [14]	23.8901	0.9224	57.19
Li et al. [15]	18.6347	0.8209	49.33
Cao et al. [19]	19.7746	0.8388	46.39
AnimeDiffusion	25.4658	0.9596	44.19

5.5 Ablation Study

We perform extensive ablation experiments to verify the effectiveness of our designed fine-tuning strategy when training AnimeDiffusion. We find that the denoising model obtained by classifier-free guidance pre-training can generate images with high diversity, but this diversity also means that the colorization results are unstable since randomness is introduced by Gaussian noise.

We believe that the pre-trained model already acquires the ability to capture the line structure and can inject different colors in the corresponding area according to the reference image. Fine-tuning is only to eliminate color gaps due to insufficient pre-training. To validate our idea, we perform image reconstruction test and reference-based line drawing



Fig. 7. Image reconstruction test for ablation study. (a) Ground truth images, (b) Results without fine-tuning, (c) Results with fine-tuning for 1 epoch, (d) Results with fine-tuning for 10 epochs.

colorization test respectively. We fine-tune AnimeDiffusion with 1 epoch and 10 epochs with a batch size of 4 for comparison.

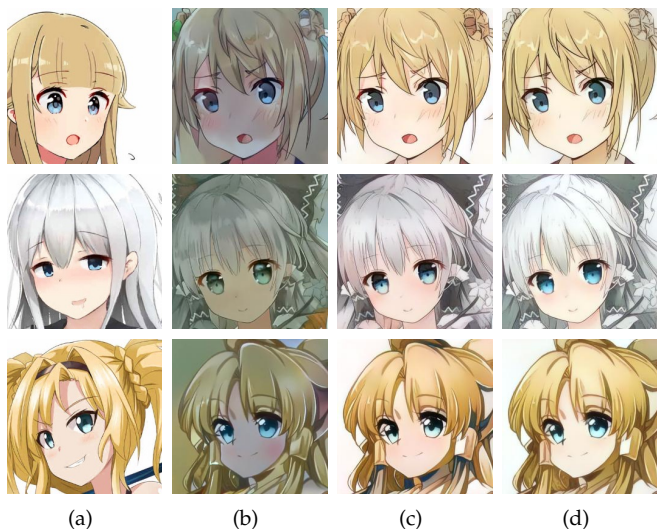


Fig. 8. Reference-based line drawing colorization test for ablation study. (a) Reference images, (b) Results without fine-tuning, (c) Results with fine-tuning for 1 epoch, (d) Results with fine-tuning for 10 epochs.

We perform the image reconstruction test by distorting the ground truth image that serves as the reference image. Comparison results are shown in Fig. 7. The results without fine-tuning can reconstruct the overall structure of original image, but the color difference is obvious. After adding the image reconstruction loss to the fine tuning, the effect is significantly improved. The results of different fine-tuning schemes are not obvious in terms of visual differences

For reference-based line drawing colorization test, we show results in Fig. 8. Although the model without fine-tuning can distinguish regions that need different colors, there is still a color gap between the generated images

and the reference images, and the result colored image looks dimmer. We think the model needs more training to generate results with accurate color information. However, training the diffusion model with classifier-free guidance is time-consuming, so we briefly fine-tune the model and get much better results.

We also compute PSNR, MS-SSIM and FID quantitative index for quantitative comparison. Results are shown in Table 2. After 10 epochs of fine-tuning, AnimeDiffusion continues to gain in FID score but little improvement in PSNR and MS-SSIM score. In fact, fine-tuning the model for 1 epoch is enough to have good colorization performance. This also validates that our fine-tuning is mainly to fix color bias based on the pre-trained model with fundamental generation ability. Our designed hybrid training strategy can make model learn better colorization ability and save the training time cost.

