# Edge Enhanced Image Style Transfer via Transformers

Chiyu Zhang[1], Jun Yang[2], Zaiyan Dai[3], Peng Cao[4]

Sichuan Normal University

[1]alienzhang19961005@gmail.com [2]jkxy_yjun@sicnu.edu.cn

{[3]daizaiyan, [4]pc}@stu.sicnu.edu.cn

## Abstract

*In recent years, arbitrary image style transfer has attracted more and more attention. Given a pair of content and style images, a stylized one is hoped that retains the content from the former while catching style patterns from the latter. However, it is difficult to simultaneously keep well the trade-off between the content details and the style features. To stylize the image with sufficient style patterns, the content details may be damaged and sometimes the objects of images can not be distinguished clearly. For this reason, we present a new transformer-based method named STT for image style transfer and an edge loss which can enhance the content details apparently to avoid generating blurred results for excessive rendering on style features. Qualitative and quantitative experiments demonstrate that STT achieves comparable performance to state-of-the-art image style transfer methods while alleviating the content leak problem.*

## 1. Introduction

Rendering a content image into the artistic style of a referenced image is the main purpose of the image style transfer. Image style transfer is an interesting topic of computer vision with a long history. About 2 decades ago, researchers [1, 2] utilize techniques, such as texture synthesis and style transfer functions, to achieve the stylizing process. But they only focus on low-level features. Thereafter Gatys et al. [3, 4] prove that the correlation between features extracted from a pre-trained VGG can be treated as the representation of style patterns, which opens the gate of neural style transfer (NST). Iterative methods [4–10] render the content images gradually by applying gradients on the input images or the noise images while the feed-forward networks [11–19] can complete the stylizing process in one feed-forward manner after training. Vivid results though the iterative and feed-forward methods may produce, are still limited to a certain number of styles or achieve inadequate style quality. Thanks to the encoder-transfer-decoder archi-



Figure 1. Visual effects of edge loss. Compared to the results from the model without edge loss (column 3), the objects of stylized images with edge loss, especially the small ones like letters or windows, are apparently more clear and distinguishable (column 4). For the convenience of comparison, the model used above is STT. The results from other methods are also blurred in this case.

tecture, arbitrary style transfer methods [20–47] are capable of rendering images into any styles. However, these models may not work well in some cases due to the limited ability to merge the content and style features. To cope with this problem, the attention mechanism [48, 49] is introduced by a few methods [37–42] to enhance the fusion effects.

Recently, An et al. [43] discover the content leak problem that the structure of results from the CNN-based methods will be dramatically changed after a few rounds of repetitive stylization process. Deng et al. [45] and Zhang et al. [47] then prove that the transformer-based methods are capable of alleviating the problem. Different from the previous methods, IEST [42] and CAST [46] leverage the contrastive learning strategy to enhance the visual quality. However, in some cases, the structure of results from previous methods still could be blurred and the objects in images are difficult to be distinguished (see Fig. 1).

Thanks to the flexibility, scalability, and ability to capture long-range dependencies, Transformers [49–51] have been widely used in all kinds of vision tasks. Owing to the self-attention mechanism, Transformer can efficiently gather global information which is important for preserving the structure of input images. Compared to the architecture of typical CNN-based methods, Transformer evades the multi-time downsample operations which may lead to

the content leak problem [43]. Therefore, the Transformer structure has a good effect on the image style transfer task.

In this work, we propose a new Transformer-based image style transfer algorithm that is capable of producing stylized results with high visual quality while preserving fine content details. We call it STT (**S**tyle **T**ransfer on **T**ransformers). Different from StyTr2 [45] and S2WAT [47], neither does STT choose to encode content and style features in different encoders as StyTr2 did, nor does STT adopts the hierarchal structure as S2WAT did. A tiny CNN-based module is equipped as a positional encoding layer (Conv PE) which extracts the positional encoding (PE) according to the semantic information. Furthermore, to enhance the content details, a novel edge loss is applied as an extra restriction when stylized images are blurred due to the overdose of style features imposed on the inputs.

The main contributions of our work are as follows:

- A new image style transfer network name STT which can stylize images with high quality while preserving fine content details.
- A novel edge loss to enhance the content details, which improves the picture clarity of the stylized images obviously.
- Extensive experiments demonstrate that STT achieves outstanding effects and is capable of generating favorable results while preserving fine content details.

## 2. Related Work

**Image Style Transfer**. Starting from Gatys et al. [3, 4], the number of methods in NST is increasingly growing with time forward and the stylizing effects have been more and more colorful. Here we make a rough classification of these models with respect to their generalization abilities, the backbone architecture, and training strategies. In generalization abilities, the categories can be divided into single style transfer [4–15], multiple style transfer [16–18], and arbitrary style transfer [20–47]. The single style transfer encodes the fixed style features into models while the multiple style transfer utilizes certain tricks, such as conditional instance normalization [16] and StyleBank [17], to support a number of styles. With a certain module to merge the features of content and style, the arbitrary style transfer is capable of handling any style transfer. As the techniques of upstream tasks like image classification and image generation have been rapidly developing these years, an increasing number of sorceries have been introduced into image style transfer. Except for the typical CNN-based methods, the Flow-based [43] and the Transformer-based [45, 47] methods also appear in recent years. The Flow-based ArtFlow [43] is proposed to solve the problem of the content leak and the Transformer-based StyTr2 [45] and S2WAT [47] are able to alleviate the problem. The encoders of StyTr2 have the traditional Transformer structure where

the shape of representations will not be changed in processing while S2WAT adopts hierarchal architecture which means the features will be downscaled gradually. Recently, the contrastive learning strategy is introduced by IEST [42] and CAST [46]. Different from other methods trained with perceptual losses or identity losses, IEST and CAST treat the contrastive loss and the adversarial loss as the optimization targets to achieve satisfying effects. However, in some cases, the results from the existing image style transfer methods may still be blurred due to no restriction on the edge of the main objects in inputs.

