

# Fast Universal Style Transfer for Artistic and Photorealistic Rendering

Jie An<sup>\*†</sup>

School of Mathematical Sciences  
Peking University

jie.an@pku.edu.cn

Haoyi Xiong<sup>\*</sup>

Big Data Lab  
Baidu Research

xionghaoyi@baidu.com

Jiebo Luo

Department of Computer Science  
University of Rochester

jluo@cs.rochester.edu

Jun Huan<sup>‡</sup>

Big Data Lab  
Baidu Research

huanjun@baidu.com

Jinwen Ma<sup>‡</sup>

School of Mathematical Sciences  
Peking University

jwma@math.pku.edu.cn

## Abstract

*Universal style transfer is an image editing task that renders an input content image using the visual style of arbitrary reference images, including both artistic and photorealistic stylization. Given a pair of images as the source of content and the reference of style, existing solutions usually first train an auto-encoder (AE) to reconstruct the image using deep features and then embeds pre-defined style transfer modules into the AE reconstruction procedure to transfer the style of the reconstructed image through modifying the deep features. While existing methods typically need multiple rounds of time-consuming AE reconstruction for better stylization, our work intends to design novel neural network architectures on top of AE for fast style transfer with fewer artifacts and distortions all in one pass of end-to-end inference. To this end, we propose two network architectures named **ArtNet** and **PhotoNet** to improve artistic and photo-realistic stylization, respectively. Extensive experiments demonstrate that ArtNet generates images with fewer artifacts and distortions against the state-of-the-art artistic transfer algorithms, while PhotoNet improves the photorealistic stylization results by creating sharp images faithfully preserving rich details of the input content. Moreover, ArtNet and PhotoNet can achieve  $3\times$  to  $100\times$  speed-up over the state-of-the-art algorithms, which is a major advantage for large content images.*

## 1. Introduction

Universal style transfer is an image editing task that renders an input content image using the visual styles of arbitrary reference images. Among a wide range of stylization tasks, two common tasks of style transfer are *artistic style transfer* and *photorealistic stylization*. More specifically, given a pair of images for content and style reference, respectively, *artistic style transfer* aims to generate an artistic rendition of the given content using the textures, colors, and shapes of the style reference, while *photorealistic stylization* creates “photographs” of the given content as if they were captured in the same settings of the style references.

Many research efforts [6, 7, 15, 36, 20, 12] have been dedicated to universal style transfer. A pioneer work is presented in [6, 7] where Gatys. *et al.* make the first attempt to connect the style representation to the Gram matrices of deep features. Their work shows that Gram matrices of deep features have an exceptional ability to encode the styles of images. Following this line of research, several recent works aiming to minimize a Gram matrices based loss function [6, 7, 27, 25] have been proposed. While these algorithms can produce high-quality stylization results, they all suffer from a high computational cost even with acceleration techniques [15, 36, 38, 4, 19]. Moreover, all these algorithms usually can work well only on a limited number of image styles.

Recently, universal style transfer methods [2, 20, 12, 21, 30, 8] have been proposed for image transfer with respect to *arbitrary* styles and contents. For example, multi-level stylization [20, 21] or iterative EM process [8] have been proposed recently. Multi-level stylization algorithms first extract multi-level features using a multi-layer auto-encoder (AE) in depth, where each layer of AE refers to a level of features, then render styles using all the features

<sup>\*</sup>Equal contribution.

<sup>†</sup>This work was done when Jie An worked as an Intern at Baidu Inc..

<sup>‡</sup>Corresponding author.

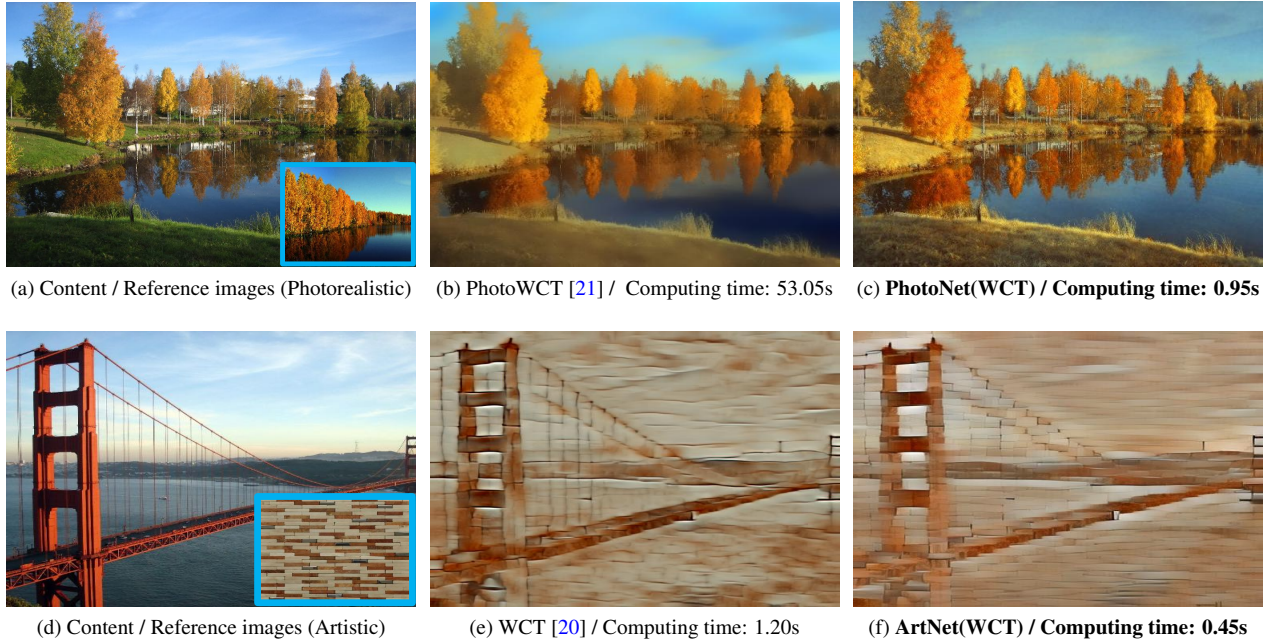


Figure 1: **Visual comparison of photorealistic and artistic style transfer.** Content images are (a) and (d); Reference style images are shown in the bottom-right corners of (a) and (d). For photorealistic stylization, PhotoWCT [21] consumes significant computing time while producing an overly smooth image shown in (b). Our proposed PhotoNet generates the image shown in (c) of rich details with only 1/50th of the computational time. Similarly, our proposed ArtNet produces the result in (f) with reduced artifacts and distortion in comparison with WCT [20] in (e) for artistic style transfer.

from high-level to low-level iteratively with style transfer modules. In Figs. 2 (a) and (b), we show the architectures of auto-encoder and multi-level stylization used in popular algorithms. In summary, multi-level stylization processes the images over trained AEs and transfer modules multiple times [20, 21] by utilizing features from high to low levels to improve universal style transfer results.

The improved quality does come with a significant drawback: multi-level stylization has a significantly large computation cost. There is an option to “turn off” the multi-level settings, at the expense of imperfect textures [2], artifacts [12, 30], and distortions [20]. In addition to multi-round stylization, post-processing for photo-realistic stylization [21] is yet another performance bottle-neck in style transfer. It is desirable to design novel approaches to utilize multi-level features of images for better image quality (*i.e.*, fewer artifacts, less distortion, and higher sharpness) while reducing the computational cost for both artistic and photorealistic rendering.

In this work, we propose a novel neural network architecture on top of common multi-layer AEs for fast universal style transfer using multi-level features with pre-defined style transfer modules (e.g., AdaIN [12], WCT [20], and PhotoWCT [21]), while avoiding the use of multi-round stylization and post-processing. We also design novel auto-encoder architectures with superior reconstruction capacity that can alleviate the artifacts and distortions with bet-

ter stylization performance. With such AE and pre-defined transfer modules, we can obtain high-quality style transferred images using multi-level features but with *single-round* computation for much reduced computation.