TABLE 2
Quantitative Evaluation for Ablation Study Results

AnimeDiffusion	PSNR↑	MS-SSIM↑	FID↓
Without Fine-tuning	12.4234	0.8079	55.1841
Fine-tuning (1 epoch)	25.4658	0.9596	44.1876
Fine-tuning (10 epochs)	25.8992	0.9600	40.4392

5.6 User Study

It is generally challenging to evaluate the visual quality of images, in particular for line drawing colorization. We randomly select 50 line drawings and 50 reference images to perform reference-based line drawing colorization using AnimeDiffusion compared with Lee et al. [14], Li et al. [15] and Cao et al. [19]. Then we conduct a user study to let participants compare colored results of these four methods, participants need subjectively evaluate the colored results according to the reference images and original color images of line drawings to choose the best one from four choices, along with a brief description of why they think this result is the visually best. A user interface of our user study is shown in Fig. 9. 20 participants take part in the user study and the percentage of each method chosen as the best is shown in Table 3. It is indicated that AnimeDiffusion has absolute advantages in human visual evaluation.

TABLE 3
User Study Result

Methods	Percentage of chosen as best
Lee et al. [14]	13.0%
Li et al. [15]	17.4%
Cao et al. [19]	21.0%
AnimeDiffusion	48.6%

6 APPLICATION

6.1 Famous Anime Character Recolorization

Although we aim to train AnimeDiffusion to perform single line drawing colorization, however, in animation creation industry, reference-based colorization method can also be

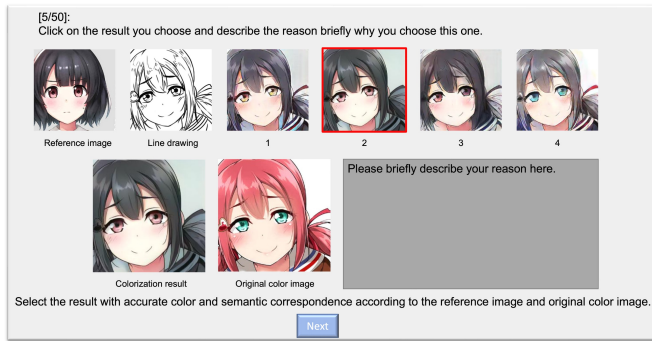


Fig. 9. A user interface of our user study.

used to recolorize a series of images or even consecutive video frames of the same anime character. Sometimes during the creation process, the same character appears different colors in different images, and using the recolorization technique, it is convenient to unify the same character’s colors based on one reference image. We apply AnimeDiffusion to perform recolorization task. In Fig. 10, for each anime character, the top row is original colored image and the bottom row is our recolorized results. In fact, we first convert original colored images to line drawings using XDoG extractor, then they are combined with the reference image together and are fed into AnimeDiffusion to generate colored results. As is shown in Fig. 10, according to one reference image, the other images of the same character can be recolorized with the same color style and accurate semantic information.

6.2 Original Anime Character Colorization

We collaborate with professional artists and use our AnimeDiffusion to assist their work for original line drawings colorization. We invite professional artist to drawing line drawings attached with one reference color image. Since AnimeDiffusion pre-trained on our collected dataset has learned the semantic information of anime character face, it can generalize well to other hand-drawn characters, as is demonstrated in Fig. 11. With our developed user interface, artist can easily do batch colorization of the same character according to one reference image. This greatly saves the artists’ creation time and helps them to complete the creation more efficiently.

6.3 Fashion Illustration Sketch Colorization

We also extend AnimeDiffusion to colorize fashion illustration sketches. Since fashion illustration is the same as animation, it is also the first to outline the line draft, and then fill in the color. We regard fashion illustration as another type of animation. As is shown in Fig. 12, given one color illustration and one sketch, AnimeDiffusion can generate colorization results which extend the range of accurate semantic correspondence to half of the body. We can not only keep the accurate color of the face, but also have good control over the clothing and torso, even the color of skin can be accurately distinguished. Fashion designers can use



Fig. 10. Illustration of famous anime character recolorization. Recolorization results have a uniform color style. The two anime characters are Hoshizora Rin and Sonoda Umi of LoveLive.



Fig. 11. Illustration of original anime character colorization. The anime character with an exaggerated hairstyle is Little Yuyuan, which is created by Ms. Yuwen Wang.

our AnimeDiffusion user interface to easily colorize hand-drawn fashion sketches, which are used for the follow-up process of garment pattern making.