**Vison Transformer**. Inherited the ability to capture the long-range dependencies from Transformers in natural language processing (NLP), vision Transformers have been developed in a wide variety of vision tasks, including image classification [52–64], object detection [65–69], semantic segmentation [70, 71], and image generation [72, 73]. In image style transfer, StyTr2 and S2WAT have demonstrated that both the traditional structure and the hierarchical architecture have a favorable effect on style transfer. In this paper, we leverage several convolutional operations to fulfill the positional encoding instead of parametric positional encoding [52] or the one with pooling operations [45].

**Utilization of Edge Maps in Style Transfer**. The operators like Laplacian, Canny, and Sobel are widely used in edge and contour detection. In image style transfer, Li et al. discover the correspondences between Laplacian deviations and image distortions and then propose Lapstyle [8], an iterative image style transfer method based on a Laplacian loss. Subsequently, Li et al. apply the Laplacian filter on the drafting and revision network and then present LapStyle [36], a feed-forward image style transfer method based on the Laplacian filter. However, the above methods are all applied to the CNN-based models. In this work, we leverage the edge maps to enhance the results from the Transformer-based STT which is capable of producing stylized images with fine content details and colorful artistic features.

## 3. Method

As shown in Fig. 2, the proposed STT has the architecture of encoder-transfer-decoder. The positional encoding (PE) is first extracted from both the content images $I_c$ and the style image $I_s$ by a module name Conv PE. Then after splitting the content images $I_c$ and style image $I_s$ into non-overlapping patches, a linear projection is equipped to transform the patches into sequences. The sum of the sequences and the PE will be treated as the inputs of the Transformer encoder. Generated from the encoder, the content features $f_c$ and style features $f_s$ then will be merged in the transfer module which is based on a Transformer decoder. Finally,

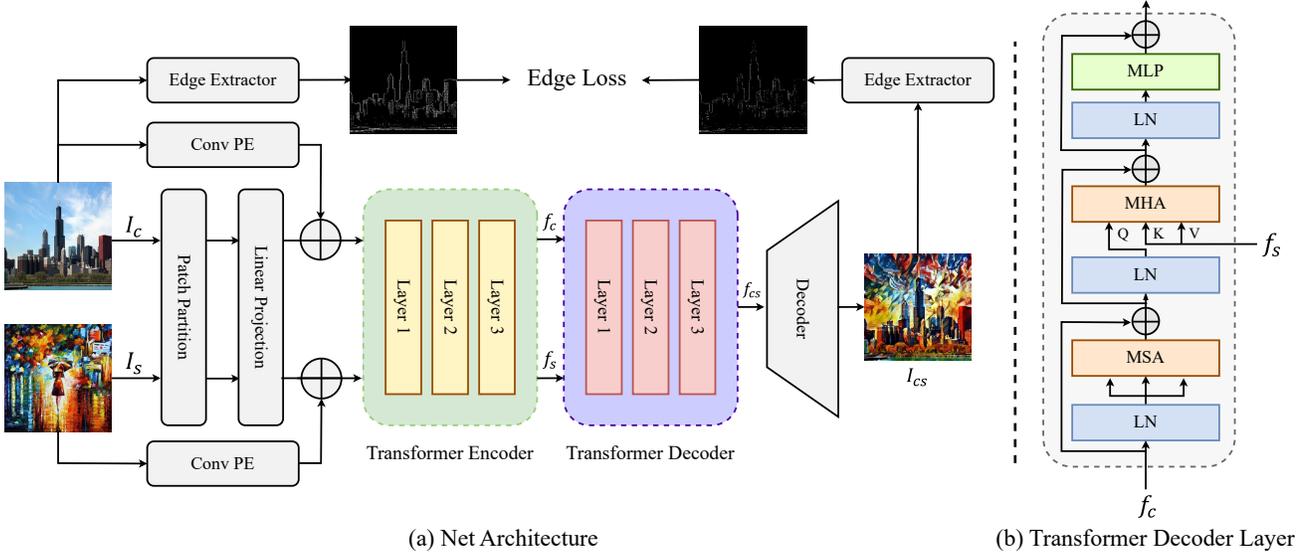(a) Net Architecture

(b) Transformer Decoder Layer

Figure 2. The net architecture of the proposed STT.

the outputs can be obtained by decoding the fused features $f_{cs}$ in a CNN-based decoder. In addition, to calculate the edge loss during the training step, the edge maps of the content images and the stylized images need to be extracted by a fine-designed edge extractor which will be introduced later.

In this part, we plan to present the overall architecture of the proposed STT first in Section 3.1 and in Section 3.2 introduce the edge extractor which is used to calculate the edge loss. Finally, the optimization strategy will be discussed in Section 3.3.

### 3.1. Overall Architecture

**Encoder**. Before the process of the encoder, a tiny CNN-based module named Conv PE is applied to the content images and style images to extract the content-aware positional encoding. As depicted in Fig. 3, Conv PE is composed by three convolutional layers, two reflections (padding) layers, and one ReLU activation layer. The main role of the reflection layers is to ensure that the size of results is consistent before and after processing while the convolutional layers are used to extract the content-aware positional encoding. As shown in Fig. 6, we find that the results from the model with Conv PE are obviously better than that of the one without.

Different from the design of StyTr2 [45] which has two independent domain-specific encoders for the content images and style images, STT treats them as normal pictures with content & style features and encodes them in one Transformer-based encoder. Given the input image in the shape of $H \times W \times 3$, the input will first be split into patches by the patch partition layer and then embedded into

sequences linearly with the shape of $\frac{HW}{8 \times 8} \times C$ (768 is the default value of $C$). Adding the positional encoding from Conv PE, the resulting sequences will be fed to a three-layer Transformer encoder. The computation process of each layer can be defined as:

$$\hat{c}^l = MSA(LN(c^{l-1})) + c^{l-1} \qquad (1)$$
$$c^l = MLP(LN(\hat{c}^l)) + \hat{c}^l \qquad (2)$$

where $\hat{c}^l$ and $c^l$ denote the outputs of MSA and MLP for layer $l$ respectively; MSA represents the module of multi-head self-attention while MLP denotes the module of multi-layer perceptron; and $LN$ means LayerNorm. After the three layers of the encoder, we obtain the content features $f_c$ and style features $f_s$ with the shape consistent before and after processing.
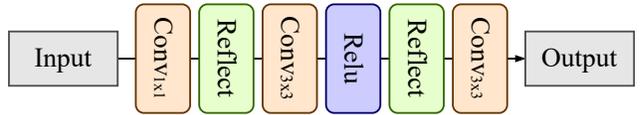


Figure 3. The illustration of Conv PE.