Specifically for artistic stylization tasks, we first introduce a feed-forward network named **ArtNet** based on feature aggregation [42] to achieve visually pleasing artistic stylization results by eliminating the artifacts and distortions of images. With a pre-trained encoder, the proposed method first extracts and aggregates features ranging from low-level and high-level of encoder parts in AE to better preserve the details of the images (see also in Fig. 2(c)). Next, to obtain better stylization, ArtNet replaces the time-consuming multi-round stylization with a single-round multi-stage decoder with transfer modules embedded, where a “*sandwich structure*” that segments every two stages of the decoder with a transfer module is adopted. Please refer to Figs. 2 (a), (b) and (c) for the detailed comparisons of the two architectures.

For photorealistic image transfer tasks, we propose **PhotoNet** that further extends ArtNet with more sophisticated connections to produce remarkable quality improvement for photorealistic stylization while avoiding the use of any post-processing. In particular, PhotoNet incorporates the additional skip connections coupled with instance normalization [37], passing the low-level features (which may be diminished after multi-layer decoding) directly to the decoder

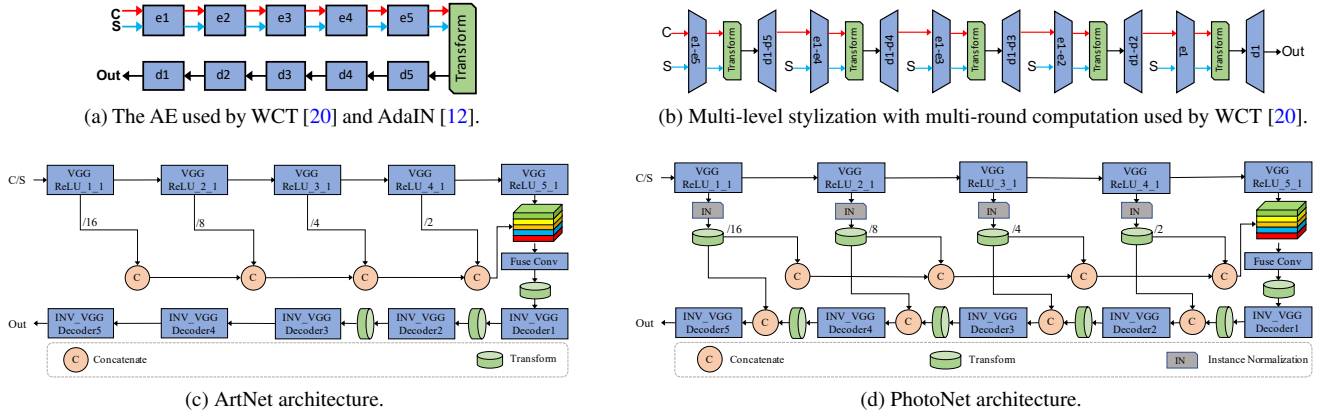


Figure 2: **Comparison of architectures.** The multi-level stylization scheme first trains an auto-encoder (AE) shown in (a) with an image reconstruction loss, then runs multi-round of AE with the WCT [20] transform for style transfer (shown in (b)). Our proposed ArtNet (c) introduces deep feature aggregation and multi-stage stylization on the decoder to better stylize images, while PhotoNet (d) utilizes additional normalized skip connections to preserve more details in style transfer. Please note that networks shown in (c) and (d) are inference-time architectures of ArtNet and PhotoNet, respectively. For training, the transform operations (green cylinders) are excluded while the networks are trained as common AEs using the reconstruction loss.

stages respectively so as to achieve exact image reconstruction. In this way with pre-defined transfer modules, PhotoNet makes the transferred images sharper with more details from the content image preserved. Figs. 2 (a), (b) and (d) offered more details.

Our extensive experiments based on subjective and objective metrics show that ArtNet outperforms the artistic stylization results of AdaIN [12] and WCT [20] when using these two as transfer modules. For photo-realistic stylization, PhotoNet with WCT [21] as transfer modules has a unique capability of rendering semantically consistent, sharp images while avoiding the artifacts or too much smoothness. In addition to superior performance against the state-of-the-art methods, PhotoNet achieves more than 50 $\times$  speed-up in terms of computational time.

Our main contributions are summarized as follows:

- We propose ArtNet for artistic style transfer based on deep feature aggregation and multi-stage stylization on the decoder. ArtNet remarkably outperforms AdaIN [12] with significantly fewer artifacts. ArtNet with the WCT module achieves visually more pleasing results in less distortion while cutting down the time-consumption of the WCT [20] to a third.
- We further propose PhotoNet with additional skip connections coupled with instance normalization. PhotoNet produces rich-detailed, locally-consistent, accurately-stylized photorealistic images with 1/50th of the time-consumption against the state-of-the-art algorithms [25, 21].

## 2. Related Work

We first review the most relevant work to our study and discuss the contribution made by our work.

**Artistic style transfer.** Prior to the adoption of deep neural networks, several classical models based on stroke rendering [9], image analogy [10, 32, 31, 5, 22], or image filtering [41] have been proposed to make a trade-off between quality, generalization, and efficiency for artistic style transfer. In addition to the work already mentioned in the introduction, numerous Gram based algorithms have been recently developed inspired by Gatys *et al.* [6, 7]. Such methods can be classified into one style per model [17, 37, 15, 36, 38, 40, 27, 18], multi-style per model [4, 1], and universal stylization methods [2, 12, 20, 8] with respect to the generalization ability.

**Photo-realistic style transfer.** In terms of methodologies, existing photo-realistic stylization methods [25, 21] either introduce smoothness-based loss term [25] or utilize post-processing to smooth the transferred images [21], which inevitably decreases the sharpness of images and increases the time-consumption significantly. In addition to style transfer, photo-realistic stylization has also been studied in image-to-image translation [14, 39, 24, 35, 33, 23, 44, 13]. The major difference between photo-realistic style transfer and image-to-image translation is that photo-realistic style transfer does not require paired training data (i.e., pre-transfer and post-transfer images). Of course, image-to-image translation can solve even more complicated task such as the man-to-woman and cat-to-dog adaptation problems.

**Discussion.** The work most relevant to our study in-

cludes WCT [20], AdaIN [12] and PhotoWCT [21], while the first two have been used for artistic stylization and the last one is for photo-realistic stylization. Specifically, ArtNet consists of a new network architecture that can incorporate the transfer modules of WCT [20] and AdaIN [12]. The proposed method can avoid multi-round stylization while ensuring the effectiveness of style transfer. PhotoNet also incorporates the transfer module of PhotoWCT [21] in our newly-proposed architecture. It does not require time-consuming multi-round stylization and post-processing. The images produced by PhotoNet have considerably higher sharpness, reduced distortion and significant reduction of computational cost.

### 3. ArtNet and PhotoNet: Architectural Approaches for Fast Universal Style Transfer

The network architectural design of ArtNet and PhotoNet consists of three elements or operations: (i) Deep Feature Aggregation, (ii) Multi-stage Stylization, and (iii) Normalized Skip Connection, as detailed below.

**Deep Feature Aggregation.** Inspired by the PSP-Net [43] and DLA [42] for semantic segmentation, we introduce the deep feature aggregation operation to concatenate the multi-level features and improve the image reconstruction quality of the auto-encoder.

**Multi-stage Stylization.** We utilize the transfer module at different stages of the decoder besides in the middle of the auto-encoder to improve style transfer effects and as a replacement of the multi-level stylization used by WCT [20] and PhotoWCT [21] to speed up the algorithm.

**Normalized Skip Connection.** We employ skip connections as used in [28] and coupled with instance normalization [37] to encourage the decoder to preserve more details when reconstructing images.

#### 3.1. ArtNet architectural design

In order to utilize the multi-level features and avoid the usage of the time-consuming multi-level stylization, we introduce two strategies to ArtNet to improve the image reconstruction quality and thereby enhance the artistic stylization performance.

As Fig. 2 (c) shows, the proposed ArtNet uses the pre-trained VGG-19 as the encoder and utilizes a structurally symmetric decoder to invert deep features back to images. We apply deep feature aggregation to concatenate and fuse multi-level features. Moreover, two additional transfer modules are placed at the end of the first two stages of the decoder to improve the stylization performance and as a replacement of the multi-level stylization strategy in WCT [20]. Fig. 8 (c) demonstrates that ArtNet outperforms the auto-encoder used by AdaIN [12] and WCT [20] regarding the quality of image reconstruction.