7 CONCLUSION AND FUTURE WORK

In this paper, we propose AnimeDiffusion, the first diffusion model tailored for anime face line drawing colorization. In order to train AnimeDiffusion we build a benchmark dataset for research purpose and also fill the gap of no available high resolution anime face dataset to evaluate line drawing colorization algorithms. To handle the high computation consumption problem of diffusion models, we design a novel hybrid training strategy which separates the image denoising task and image reconstruction task. Through extensive experiments and a user study, AnimeDiffusion has demonstrated better performance both qualitatively and quantitatively, outperforming other state-of-the-art GANs-based methods, with higher image quality and semantic



Fig. 12. Illustration of fashion illustration sketch colorization. (a) Reference images, (b) Sketches, (c) Colorization results

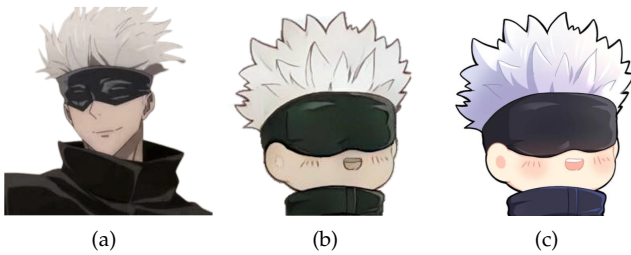


Fig. 13. Limitation of our approach. (a) Reference images, (b) Colorization results, (c) Ground truth. Due to the large stylistic differences between the reference image and the line drawing, the color of the mouth is not available in the reference image, while the teeth are not correctly identified in our model in the line drawing, and the color of the teeth is not accurately reflected in our colorization results.

color information. To the best of our knowledge, AnimeDiffusion is the first learning based work can accurately colorize anime face line drawing with heterochromatic pupils according to reference color image, without other special module for processing eyes or pupils in anime face.

However, there is a limitation in our method. Our model uses paired training data in training, and there is some style correlation between the reference image and the line drawings. For special style line drawings, such as Chibi cartoons as shown in Fig 13(c), if the corresponding semantic information does not exist in the reference image, the

colorization result of the line drawing may appear to be inconsistent with the real image. In the future work, We will work on multi-modal input line drawing colorization such as combining text information and reference image together to make the interactive way of colorization more rich. This will greatly reduce the manual tasks of animators and improve the creation efficiency and colorization effect of the animation creation industry.

ACKNOWLEDGMENTS

We thank Mr. Henry Tian, Ms. Mandy Wong, Ms. Rachel Liu and Mr. Leonard Chen for their help with anime knowledge and image samples selection. We thank Ms. Xiao Meng and Ms. Yuwen Wang for helping us create wonderful hand-painted line drawings.