**Decoder**. Instead of upsampling the stylized features $f_{cs}$ to the original size in a single projection as [74] once did, STT follows [21, 37, 41] to utilize a mirrored VGG to decode the stylized features gradually. Before the step of decoding, the stylized features need to be reshaped first for the sequence-like shape $\frac{HW}{8 \times 8} \times C$. After the three stages of upsampling and refining, we obtain the stylized images $I_{cs}$ with the shape of $H \times W \times 3$.

**Transfer Module**. The transfer module is used to merge the content features $f_c$ and style features $f_s$. We introduce the transfer module of S2WAT [47] as the means to fuse the features. As depicted in Fig. 2 (b), the transfer module consists of three layers of the Transformer decoder layers, and each layer is mainly composed by an MSA module, an MHA module, and an MLP module. The computational process can be defined as:

$$\hat{x}^l = MSA(LN(x^{l-1})) + x^{l-1}$$
$$Q = LN(\hat{x}^l) \cdot W_Q$$
$$K, V = y \cdot W_K, \ y \cdot W_V$$
$$\tilde{x}^l = MHA(Q, K, V) + \hat{x}^l$$
$$x^l = MLP(LN(\tilde{x}^l)) + \tilde{x}^l \tag{3}$$

where $\hat{x}^l$, $\tilde{x}^l$, and $x^l$ represent the results of MSA, MHA, and MLP for layer $l$, respectively; $y$ denotes the style features; $W_Q$, $W_K$, and $W_V$ are the projection matrices for $Q$, $K$, and $V$; $Q$, $K$, and $V$ denote the query, key, and value vectors. Leveraging the fusion effects of the cross attention, the stylized features $f_{cs}$ can be received.

### 3.2. Edge Extractor

To make the content structure of the stylized images clear, we design a novel edge loss to enhance the edge of the objects in output images. Before calculating the edge loss, the edge maps suitable for style transfer need to be captured by a fine-designed edge extractor. Different from the tasks like edge detection or contour extraction, the content details of the outputs in image style transfer are probably not the same as that of the content images, especially the background which may has the artistic patterns from the style images. The results will be blurred if we take the similarity between the edge maps from the content images and the stylized images as the optimization target directly. Therefore one of the problems that need to be solved is to filter out the place where the main structure of the content images does not exist. A mask operation is introduced to cope with this problem. As shown in (6), all of the edges in the edge maps of the stylized images ($edg'$-$I_{cs}$) that are not exist in the corresponding place of the edge maps of content images ($edg$-$I_c$) will be masked out. Furthermore, we also set a threshold to exclude the weak responses of edge maps which may play a role as noise. The overall computational process can be defined as:

$$edg\text{-}I_c = threshold(lap(I_c), \tau) \tag{4}$$
$$edg'\text{-}I_{cs} = threshold(lap(I_{cs}), \tau) \tag{5}$$
$$edg\text{-}I_{cs} = mask(edg'\text{-}I_{cs}, edg\text{-}I_c) \tag{6}$$

where $edg$-$I_c$ and $edg$-$I_{cs}$ are the edge maps of the content and stylized images respectively; $lap$ denotes the Laplacian

operator and $threshold$ represents the function that sets 0 to the responses where the value is smaller than the threshold parameter $\tau$ (0.2 is set as the default value). After the above steps, we obtain the refined edge maps to be used in calculating the edge loss.

### 3.3. Network Optimization

The main purpose of image style transfer is to maintain the structure of the content images while transferring the artistic patterns to the stylized results from the style images. To achieve this target, we follow [21] to construct two perceptual losses to measure the content differences between the stylized images and the content images as well as the style differences between the stylized images and the style images. Furthermore, we also adopt the identity losses [37] to enrich the content details and style patterns of the stylized images. Finally, the proposed edge loss is equipped to enhance the content structure further. As shown in (7), the whole loss function can be defined as:

$$\mathcal{L}_{total} = \lambda_c \mathcal{L}_{content} + \lambda_s \mathcal{L}_{style} +$$
$$\lambda_{id1} \mathcal{L}_{id1} + \lambda_{id2} \mathcal{L}_{id2} + \lambda_{edg} \mathcal{L}_{edg} \tag{7}$$

where the $\lambda_c$, $\lambda_s$, $\lambda_{id1}$, $\lambda_{id2}$, and $\lambda_{edg}$ are the weights of losses; $\mathcal{L}_{content}$ and $\mathcal{L}_{style}$ denote the perceptual losses; $\mathcal{L}_{id1}$ and $\mathcal{L}_{id2}$ are the identity losses; $\mathcal{L}_{edg}$ represents the edge loss and we only apply the edge loss in the situation when the results are apparently blurred. We set $\lambda_c$, $\lambda_s$, $\lambda_{id1}$, $\lambda_{id2}$, and $\lambda_{edg}$ to 1, 3, 50, 1, and 5000 to alleviate the impact of magnitude differences.