#### 3.2. PhotoNet architectural design

The vanilla auto-encoder and ArtNet are not sufficient to reconstruct large amounts of fine details due to the distortion of lines and shapes in the inverted images (Fig. 8). Although such distortions make the generated images look more like art creations in the artistic stylization task, they seriously harm the photorealistic stylization effects due to the decrease of the visual authenticity of images.

On the basis of ArtNet, our PhotoNet additionally adopts normalized skip connections to straightforwardly introduce low-level information from the encoder to the corresponding decoder stages to improve the image reconstruction effects. Moreover, we place transfer modules at normalized skip connections and every stage of the decoder to improve the stylization quality as shown in Fig. 2 (d). As demonstrated in Fig. 8 (e), the proposed PhotoNet outperforms the auto-encoder used by AdaIN [12] and PhotoWCT [21] in faithfully inverting deep features back to images. With PhotoNet coupled with the ZPA transform in WCT [20]/PhotoWCT [21] as the feature stylization module, PhotoNet avoids the use of post-processing as well as time-consuming optimization and achieves visually pleasing transfer results with rich details and is more than  $600\times$  faster than Luan *et al.* [25] and  $50\times$  faster than Li *et al.* [21] on large images.

#### 3.3. Decoder training

We train the decoder of ArtNet and PhotoNet to invert deep features back to images with the Frobenius norm of the original and inverted image as the reconstruction loss,

$$\mathcal{L}_{recon} = \|I_{in} - I_{out}\|_F, \quad (1)$$

where  $I_{in}$  denotes the input image,  $I_{out}$  is reconstructed image, and  $\|\cdot\|_F$  represents the Frobenius norm. Inspired by Li *et al.* [20], we introduce the perceptual loss term [15] to improve the reconstruction quality of the decoder,

$$\mathcal{L}_{percep} = \sum_{i=1}^5 \|\Phi_i(I_{in}) - \Phi_i(I_{out})\|_F, \quad (2)$$

where  $\Phi_i(\cdot)$  denotes the output of the  $i^{th}$  stage of the ImageNet [3] pre-trained VGG-19 [34]. The overall loss function is,

$$\mathcal{L} = \alpha \mathcal{L}_{recon} + (1 - \alpha) \mathcal{L}_{percep}, \quad (3)$$

where the  $\alpha$  balances two loss terms. During training, all transfer modules in ArtNet and PhotoNet are skipped and the whole framework is trained in an image reconstruction manner.

## 4. Experiments and Empirical Validations

In this section, we present the stylization results of our proposed algorithms in comparison with those of the state-



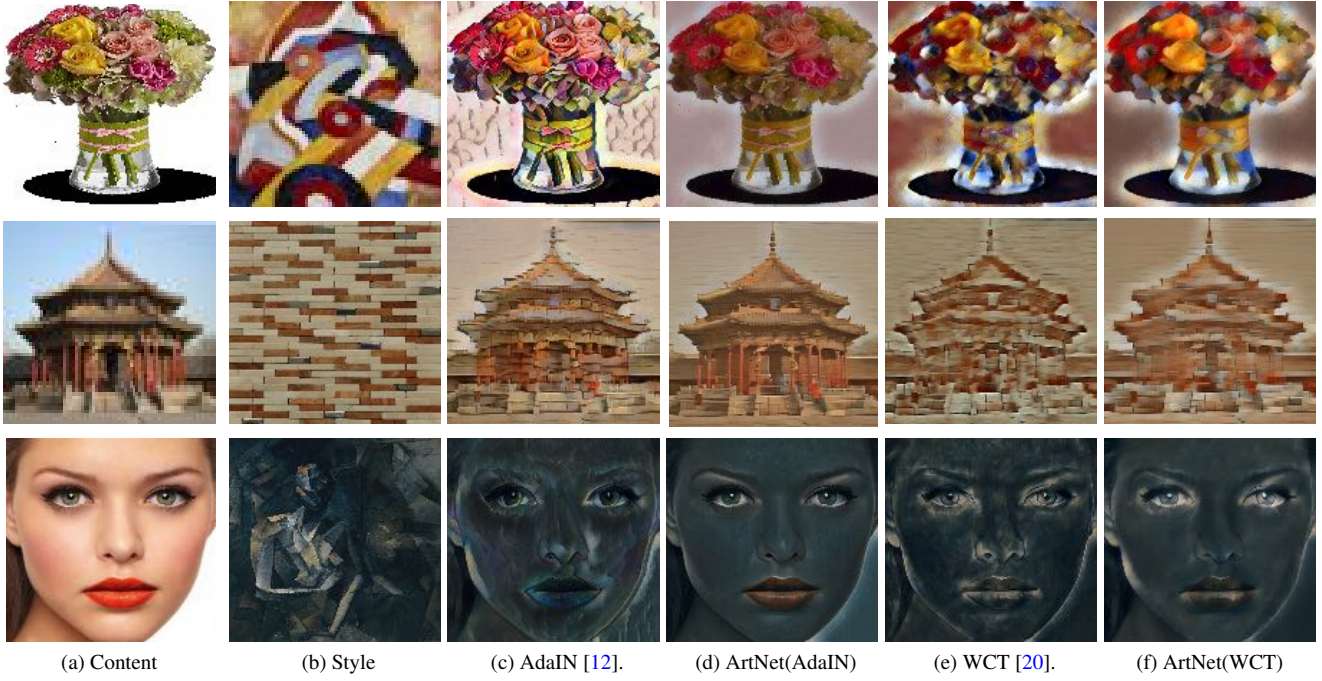


Figure 3: Results of the contrast experiments against baseline artistic style transfer methods.

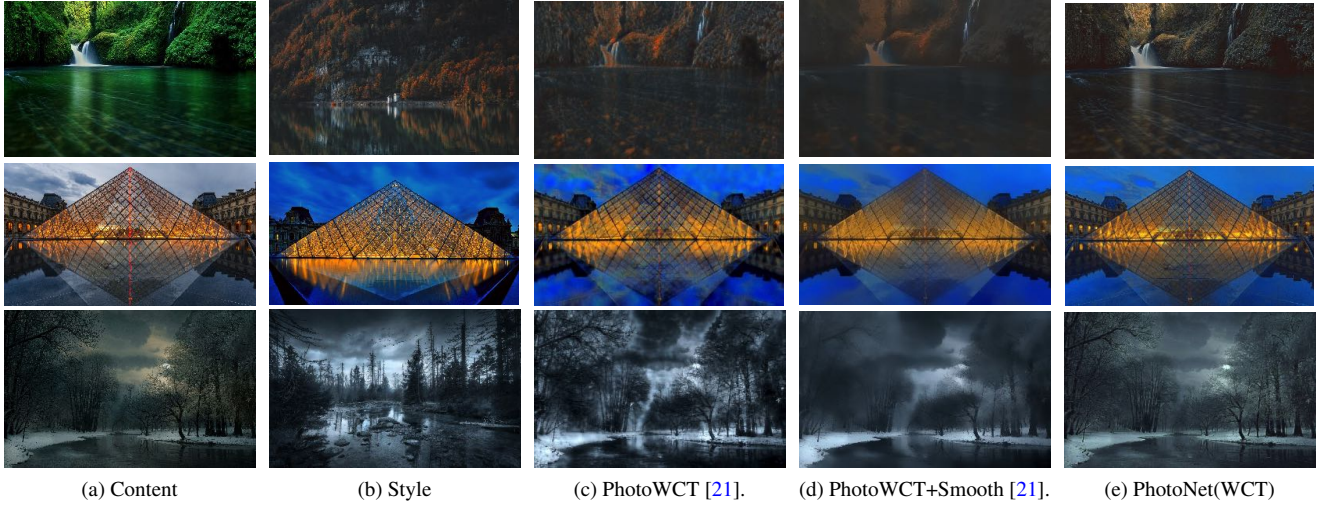


Figure 4: Results of the contrast experiments against baseline photorealistic style transfer method.

of-the-art approaches and further provide a comprehensive empirical analysis to substantiate our observations.