REFERENCES

- [1] X. Dong, W. Li, X. Hu, X. Wang, and Y. Wang, "A colorization framework for monochrome-color dual-lens systems using a deep convolutional network," *IEEE Transactions on Visualization and Computer Graphics*, vol. 28, no. 3, pp. 1469–1485, 2022.
- [2] Y. Xiao, J. Wu, J. Zhang, P. Zhou, Y. Zheng, C.-S. Leung, and L. Kavan, "Interactive deep colorization and its application for image compression," *IEEE Transactions on Visualization and Computer Graphics*, vol. 28, no. 3, pp. 1557–1572, 2022.
- [3] R. Hu, Z. Ye, B. Chen, O. van Kaick, and H. Huang, "Self-supervised color-concept association via image colorization," *IEEE Transactions on Visualization and Computer Graphics*, vol. 29, no. 1, pp. 247–256, 2023.
- [4] F. Fang, T. Wang, T. Zeng, and G. Zhang, "A superpixel-based variational model for image colorization," *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, no. 10, pp. 2931–2943, 2020.
- [5] Y. Qu, T.-T. Wong, and P.-A. Heng, "Manga colorization," *ACM Transactions on Graphics (TOG)*, vol. 25, no. 3, pp. 1214–1220, 2006.
- [6] C. Furusawa, K. Hiroshiba, K. Ogaki, and Y. Odagiri, "Comicolorization: semi-automatic manga colorization," in *SIGGRAPH Asia 2017 Technical Briefs*, 2017, pp. 1–4.
- [7] D. Šykora, J. Dingliana, and S. Collins, "Lazybrush: Flexible painting tool for hand-drawn cartoons," in *Computer Graphics Forum*, vol. 28, no. 2. Wiley Online Library, 2009, pp. 599–608.
- [8] S.-Y. Chen, J.-Q. Zhang, L. Gao, Y. He, S. Xia, M. Shi, and F.-L. Zhang, "Active colorization for cartoon line drawings," *IEEE Transactions on Visualization and Computer Graphics*, vol. 28, no. 2, pp. 1198–1208, 2022.
- [9] D. Varga, C. A. Szabo, and T. Sziranyi, "Automatic cartoon colorization based on convolutional neural network," in *Proceedings of the 15th International Workshop on Content-Based Multimedia Indexing*, 2017, pp. 1–6.
- [10] Y. Ci, X. Ma, Z. Wang, H. Li, and Z. Luo, "User-guided deep anime line art colorization with conditional adversarial networks," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 1536–1544.
- [11] L. Zhang, C. Li, T.-T. Wong, Y. Ji, and C. Liu, "Two-stage sketch colorization," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 6, pp. 1–14, 2018.
- [12] H. Kim, H. Y. Jhoo, E. Park, and S. Yoo, "Tag2pix: Line art colorization using text tag with secant and changing loss," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9056–9065.
- [13] C. Zou, H. Mo, C. Gao, R. Du, and H. Fu, "Language-based colorization of scene sketches," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–16, 2019.
- [14] J. Lee, E. Kim, Y. Lee, D. Kim, J. Chang, and J. Choo, "Reference-based sketch image colorization using augmented-self reference and dense semantic correspondence," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5801–5810.
- [15] Z. Li, Z. Geng, Z. Kang, W. Chen, and Y. Yang, "Eliminating gradient conflict in reference-based line-art colorization," *arXiv preprint arXiv:2207.06095*, 2022.