**Perceptual Loss**. Similar to [21], we leverage a pretrained VGG19 to extract the feature maps of the content and style images which are used to calculate the perceptual losses. In our model, the layer $Relu\_4\_1$ and $Relu\_5\_1$ are used to calculate the content perceptual loss while the layer $Relu\_1\_1$, $Relu\_2\_1$, $Relu\_3\_1$, $Relu\_4\_1$, and $Relu\_5\_1$ are used to calculate the style perceptual loss. One thing that needs to be attended to is that the mean-variance channel-wise normalization is applied on the feature maps before the calculation of the content perceptual loss. The perceptual losses can be defined as:

$$\mathcal{L}_{content} = \sum_{l \in C} \|\overline{\phi_l(I_{cs})} - \overline{\phi_l(I_c)}\|_2 \tag{8}$$

$$\mathcal{L}_{style} = \sum_{l \in L} \|\mu(\phi_l(I_{cs})) - \mu(\phi_l(I_s))\|_2 +$$
$$\|\sigma(\phi_l(I_{cs})) - \sigma(\phi_l(I_s))\|_2 \tag{9}$$

where the $C$ and $L$ are the layers of the pretrained VGG which are concerned to calculate the content and style perceptual losses respectively; $\phi_l$ denotes the feature maps of the $l$-th layer in the pretrained VGG; $\mu$ and $\sigma$ are the mean and variance of the features; and the overline represents the mean-variance channel-wise normalization.

**Identity Loss**. Following the work of [37], a pair of identity losses are constructed to learn the relationship between the content and style representations. The identity losses are defined as :

$$\mathcal{L}_{id1} = \|I_{cc} - I_c\|_2 + \|I_{ss} - I_s\|_2 \qquad (10)$$

$$\mathcal{L}_{id2} = \sum_{l \in L} \|\phi_l(I_{cc}) - \phi_l(I_c)\|_2 + \|\phi_l(I_{ss}) - \phi_l(I_s)\|_2 \qquad (11)$$

where $I_{cc}$ ($I_{ss}$) denotes the stylized images from a common pair of content (style) images. Specifically, the original content (style) image is expected when we feed two of the same content (style) images to the model. As shown in (11), this operation is also applied on feature maps from the pretrained VGG. And the layer $Relu\_1\_1$, $Relu\_2\_1$, $Relu\_3\_1$, $Relu\_4\_1$, and $Relu\_5\_1$ are used to calculate the second identity loss.

**Edge Loss**. To enhance the edge of the objects when the original results from STT are obviously blurred, we design an edge loss to cope with this problem. As depicted in Section 3.2, the edge maps are computed by the Laplacian operator first and then refined by a threshold function and a mask operation successively. After we obtain the refined edge maps, the edge loss can be computed in the following process:

$$\mathcal{L}_{edg} = \|edg\text{-}I_c - edg\text{-}I_{cs}\|_2 \qquad (12)$$

where $edg\text{-}I_c$ and $edg\text{-}I_{cs}$ are the refined edge maps of the content and stylized images respectively. As shown in Fig. 1 and Fig. 7 (columns 3 and 6), applying the edge loss on STT can obviously improve the edges of blurred results.

# 4. Experiments

## 4.1. Implementation Details

**Datasets**. MS-COCO [75] is used as the content dataset while WikiArt [76] is used as the style dataset. We randomly select 80000 images of each dataset to build the training datasets. During the process of training, the input image will be resized to 512 on the shorter side first and then randomly cropped into $224 \times 224$. While in the process of testing, inputs of any size are accepted.

**Training Information**. Pytorch framework is used to implement STT and 40000 iterations are taken to complete the training. With a batch size of 4 and an initial learning rate of 1e-4, we use an Adam optimizer [77] to train the network and the warmup strategy [78] to adjust the learning rate. The training step is taken about 10 hours on a single Tesla V100 GPU. We also calculate the reference time (see the last row of Table 1) of different image style transfer models with one Tesla P100 GPU.

## 4.2. Style Transfer Results

In order to demonstrate the style transfer effect of the proposed STT, we make a comparison between the results from the proposed STT and the state-of-the-art arbitrary style transfer methods, including AdaIn [21], WCT [22], SANet [37], MCC [41], ArtFlow [43], IEST [42], CAST [46], StyTr2 [45], and S2WAT [47] .

**Qualitative Comparison**. The results of the qualitative comparison are presented in Fig. 4. Although the different methods fulfill the image style transfer in different ways, they all achieve colorful results. Due to the over-simplified alignment of the second-order statistics, AdaIN can not draw sufficient style patterns on the content images. By applying the alignment process on the style feature space with whitening and coloring operations, WCT attracts more artistic characteristics but damages the content details. Inspired by the attention mechanism, SANet transfers adequate style features to the content images but the structure is not ideal sometimes. MCC suffers from an overflow issue for the lack of linear operations. In conjunction with the projection flow network, ArtFlow is capable of producing content-unbiased results but sometimes may generate undesired patterns on the borders. Different from other methods which train the models with perceptual losses or identity losses, IEST and CAST adopt the contrastive learning strategy and make favorable effects sometimes. But in some cases, the results fail to obtain plentiful style representations. Transformer-based methods find a better balance between content and style. With the Transformer-based encoder and transfer module, StyTr2 and S2WAT both achieve satisfying effects while S2WAT may lose some style patterns and StyTr2 drops content details in some places. As shown in the last column of Fig. 4, STT preserves the fine content details while sufficient artistic characteristics are transferred.