#### 4.1. Results on artistic style transfer

In order to demonstrate the effectiveness of the proposed ArtNet, we conduct contrast experiments on AdaIN [12] and WCT [20], where we replace the AE part of these two approaches with ArtNet and keep other part fixed. As Fig. 3 shows, the results by [12] contains unpleasant artifacts such as unfaithful edge lines in brick edges (Row 2), the face (Row 3), and the background of flowers. Such

significant artifacts are clearly eliminated by the proposed method when integrating ArtNet with AdaIN [12] as the transform. The WCT [20] method generates images with significant distortions and a lack of local similarity, fragmented color blobs (Row 1) and twisted lines (Row 2). In contrast, ArtNet with WCT [21] as transform creates images with straight lines (Row 2), clear color blobs (Row 1) and a clean face (Row 3). To demonstrate the exceptional performance of ArtNet coupled with the ZCA in WCT [20] as the transfer module, we compare the artistic stylization results against the state-of-the-art methods. As shown in

Table 1: **Quantitative evaluation results for stylization methods.** The FID and TV scores are only applicable to the photo-realistic methods. A lower FID score means the evaluated method creates an image with a more similar style to the reference style image. A higher total variation score indicates that the measured image has more details.

Method	AdaIN [12] vs. ArtNet(AdaIN)	WCT [20] vs. ArtNet(WCT)	PhotoWCT [21] vs. PhotoNet(WCT)
Preference $\uparrow$	14.58% / <b>85.42%</b>	3.13% / <b>96.87%</b>	6.25% / <b>93.75%</b>
TV score $\uparrow$	-	-	5.15 / <b>7.11</b>
FID score $\downarrow$	-	-	169.22 / <b>167.06</b>

Table 2: **Computing-time comparison for the artistic and photo-realistic style transfer methods.**

Method	Gatys <i>et al.</i> [7]	WCT [20]	ArtNet(WCT)	Luan <i>et al.</i> [25]	PhotoWCT [21]	PhotoNet(WCT)
$256 \times 128$	6.01	0.36	<b>0.34</b>	114.11	4.07	<b>0.76</b>
$512 \times 256$	17.84	0.71	<b>0.42</b>	293.28	20.72	<b>0.86</b>
$768 \times 384$	38.18	1.20	<b>0.45</b>	628.24	53.05	<b>0.95</b>
$1024 \times 512$	66.24	1.88	<b>0.52</b>	947.61	133.90	<b>1.06</b>

Fig. 5, the method by Gatys *et al.* [7] generates images with areas of artifacts (badly stylized areas in the image of the top row and overexposed background of the human in the image in the bottom row.). The results of AdaIN [12] contain evident artifacts that render inaccurate shapes and color blobs in flowers (Row 1) and unnatural hair of the girl (Row 2). The WCT [20] method distorts and twists the shapes, lines, and color blobs of the transferred images, while AvatarNet [30] creates artifacts by rendering artistically matching patches to the transferred images but disregarding their semantics. This is demonstrated by Fig. 5 (f) where AvatarNet [30] renders “red eyes” (Row 1) as well as “blue eyes and red lips” (Row 2) arbitrarily in the generated images.

## 4.2. Results on photorealistic style transfer

We verify the effectiveness of the proposed PhotoNet (using WCT [20] as transform) by the comparison with the photorealistic stylization results of PhotoWCT [21]. As shown in Fig. 4 (c), the results of PhotoWCT [21] without the post-processing for smoothing contain apparent distortions such that the sky in images in the second and third rows are distorted and twisted. The transferred images of PhotoWCT [21] with the smoothing operation turned on are overly smooth and have low sharpness in details such that grasses (Row 1), steel frame of the Louvre (Row 2), and trees (Row 3) have been smoothed out and lost their details. We additionally compare the results of PhotoNet with the algorithm by Luan *et al.* [25] and PhotoWCT [21] approach. Fig. 6 shows that PhotoNet renders the leaves in trees (Row 2), the texture of woods and details of the camera lens (Row 1) considerably sharper and transfer styles much more faithfully, demonstrating that PhotoNet outperforms the compared methods in generating visually pleasing and sharp images.

## 4.3. Computational time comparison

We conduct a computing time comparison against the state-of-the-art methods to demonstrate the efficiency of

the proposed network architectures. All approaches are tested on the same computing platform which includes an NVIDIA Tesla P100 with 16GB RAM. We compare the computing time on content and style images with different resolutions.

**ArtNet.** As Table 2 shows, Gatys *et al.* [7] is slow due to the optimization process, WCT [20] method is considerably faster but not efficient enough due to the usage of multi-level stylization (especially for high-resolution images), while ArtNet improves the inference speed of WCT [20] by three times on large images, thanks to the avoidance of multi-level stylization. Moreover, ArtNet generates stylized images with fewer artifacts compared with AdaIN [12] with a minor additional time cost.

**PhotoNet.** In order to verify the superior efficiency of the proposed PhotoNet, we conduct experiments against baseline photorealistic stylization methods in terms of the computing time. As Table 2 demonstrated, PhotoNet is hundreds of times faster than the method of Luan *et al.* [25] and tens of times faster than PhotoWCT [21]. It is even more time-efficient on high-resolution images. It is worth mentioning that the method by Luan *et al.* [25] requires additional segmentation masks to assist stylization, which costs more computing time.

## 4.4. Empirical validation

In this section, we try to provide more insights into the proposed methods through three empirical studies: (1) We first use quantitative analysis to demonstrate the quality of stylized images using a user study together with Fréchet Inception Distance (FID) [11] and total variation [29] (TV) scores, where we are specifically interested in the sharpness of the generated images; (2) We then conduct an ablation study to validate the necessary/contribution of each component (*e.g.*, feature aggregation, multi-stage stylization, and normalized skip connections) included in PhotoNet and ArtNet; (3) We add a case study on the image reconstruction performance of AEs used for style transfer to make sense of our simple intuition that lower AE image reconstruction er-



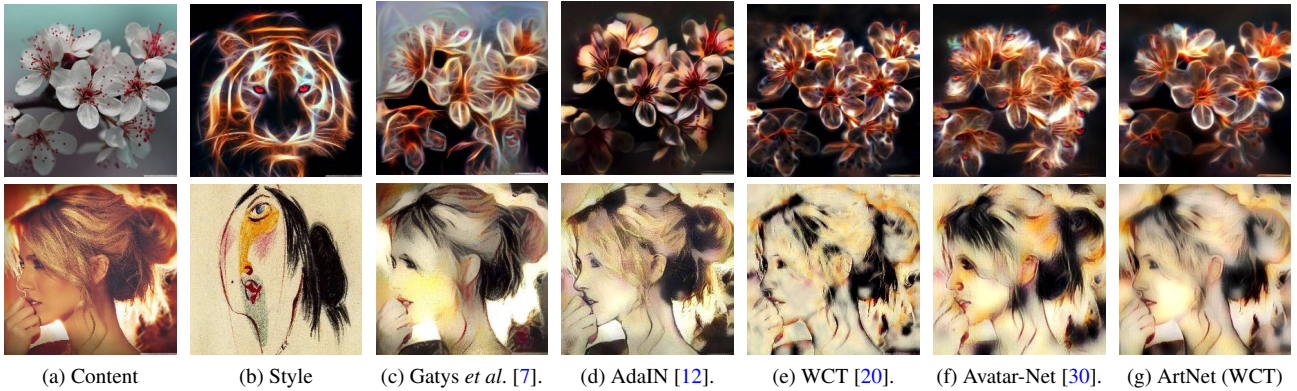


Figure 5: **Artistic stylization results.**

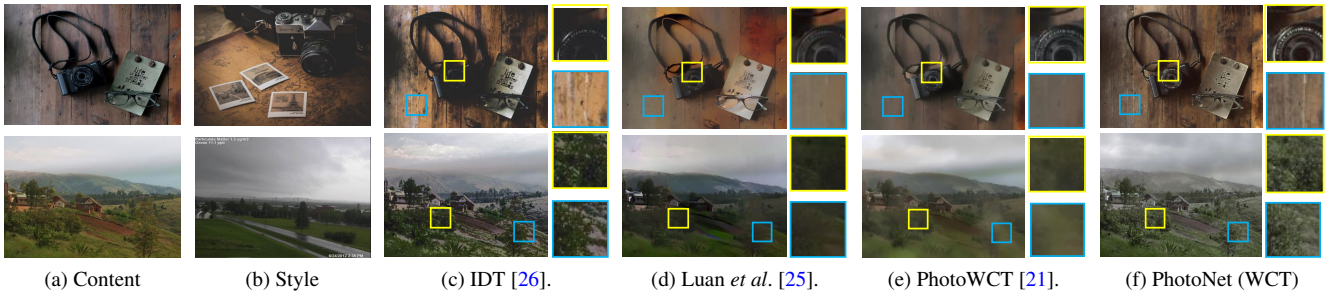


Figure 6: **Photo-realistic stylization results.** Note that the method of Luan *et al.* [25] requires additional segmentation masks for stylization while other compared algorithms do not.

ror leads to better stylization performance.