- [16] L. Zhang, Y. Ji, X. Lin, and C. Liu, "Style transfer for anime sketches with enhanced residual u-net and auxiliary classifier gan," in *2017 4th IAPR Asian conference on pattern recognition (ACPR)*. IEEE, 2017, pp. 506–511.
- [17] T.-H. Sun, C.-H. Lai, S.-K. Wong, and Y.-S. Wang, "Adversarial colorization of icons based on contour and color conditions," in *Proceedings of the 27th ACM International Conference on Multimedia*, 2019, pp. 683–691.
- [18] X. Liu, W. Wu, C. Li, Y. Li, and H. Wu, "Reference-guided structure-aware deep sketch colorization for cartoons," *Computational Visual Media*, vol. 8, no. 1, pp. 135–148, 2022.
- [19] Y. Cao, H. Tian, and P. Mok, "Attention-aware anime line drawing colorization," *arXiv preprint arXiv:2212.10988*, 2022.
- [20] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International Conference on Machine Learning*. PMLR, 2015, pp. 2256–2265.
- [21] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon, "Sdedit: Guided image synthesis and editing with stochastic differential equations," in *International Conference on Learning Representations*, 2021.
- [22] G. Kim, T. Kwon, and J. C. Ye, "Diffusionclip: Text-guided diffusion models for robust image manipulation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2426–2435.
- [23] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 684–10 695.
- [24] L. Zhang and M. Agrawala, "Adding conditional control to text-to-image diffusion models," 2023.
- [25] Anonymous, D. community, and G. Branwen, "Danbooru2020: A large-scale crowdsourced and tagged anime illustration dataset," <https://www.gwern.net/Danbooru2020>, January 2021, accessed: DATE. [Online]. Available: <https://www.gwern.net/Danbooru2020>
- [26] K. Sato, Y. Matsui, T. Yamasaki, and K. Aizawa, "Reference-based manga colorization by graph correspondence using quadratic programming," in *SIGGRAPH Asia 2014 Technical Briefs*, 2014, pp. 1–4.
- [27] M. Shi, J.-Q. Zhang, S.-Y. Chen, L. Gao, Y. Lai, and F.-L. Zhang, "Reference-based deep line art video colorization," *IEEE Transactions on Visualization and Computer Graphics*, 2022.
- [28] Z. Dou, N. Wang, B. Li, Z. Wang, H. Li, and B. Liu, "Dual color space guided sketch colorization," *IEEE Transactions on Image Processing*, vol. 30, pp. 7292–7304, 2021.
- [29] A. Maejima, H. Kubo, S. Shinagawa, T. Funatomi, T. Yotsukura, S. Nakamura, and Y. Mukaigawa, "Anime character colorization using few-shot learning," in *SIGGRAPH Asia 2021 Technical Communications*, 2021, pp. 1–4.
- [30] Y.-k. Li, Y.-H. Lien, and Y.-S. Wang, "Style-structure disentangled features and normalizing flows for diverse icon colorization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 244–11 253.
- [31] T. Xiao, S. Liu, S. De Mello, Z. Yu, J. Kautz, and M.-H. Yang, "Learning contrastive representation for semantic correspondence," *International Journal of Computer Vision*, vol. 130, no. 5, pp. 1293–1309, 2022.
- [32] P. Zhang, B. Zhang, D. Chen, L. Yuan, and F. Wen, "Cross-domain correspondence learning for exemplar-based image translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5143–5153.
- [33] H. Li, B. Sheng, P. Li, R. Ali, and C. P. Chen, "Globally and locally semantic colorization via exemplar-based broad-gan," *IEEE Transactions on Image Processing*, vol. 30, pp. 8526–8539, 2021.
- [34] M. He, J. Liao, D. Chen, L. Yuan, and P. V. Sander, "Progressive color transfer with dense semantic correspondences," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 2, pp. 1–18, 2019.
- [35] B. Zhang, M. He, J. Liao, P. V. Sander, L. Yuan, A. Bermak, and D. Chen, "Deep exemplar-based video colorization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8052–8061.
- [36] P. Lu, J. Yu, X. Peng, Z. Zhao, and X. Wang, "Gray2colornet: Transfer more colors from reference image," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 3210–3218.
- [37] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [38] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in Neural Information Processing Systems*, vol. 34, pp. 8780–8794, 2021.
- [39] J. Ho, C. Saharia, W. Chan, D. J. Fleet, M. Norouzi, and T. Salimans, "Cascaded diffusion models for high fidelity image generation." *J. Mach. Learn. Res.*, vol. 23, pp. 47–1, 2022.
- [40] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.
- [41] C. Saharia, W. Chan, H. Chang, C. Lee, J. Ho, T. Salimans, D. Fleet, and M. Norouzi, "Palette: Image-to-image diffusion models," in *ACM SIGGRAPH 2022 Conference Proceedings*, 2022, pp. 1–10.
- [42] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [43] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*. Springer, 2016, pp. 694–711.
- [44] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 172–189.
- [45] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1501–1510.
- [46] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," *arXiv preprint arXiv:1802.05957*, 2018.
- [47] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International conference on machine learning*. PMLR, 2017, pp. 214–223.
- [48] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," *Advances in neural information processing systems*, vol. 30, 2017.
- [49] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, "Least squares generative adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2794–2802.
- [50] H. Winnemöller, J. E. Kyprianidis, and S. C. Olsen, "Xdog: An extended difference-of-gaussians compendium including advanced image stylization," *Computers & Graphics*, vol. 36, no. 6, pp. 740–753, 2012.
- [51] F. Bookstein, "Principal warps: thin-plate splines and the decomposition of deformations," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 11, no. 6, pp. 567–585, 1989.
- [52] K. Akita, Y. Morimoto, and R. Tsuruno, "Colorization of line drawings with empty pupils," in *Computer Graphics Forum*, vol. 39, no. 7. Wiley Online Library, 2020, pp. 601–610.