**Quantitative Comparison**. In this part, the content differences between the stylized images and the content images are computed as an indirect metric to measure the content quality while the style differences between the results and the style images are calculated as an implicit metric to evaluate the style quality. The identity losses are also taken into consideration playing a role as the auxiliary metrics to judge the ability to preserve content/style features. As shown in Table 1, S2WAT achieves the lowest content loss while STT and SANet outperform the other methods on style quality. Compared with the CNN-based models, the Transformer-based methods have obvious advantages in identity losses. Due to the ability of completely reversible transformation, ArtFlow does not use identity losses. Although ArtFlow can produce content-unbiased results, STT outperforms it on style quality. In summary, STT can preserve both the
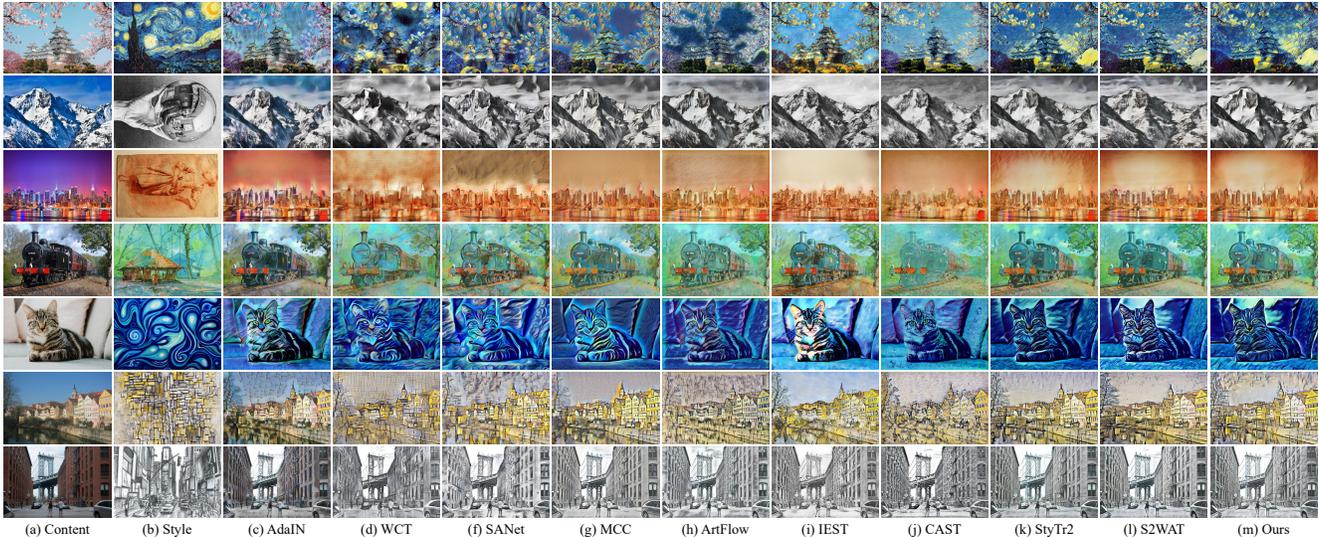
Figure 4. The visual comparison of the state-of-the-art arbitrary style transfer algorithms.

content details from the content image as well as the style patterns from the style images.

## 4.3. Content Leak

After repeated stylization with the same pair of content and style images, CNN-based methods will suffer the problem of content leak that the content structure will drop gradually as the number of experimental rounds grows. An et al. [43] utilize the projection flow network, a kind of network which is able to achieve completely reversible transformation, to settle the content leak problem. However, strict reversibility may have an undesired impact on the stylizing process. With the ability to capture long-range dependencies, StyTr2 [45] and S2WAT [47] are demonstrated to be capable of alleviating the content leak problem.

To examine the stylizing effects on the content leak issue, we make a comparison with the CNN-based method [21, 22, 37, 41, 42, 46], the Flow-based method [43], and the Transformer-based methods [45, 47]. As depicted in Fig. 5, the results from the 1st and the 20th rounds of repeated stylization have been presented. All the methods can keep the content details well after the 1st stylizing process except that the results from AdaIN and ArtFlow are to some degree lack of style features. However, after the 20th round of the stylizing process, the CNN-based methods fail to preserve the content structure and the results are apparently blurred. Compared to the completely content-unbiased ArtFlow, the Transformer-based StyTr2, S2WAT, and the proposed STT still drop the content details slightly but the results are obviously superior to that of the CNN-based methods. Therefore, the proposed STT can preserve both the content structure and the style features while capable of alleviating the content leak problem.

## 4.4. Ablation Study

**Conv PE**. Positional encodings (PE) are important for Transformer-based models, which provide information on locations. There are two types of absolute positional encoding (APE) that are widely used: functional [49] and parametric [51] positional encoding. As Deng et al. [45] have discussed in StyTr2, the functional APE, such as the sinusoidal APE, will result in vertical track artifacts due to the large positional deviation. And we examine the parametric APE whose results are shown in Fig. 6 (column 4). Some undesired patterns that do not vary substantially with the inputs appear on the outputs. Due to the unsatisfactory performance of the functional APE and parametric APE, we propose a positional encoding based on convolutional operations (Conv PE), and the results are presented in Fig. 6 (column 5). Because the CAPE needs to work with the transfer module while the transfer module of STT does not have the interface of PE, we do not conduct the experiments on CAPE.

As depicted in Fig. 6, the strokes of the results from the model without PE are obviously thicker than that from the model with Conv PE. Furthermore, there are a few vertical track artifacts on the edge of objects in images (row 1 column 3). For the results from the model with parametric APE, the background is blurry and a sort of undesired pattern makes the pictures unsightly. By contrast, the results from STT fix these problems and preserve both the content details and style features.

**Edge Loss**. When the results of image style transfer are blurred, applying the edge loss on STT can improve picture clarity obviously. As depicted in Fig. 1, the model without the edge loss erases the majority of content details in the

| Method | Ours | S2WAT | StyTr2 | CAST | IEST | ArtFlow | MCC | SANet | WCT | AdaIN |
|---|---|---|---|---|---|---|---|---|---|---|
| *Content Loss* ↓ | 2.18 | **1.66** | 1.83 | 2.07 | 1.81 | 1.93 | 1.92 | 2.16 | 2.56 | <u>1.71</u> |
| *Style Loss* ↓ | <u>1.35</u> | 1.74 | 1.52 | 4.33 | 2.72 | 1.90 | 1.70 | **1.11** | 2.23 | 3.50 |
| *Identity Loss 1* ↓ | **0.16** | **0.16** | <u>0.26</u> | 1.94 | 0.91 | 0.00 | 1.07 | 0.81 | 3.01 | 2.54 |
| *Identity Loss 2* ↓ | <u>1.55</u> | **1.38** | 3.10 | 18.72 | 7.16 | 0.00 | 7.72 | 6.03 | 21.88 | 17.97 |
| *Time(seconds)* ↓ | 0.270 | 0.558 | 0.237 | **0.042** | <u>0.061</u> | 0.325 | 0.078 | <u>0.061</u> | 0.590 | **0.042** |