#### 4.4.1 Quantitative evaluation

In this study, artistic style transfer methods are evaluated on a dataset consisting of 12 content images and 16 style images, where each content image is transferred into every style. We compute the evaluation metrics and conduct the user study based on totally 192 generated images. As for photo-realistic cases, the evaluation is based on 38 content images and their corresponding styles.

**User study.** We conducted a user study to subjectively demonstrate the effectiveness of the proposed ArtNet and PhotoNet. We randomly select 12 content and style image pairs to evaluate the artistic style transfer methods and use 6 pairs to measure photo-realistic approaches. For each content and style pairs, we display the results of AdaIN [12]/ArtNet(AdaIN), WCT [20]/ArtNet(WCT), and PhotoWCT [21]/PhotoNet(WCT) side-by-side and let the subject choose the better one in terms of less artifact, less distortion, and more details, respectively. We collect 16 responses and a total of 288 votes. The preference percentage of the choices are summarized in Table 1, which demonstrates that using ArtNet improves over the stylization results of WCT [20] in terms of less distortion and AdaIN [12] in terms of fewer artifacts, while PhotoNet improves over the results of PhotoWCT [21] in terms of more sharp details.

**FID [11].** We compute the FID score between the reference style images and transferred images by PhotoNet and state-of-the-art methods, *i.e.*, PhotoWCT [21] for comparison. As Table 1 shows, PhotoNet using WCT [20] as the transfer module outperforms PhotoWCT [21] with a higher FID score (*i.e.*, better stylization). Note that FID was originally used to validate the image quality for domain adaption and image translation, which is close to photorealistic stylization.

**Total variation [29].** We compare the total variation scores of the results by PhotoNet and PhotoWCT [21] methods. As demonstrated in Table 1, images generated by PhotoNet are of higher total variation scores (*i.e.*, more sharpness and details) than PhotoWCT [21].

#### 4.4.2 Ablation study

Note that we denote the use of feature aggregation module as *FA*, the use of normalized skip connections as *NS*, and *MST-X* for the multi-stage style transfer module, where  $X = 3, 5, \infty$  refers to the incorporation of transfer modules (like WCT [20] or AdaIN [12]) in the first three stage of the decoder, all stages of the decoder, and all normalized skip connections, respectively.

**ArtNet.** As was shown in Fig. 7 Row 1, FA and MST-3 improve the stylization effects in succession. However, as shown in Fig. 7 (e), placing the transfer module in the low-level stages of the decoder hurt the results of the stylization.

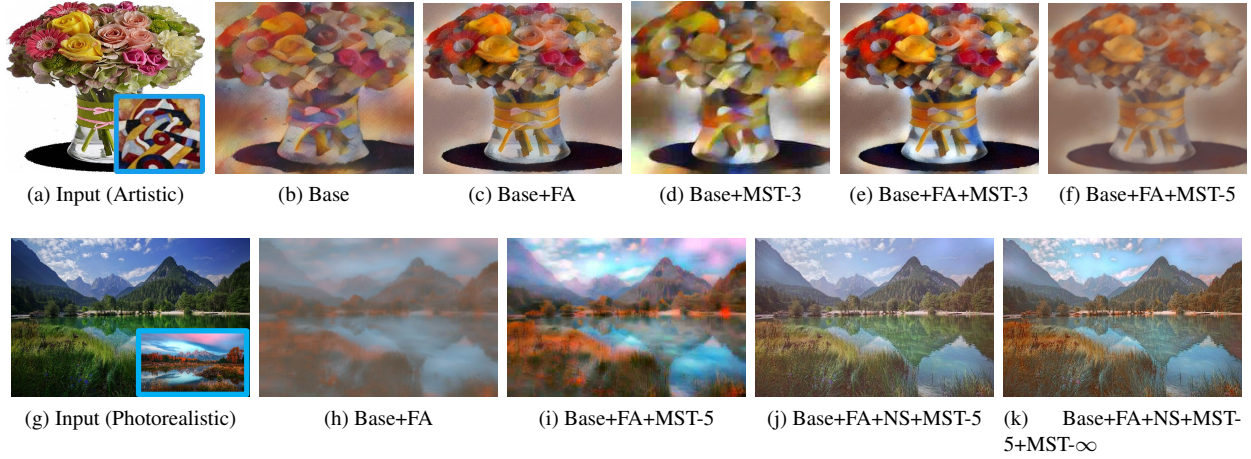


Figure 7: Ablation study of ArtNet and PhotoNet.

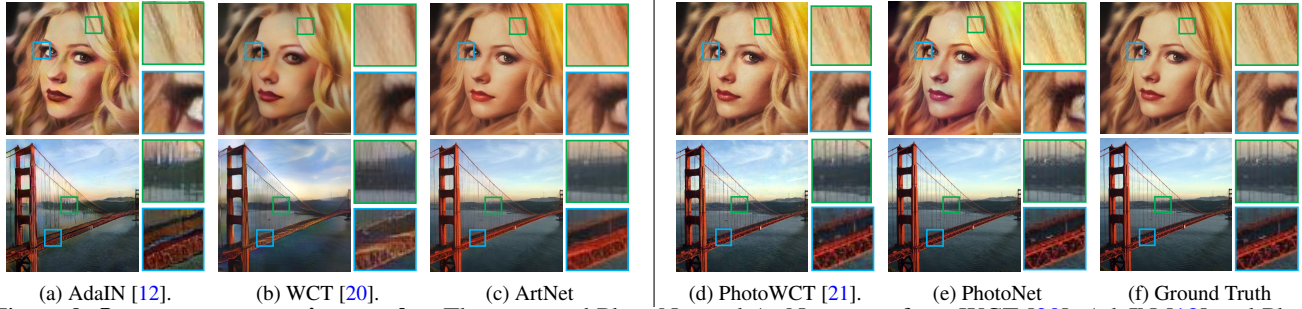


Figure 8: **Image reconstruction results.** The proposed PhotoNet and ArtNet outperform WCT [20], AdaIN [12] and PhotoWCT [21] in preserving details (*i.e.*, hairs, eye-slashes, and the local structure/textures in both artistic and photorealistic settings).

Table 3: **Image reconstruction evaluation.** All the listed algorithms are evaluated on the  $512 \times 512$  images.

Method	AdaIN	WCT	PhotoWCT	ArtNet	PhotoNet
Error	86611.57	83196.87	76755.60	<b>75021.59</b>	<b>74643.66</b>

**PhotoNet.** We present the ablation study results of PhotoNet in Fig. 7 Row 2, which demonstrates the effectiveness of NS. Moreover, Fig. 7 (i) shows that conducting style transfer in normalized skip connections can further improve the visual effects of the transferred images (red flowers and blue sky in the image are highlighted by MST- $\infty$ ).

#### 4.4.3 Image reconstruction

One major finding of our work is that lower image reconstruction error of AEs leads to better stylization performance. We present the images generated by vanilla AEs used in AdaIN [12], WCT [20], PhotoWCT [21], as well as ArtNet and PhotoNet in Fig. 8. Images generated by AdaIN [12] contains significant artifacts in both two cases, *e.g.*, changes in the lip color of the girl in the top row and twisted cables/steel frame of the bridge in the bottom row. The WCT [20] method distorts images where the hairlines of the girl and cables of the bridge are blurred. ArtNet

renders the clearest reconstruction result among the compared artistic stylization methods. As for photorealistic cases, PhotoNet outperforms PhotoWCT [21] in terms of the sharpness in details, *e.g.*, the steel frame of the bridge, hairlines, and eyelash of the girl contain clearer details. We quantitatively evaluate the performance of the proposed algorithms against the state-of-the-art baseline methods by computing the mean squared error as defined by Eq. 4 between the original and the reconstructed images on a randomly selected dataset.

$$error = \sum_{i=1}^N (\|I_{in} - I_{out}\|_F) / N, \quad (4)$$

where  $N$  denotes the number of the selected images and  $N = 13$  here. As shown in Table 3, the AE used in ArtNet and PhotoNet achieve better image reconstruction performance compared to all other methods.

## 5. Conclusions

In this paper, we present two network architectures to address artistic and photorealistic style transfer, respectively.



ArtNet outperforms the artistic stylization results of the existing methods by introducing a feature aggregation operation and a multi-stage stylization module, which also avoids the use of multi-round computation for stylization to speed up the transfer process. In addition, PhotoNet utilizes normalized skip connections to preserve details of the transferred images, thus generating rich-detailed and well-stylized images. Our extensive experiments include visual comparisons, quantitative comparisons, and a thorough ablation study to show that the proposed approaches have the ability to remarkably improve the stylization effects for both artistic and photo-realistic stylization while reducing the time consumption dramatically especially for photo-realistic transfer algorithms. In the future, we will try to combine the proposed networks with newly proposed style transfer modules such as Avatar-Net [30] and the method by Gu *et al.* [8] to further improve the style transfer results.

## References

- [1] D. Chen, L. Yuan, J. Liao, N. Yu, and G. Hua. Stylebank: an explicit representation for neural image style transfer. In *CVPR*, 2017. 3
- [2] T. Q. Chen and M. Schmidt. Fast patch-based style transfer of arbitrary style. *arXiv preprint arXiv:1612.04337*, 2016. 1, 2, 3
- [3] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: a large-scale hierarchical image database. In *CVPR*, 2009. 4
- [4] V. Dumoulin, J. Shlens, and M. Kudlur. A learned representation for artistic style. In *ICLR*, 2017. 1, 3
- [5] O. Frigo, N. Sabater, J. Delon, and P. Hellier. Split and match: example-based adaptive patch sampling for unsupervised style transfer. In *CVPR*, 2016. 3
- [6] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015. 1, 3
- [7] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *CVPR*, 2016. 1, 3, 6, 7
- [8] S. Gu, C. Chen, J. Liao, and L. Yuan. Arbitrary style transfer with deep feature reshuffle. In *CVPR*, 2018. 1, 3, 9
- [9] A. Hertzmann. Painterly rendering with curved brush strokes of multiple sizes. In *SIGGRAPH*, 1998. 3
- [10] A. Hertzmann, C. E. Jacobs, N. Oliver, B. Curless, and D. H. Salesin. Image analogies. In *SIGGRAPH*, 2001. 3
- [11] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, 2017. 6, 7
- [12] X. Huang and S. J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 1, 2, 3, 4, 5, 6, 7, 8, 11, 14, 15, 16
- [13] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz. Multimodal unsupervised image-to-image translation. In *ECCV*, 2018. 3
- [14] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 3
- [15] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 1, 3, 4
- [16] D. P. Kingma and J. Ba. Adam: a method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 11
- [17] C. Li and M. Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *CVPR*, 2016. 3
- [18] S. Li, X. Xu, L. Nie, and T.-S. Chua. Laplacian-steered neural style transfer. In *ACM MM*, 2017. 3
- [19] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Diversified texture synthesis with feed-forward networks. In *CVPR*, 2017. 1
- [20] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Universal style transfer via feature transforms. In *NIPS*, 2017. 1, 2, 3, 4, 5, 6, 7, 8, 11, 17, 18, 19
- [21] Y. Li, M.-Y. Liu, X. Li, M.-H. Yang, and J. Kautz. A closed-form solution to photorealistic image stylization. In *ECCV*, 2018. 1, 2, 3, 4, 5, 6, 7, 8, 11, 20, 21
- [22] J. Liao, Y. Yao, L. Yuan, G. Hua, and S. B. Kang. Visual attribute transfer through deep image analogy. *arXiv preprint arXiv:1705.01088*, 2017. 3
- [23] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *NIPS*, 2017. 3
- [24] M.-Y. Liu and O. Tuzel. Coupled generative adversarial networks. In *NIPS*, 2016. 3
- [25] F. Luan, S. Paris, E. Shechtman, and K. Bala. Deep photo style transfer. In *CVPR*, 2017. 1, 3, 4, 6, 7
- [26] F. Pitie, A. C. Kokaram, and R. Dahyot. N-dimensional probability density function transfer and its application to color transfer. In *ICCV*, 2005. 7
- [27] E. Risser, P. Wilmot, and C. Barnes. Stable and controllable neural texture synthesis and style transfer using histogram losses. *arXiv preprint arXiv:1701.08893*, 2017. 1, 3

- [28] O. Ronneberger, P. Fischer, and T. Brox. U-net: convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-assisted Intervention*, 2015. 4
- [29] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992. 6, 7
- [30] L. Sheng, Z. Lin, J. Shao, and X. Wang. Avatar-net: multi-scale zero-shot style transfer by feature decoration. In *CVPR*, 2018. 1, 2, 6, 7, 9
- [31] Y. Shih, S. Paris, C. Barnes, W. T. Freeman, and F. Durand. Style transfer for headshot portraits. *ACM Transactions on Graphics*, 33(4):148, 2014. 3
- [32] Y. Shih, S. Paris, F. Durand, and W. T. Freeman. Data-driven hallucination of different times of day from a single outdoor photo. *ACM Transactions on Graphics*, 32(6):200, 2013. 3
- [33] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. In *CVPR*, 2017. 3
- [34] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4
- [35] Y. Taigman, A. Polyak, and L. Wolf. Unsupervised cross-domain image generation. In *ICLR*, 2017. 3
- [36] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. S. Lempitsky. Texture networks: feed-forward synthesis of textures and stylized images. In *ICML*, 2016. 1, 3
- [37] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: the missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 2, 3, 4
- [38] D. Ulyanov, A. Vedaldi, and V. S. Lempitsky. Improved texture networks: maximizing quality and diversity in feed-forward stylization and texture synthesis. In *CVPR*, 2017. 1, 3
- [39] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *CVPR*, 2018. 3
- [40] X. Wang, G. Oxholm, D. Zhang, and Y.-F. Wang. Multimodal transfer: a hierarchical deep convolutional neural network for fast artistic style transfer. In *CVPR*, 2017. 3
- [41] H. Winnemöller, S. C. Olsen, and B. Gooch. Real-time video abstraction. *ACM Transactions on Graphics*, 25(3):1221–1226, 2006. 3
- [42] F. Yu, D. Wang, E. Shelhamer, and T. Darrell. Deep layer aggregation. In *CVPR*, 2018. 2, 4
- [43] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *CVPR*, 2017. 4
- [44] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 3

# Supplementary Material

## A. Network Training Setting

We train the ArtNet and PhotoNet with the reconstruction and perceptual loss functions,

$$\mathcal{L} = \alpha \cdot \mathcal{L}_{recon} + (1 - \alpha) \cdot \mathcal{L}_{precep}, \quad (5)$$

where  $\alpha$  is used to balance tow loss terms. We set  $\alpha = 0.5$  during training. In addition, we use the Adam method [16] and set the learning rate to be  $1e^{-4}$ . We train both the ArtNet and PhotoNet for 5 epoches with the fixed learning rate. The training process spends about eight hours on a NVIDIA Tesla P100 GPU with 16GB GPU RAM.

## B. User Control

We conduct extensive experiments to demonstrate the ability of the proposed ArtNet and PhotoNet that enable flexible user control of the stylization effects as [21, 12, 20]

do. We introduce an user control module at the end of each style transfer module, which mixes the transferred and the content features as,

$$\mathcal{F}_{out} = \beta \cdot \mathcal{F}_{transferred} + (1 - \beta) \cdot \mathcal{F}_{content}, \quad (6)$$

where  $\mathcal{F}$  represents deep feature maps,  $\beta$  is a factor to let the user to control the degree of stylization effects. We present the artistic style transfer results with  $\beta$  ranging from 0.2 to 1.0 in Figs. 910 and the photorealistic stylization results with the same user control setting in Fig. 11. The generated images demonstrate that the proposed ArtNet and PhotoNet achieve incremental style transfer effects with changing control factors, which facilitates the flexible user control to the degree of the proposed style transfer algorithms.

## C. Style Transfer Results

We present more style transfer results of the proposed ArtNet and PhotoNet in comparison to the AdaIN [12], WCT [20], and PhotoWCT [21] respectively.