Table 1. Quantitative comparison between the results from different image style transfer methods. The loss values above are all computed on 400 random samples average and the reference time is calculated on a hundred random samples in a resolution of $512 \times 512$. The bold font marks the best values while the underline shows the second-best values.
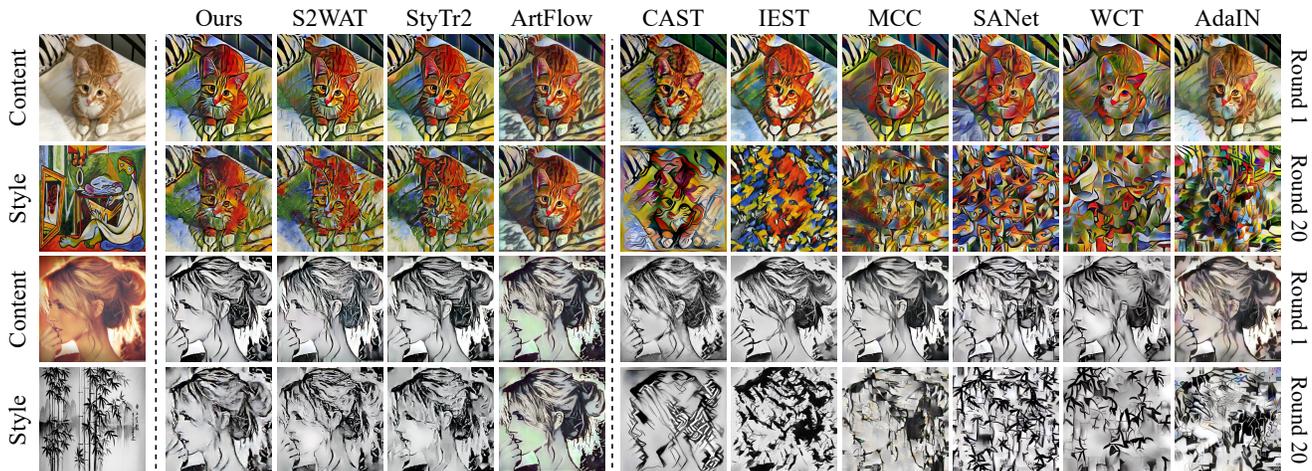


Figure 5. Visualization of the content leak problem.



Figure 6. Comparison between the results from different types of PE.

content images, such as the windows on buildings (row 1 column 3) and the letters on the billboards (row 2 column 3). In contrast, these details are well preserved when the edge loss is equipped (column 4).

Besides the comparison between the models with and without the edge loss, we also compare the operators to extract the edge maps which are the important step to form the edge loss. As depicted in Fig. 7 (a), the operator of Canny, Sobel, and Laplacian are taken into consideration. A kind of hollow stroke appears on the results based on the Canny operator (see column 4) while the results on the Laplacian operator can produce natural and fine strokes. The clearest result though the model based on the Sobel operator can generate, unpleasant patterns, such as the vertical/horizontal tracks and the blurred strokes, appears in the stylized im-

ages (see column 5). For the outputs based on the Laplacian operator which is applied to the edge loss, the strokes are natural and the structure of objects is clear which demonstrate the performance of the edge loss.

In addition, we also provide the edge maps calculated by the edge extractor where a phenomenon can be easily found that the edges of the results from the model applying the edge loss will be much richer than that of the results from models without the edge loss.

## 5. Conclusion

In this work, we proposed a Transformer-based method named STT for arbitrary image style transfer. The proposed STT has a Transformer-based encoder that can encode both the content and style images capturing the long-range information between them. A content-aware positional encoding scheme (Conv PE) based on the convolutional operations is applied to the encoder to provide the positional information. To overcome the problem that the results of image style transfer are blurred in some cases, a novel edge loss is presented to improve the clarity of the stylized images. As another new method based on Transformer, STT is capable of producing vivid stylized images with fine content details and sufficient style features while alleviating the content leak problem.
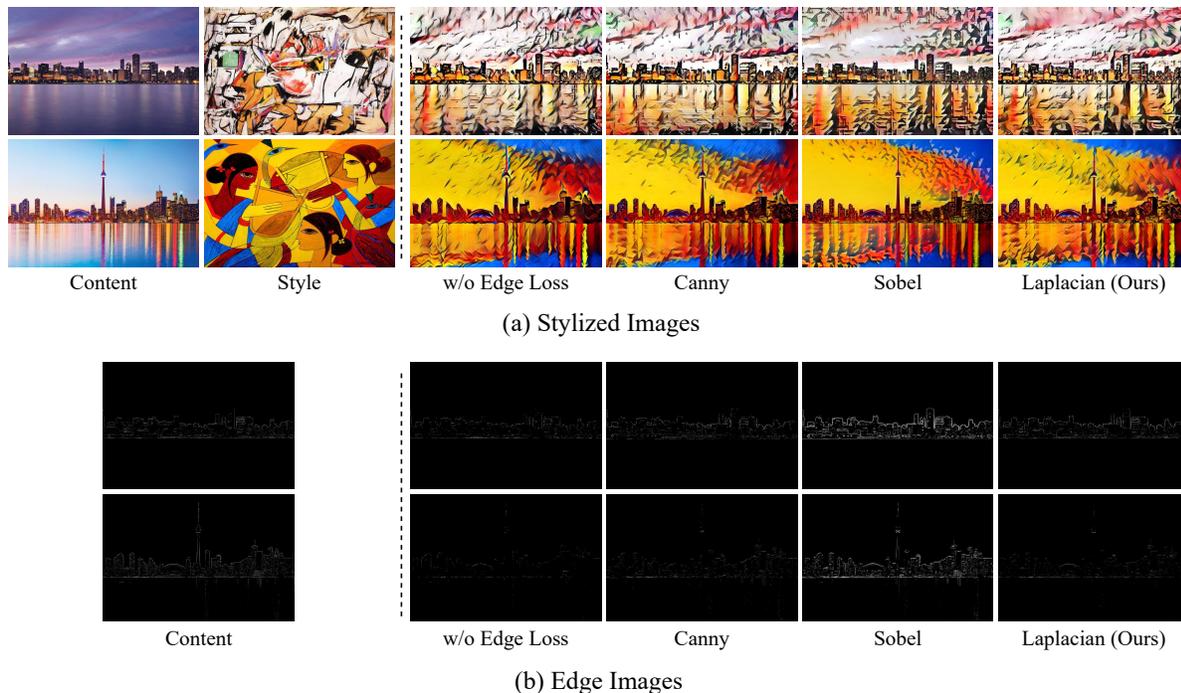
(a) Stylized Images