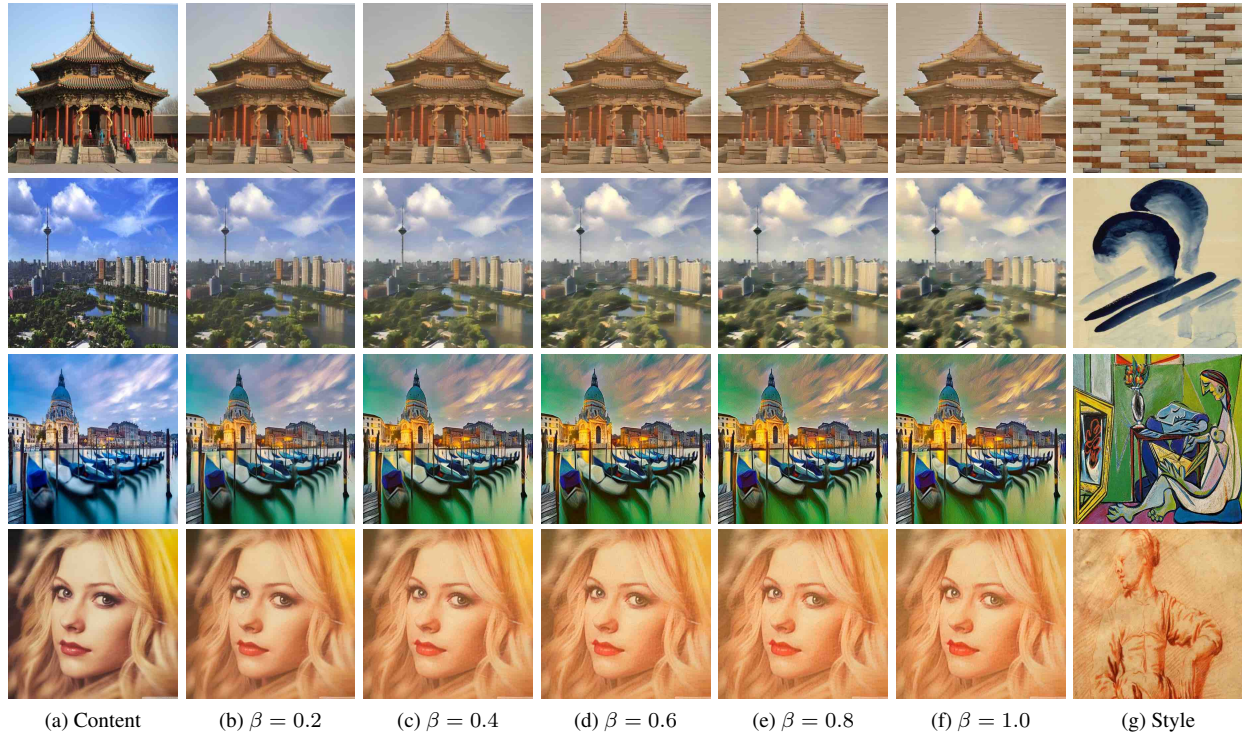


Figure 9: Artistic style transfer results by the ArtNet(AdaIN) with control factor  $\beta$  ranging from 0.2 to 1.0.

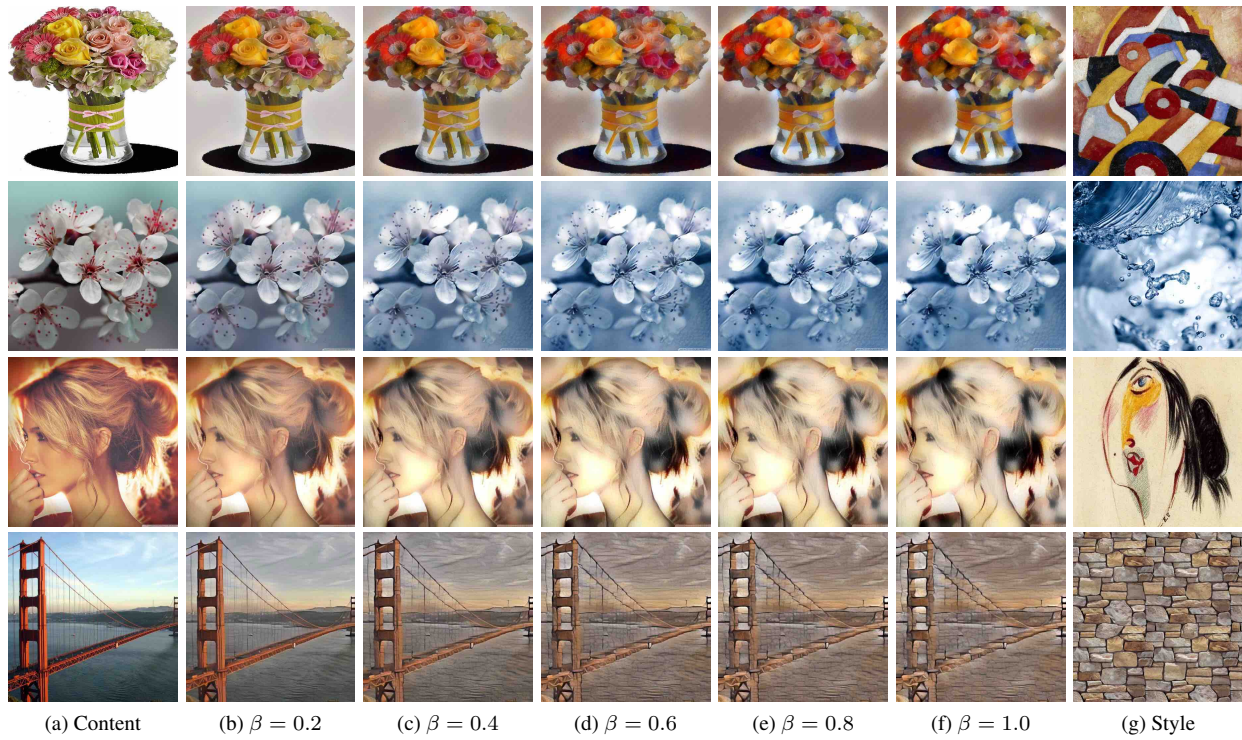


Figure 10: Artistic style transfer results by the ArtNet(WCT) with control factor  $\beta$  ranging from 0.2 to 1.0.

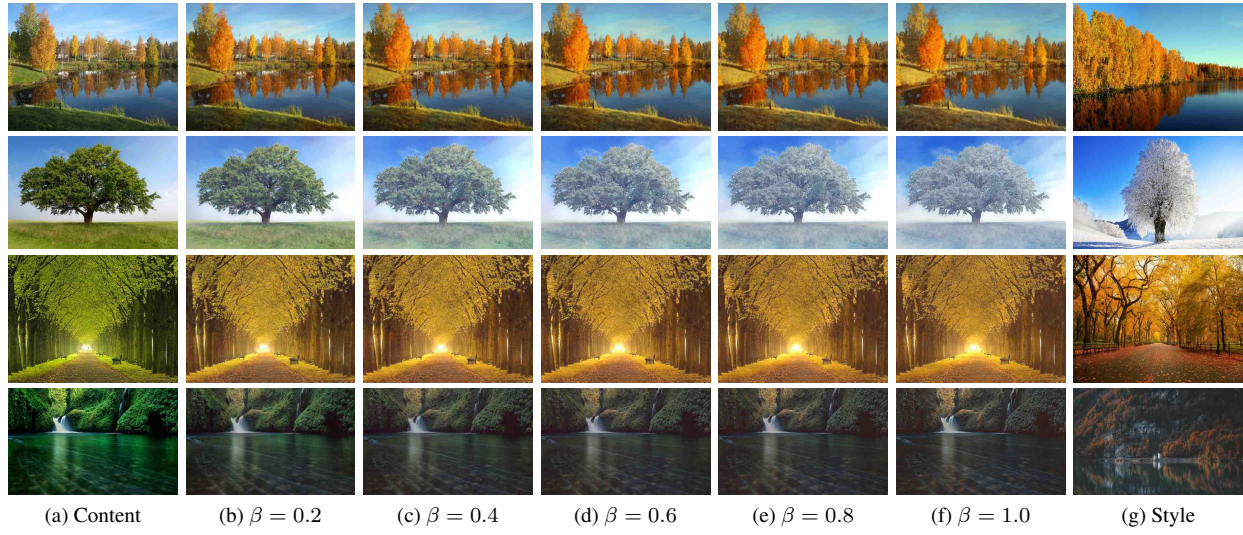


Figure 11: Photorealistic style transfer results by the PhotoNet(WCT) with control factor  $\beta$  ranging from 0.2 to 1.0.