(b) Edge Images

Figure 7. Comparison between the results using different edge detection operators.

# References

[1] Alexei A Efros and William T Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 341–346, 2001. 1

[2] Stefan Bruckner and M Eduard Gröller. Style transfer functions for illustrative volume rendering. In *Computer Graphics Forum*, volume 26, pages 715–724. Wiley Online Library, 2007. 1

[3] Leon Gatys, Alexander S Ecker, and Matthias Bethge. Texture synthesis using convolutional neural networks. *Advances in neural information processing systems*, 28, 2015. 1, 2

[4] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015. 1, 2

[5] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016. 1, 2

[6] Leon A Gatys, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. Preserving color in neural artistic style transfer. *arXiv preprint arXiv:1606.05897*, 2016. 1, 2

[7] Eric Risser, Pierre Wilmot, and Connelly Barnes. Stable and controllable neural texture synthesis and style transfer using histogram losses. *arXiv preprint arXiv:1701.08893*, 2017. 1, 2

[8] Shaohua Li, Xinxing Xu, Liqiang Nie, and Tat-Seng Chua. Laplacian-steered neural style transfer. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1716–1724, 2017. 1, 2

[9] Yanghao Li, Naiyan Wang, Jiaying Liu, and Xiaodi Hou. Demystifying neural style transfer. *arXiv preprint arXiv:1701.01036*, 2017. 1, 2

[10] Jiahao Lu. Transformer-based neural texture synthesis and style transfer. In *2022 4th Asia Pacific Information Technology Conference*, pages 88–95, 2022. 1, 2

[11] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 1, 2

[12] Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, and Victor Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. *arXiv preprint arXiv:1603.03417*, 2016. 1, 2

[13] Chuan Li and Michael Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *European conference on computer vision*, pages 702–716. Springer, 2016. 1, 2

[14] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 1, 2

[15] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6924–6932, 2017. 1, 2

[16] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. *arXiv preprint arXiv:1610.07629*, 2016. 1, 2

[17] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1897–1906, 2017. 1, 2

[18] Minxuan Lin, Fan Tang, Weiming Dong, Xiao Li, Changsheng Xu, and Chongyang Ma. Distribution aligned multimodal and multi-domain image stylization. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 17(3):1–17, 2021. 1, 2

[19] Hang Zhang and Kristin Dana. Multi-style generative network for real-time transfer. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 1

[20] Tian Qi Chen and Mark Schmidt. Fast patch-based style transfer of arbitrary style. *arXiv preprint arXiv:1612.04337*, 2016. 1, 2

[21] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 1, 2, 3, 4, 5, 6

[22] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. *Advances in neural information processing systems*, 30, 2017. 1, 2, 5, 6

[23] Lu Sheng, Ziyi Lin, Jing Shao, and Xiaogang Wang. Avatar-net: Multi-scale zero-shot style transfer by feature decoration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8242–8250, 2018. 1, 2

[24] Shuyang Gu, Congliang Chen, Jing Liao, and Lu Yuan. Arbitrary style transfer with deep feature reshuffle. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8222–8231, 2018. 1, 2

[25] Yijun Li, Ming-Yu Liu, Xueting Li, Ming-Hsuan Yang, and Jan Kautz. A closed-form solution to photorealistic image stylization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 453–468, 2018. 1, 2

[26] Ming Lu, Hao Zhao, Anbang Yao, Yurong Chen, Feng Xu, and Li Zhang. A closed-form solution to universal style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5952–5961, 2019. 1, 2

[27] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast image and video style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3809–3817, 2019. 1, 2

[28] Huan Wang, Yijun Li, Yuehai Wang, Haoji Hu, and Ming-Hsuan Yang. Collaborative distillation for ultra-resolution universal style transfer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1860–1869, 2020. 1, 2

[29] Jie An, Haoyi Xiong, Jun Huan, and Jiebo Luo. Ultrafast photorealistic style transfer via neural architecture search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10443–10450, 2020. 1, 2

[30] Zhijie Wu, Chunjin Song, Yang Zhou, Minglun Gong, and Hui Huang. Efanet: Exchangeable feature alignment network for arbitrary style transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12305–12312, 2020. 1, 2

[31] Jan Svoboda, Asha Anoosheh, Christian Osendorfer, and Jonathan Masci. Two-stage peer-regularized feature recombination for arbitrary image style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13816–13825, 2020. 1, 2

[32] Xiao-Chang Liu, Xuan-Yi Li, Ming-Ming Cheng, and Peter Hall. Geometric style transfer. *arXiv preprint arXiv:2007.05471*, 2020. 1, 2

[33] Yongcheng Jing, Xiao Liu, Yukang Ding, Xinchao Wang, Errui Ding, Mingli Song, and Shilei Wen. Dynamic instance normalization for arbitrary style transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4369–4376, 2020. 1, 2

[34] Zhizhong Wang, Lei Zhao, Haibo Chen, Lihong Qiu, Qihang Mo, Sihuan Lin, Wei Xing, and Dongming Lu. Diversified arbitrary style transfer via deep feature perturbation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7789–7798, 2020. 1, 2

[35] Jie An, Tao Li, Haozhi Huang, Li Shen, Xuan Wang, Yongyi Tang, Jinwen Ma, Wei Liu, and Jiebo Luo. Real-time universal style transfer on high-resolution images via zero-channel pruning. *arXiv preprint arXiv:2006.09029*, 2020. 1, 2

[36] Tianwei Lin, Zhuoqi Ma, Fu Li, Dongliang He, Xin Li, Errui Ding, Nannan Wang, Jie Li, and Xinbo Gao. Drafting and revision: Laplacian pyramid network for fast high-quality artistic style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5141–5150, 2021. 1, 2

[37] Dae Young Park and Kwang Hee Lee. Arbitrary style transfer with style-attentional networks. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5880–5888, 2019. 1, 2, 3, 4, 5, 6

[38] Yuan Yao, Jianqiang Ren, Xuansong Xie, Weidong Liu, Yong-Jin Liu, and Jun Wang. Attention-aware multi-stroke style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1467–1475, 2019. 1, 2

[39] Yingying Deng, Fan Tang, Weiming Dong, Wen Sun, Feiyue Huang, and Changsheng Xu. Arbitrary style transfer via multi-adaptation network. In *Proceedings of the 28th ACM international conference on multimedia*, pages 2719–2727, 2020. 1, 2

[40] Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Meiling Wang, Xin Li, Zhengxing Sun, Qian Li, and Errui Ding. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6649–6658, 2021. 1, 2

[41] Yingying Deng, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, and Changsheng Xu. Arbitrary video style transfer via multi-channel correlation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1210–1217, 2021. 1, 2, 3, 5, 6

9

[42] Haibo Chen, Zhizhong Wang, Huiming Zhang, Zhiwen Zuo, Ailin Li, Wei Xing, Dongming Lu, et al. Artistic style transfer with internal-external learning and contrastive learning. *Advances in Neural Information Processing Systems*, 34:26561–26573, 2021. 1, 2, 5, 6

[43] Jie An, Siyu Huang, Yibing Song, Dejing Dou, Wei Liu, and Jiebo Luo. Artflow: Unbiased image style transfer via reversible neural flows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 862–871, 2021. 1, 2, 5, 6

[44] Xiaolei Wu, Zhihao Hu, Lu Sheng, and Dong Xu. Styleformer: Real-time arbitrary style transfer via parametric style composition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14618–14627, 2021. 1, 2

[45] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11326–11336, 2022. 1, 2, 3, 5, 6

[46] Yuxin Zhang, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, Tong-Yee Lee, and Changsheng Xu. Domain enhanced arbitrary image style transfer via contrastive learning. *arXiv preprint arXiv:2205.09542*, 2022. 1, 2, 5, 6

[47] Chiyu Zhang, Jun Yang, Lei Wang, and Zaiyan Dai. S2wat: Image style transfer via hierarchical vision transformer using strips window attention. *arXiv preprint arXiv:2210.12381*, 2022. 1, 2, 4, 5, 6

[48] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014. 1

[49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 1, 6

[50] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018. 1

[51] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1, 6

[52] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2

[53] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. 2

[54] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 558–567, 2021. 2

[55] Xiangxiang Chu, Bo Zhang, Zhi Tian, Xiaolin Wei, and Huaxia Xia. Do we really need explicit position encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 3(8), 2021. 2

[56] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *Advances in Neural Information Processing Systems*, 34:15908–15919, 2021. 2

[57] Yifan Xu, Zhijie Zhang, Mengdan Zhang, Kekai Sheng, Ke Li, Weiming Dong, Liqing Zhang, Changsheng Xu, and Xing Sun. Evo-vit: Slow-fast token evolution for dynamic vision transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2964–2972, 2022. 2

[58] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. *arXiv preprint arXiv:2204.01697*, 2022. 2

[59] Benjamin Graham, Alaaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou, and Matthijs Douze. Levit: a vision transformer in convnet's clothing for faster inference. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12259–12269, 2021. 2

[60] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 2

[61] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6824–6835, 2021. 2

[62] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Improved multiscale vision transformers for classification and detection. *arXiv preprint arXiv:2112.01526*, 2021. 2

[63] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 2

[64] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12009–12019, 2022. 2

[65] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 2

[66] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable trans-

formers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 2

[67] Zhigang Dai, Bolun Cai, Yugeng Lin, and Junying Chen. Up-detr: Unsupervised pre-training for object detection with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1601–1610, 2021. 2

[68] Josh Beal, Eric Kim, Eric Tzeng, Dong Huk Park, Andrew Zhai, and Dmitry Kislyuk. Toward transformer-based object detection. *arXiv preprint arXiv:2012.09958*, 2020. 2

[69] Yanghao Li, Saining Xie, Xinlei Chen, Piotr Dollar, Kaiming He, and Ross Girshick. Benchmarking detection transfer learning with vision transformers. *arXiv preprint arXiv:2111.11429*, 2021. 2

[70] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8741–8750, 2021. 2

[71] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. 2

[72] Yifan Jiang, Shiyu Chang, and Zhangyang Wang. Transgan: Two transformers can make one strong gan. *arXiv preprint arXiv:2102.07074*, 1(3), 2021. 2

[73] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022. 2

[74] Kwonjoon Lee, Huiwen Chang, Lu Jiang, Han Zhang, Zhuowen Tu, and Ce Liu. Vitgan: Training gans with vision transformers. *arXiv preprint arXiv:2107.04589*, 2021. 3

[75] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 5

[76] Fred Phillips and Brandy Mackintosh. Wiki art gallery, inc.: A case for critical thinking. *Issues in Accounting Education*, 26(3):593–608, 2011. 5

[77] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[78] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *International Conference on Machine Learning*, pages 10524–10533. PMLR, 2020. 5