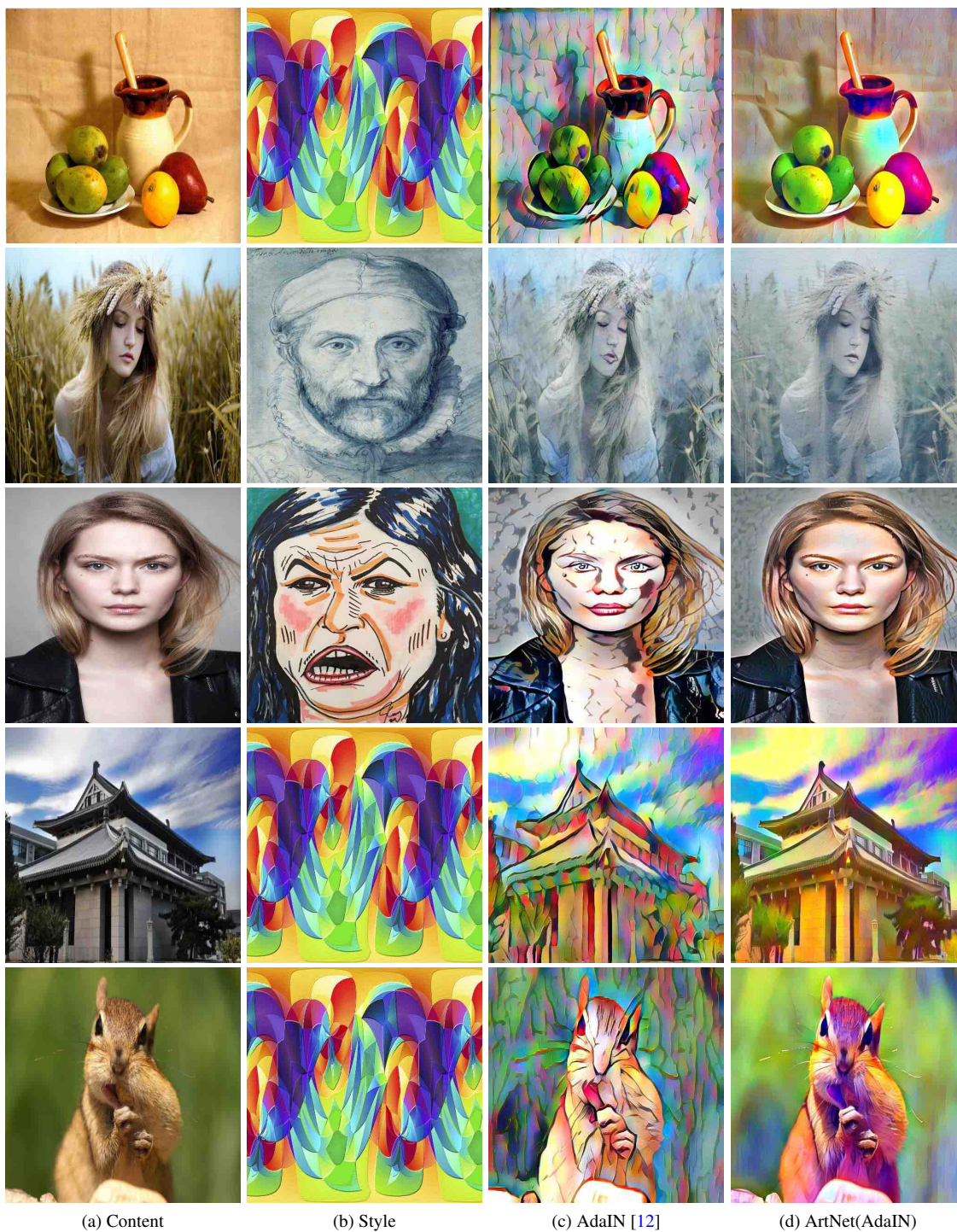


Figure 12: Artistic style transfer comparison between the ArtNet(AdaIN) and AdaIN [12].



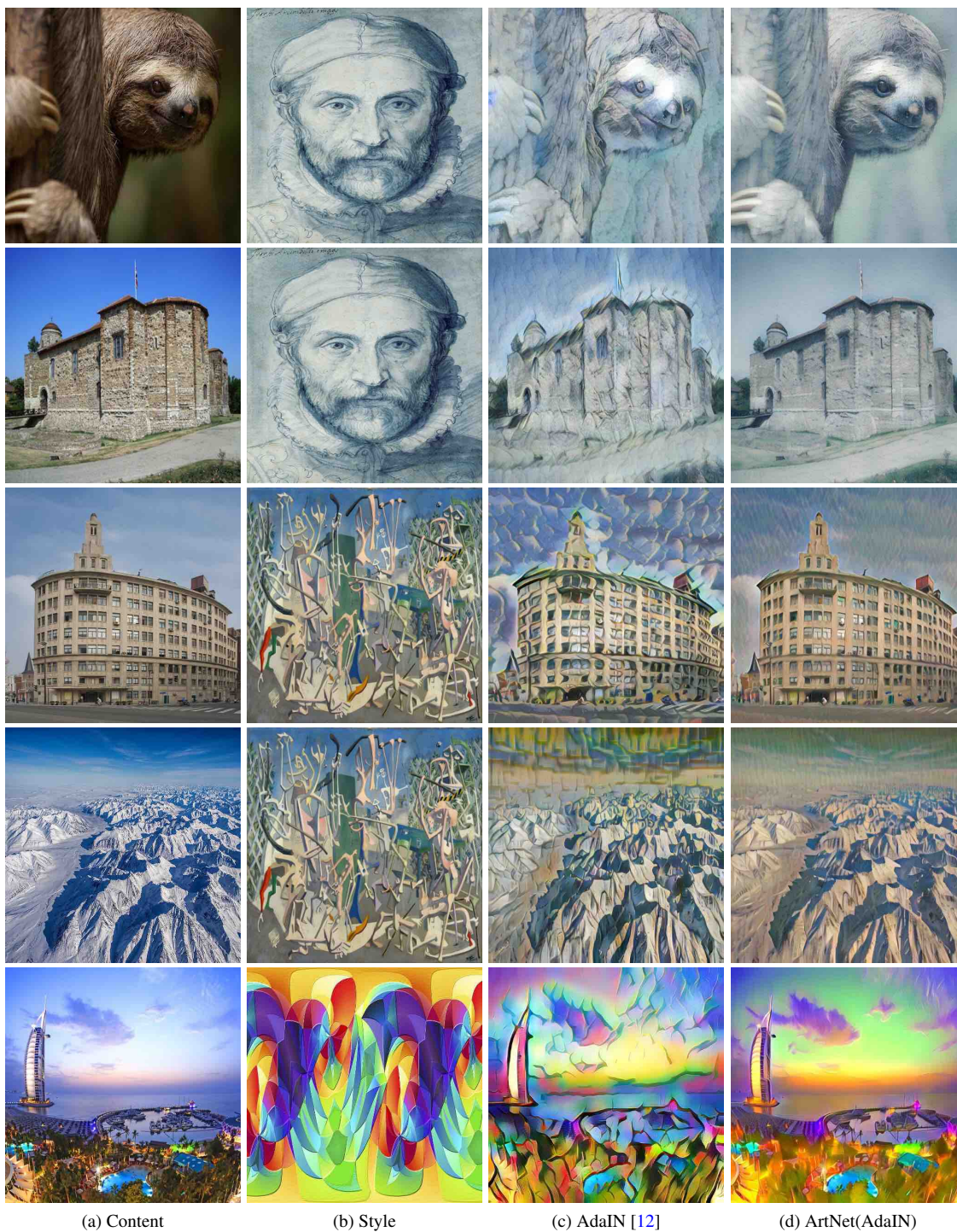


Figure 13: Artistic style transfer comparison between the ArtNet(AdaIN) and AdaIN [12].



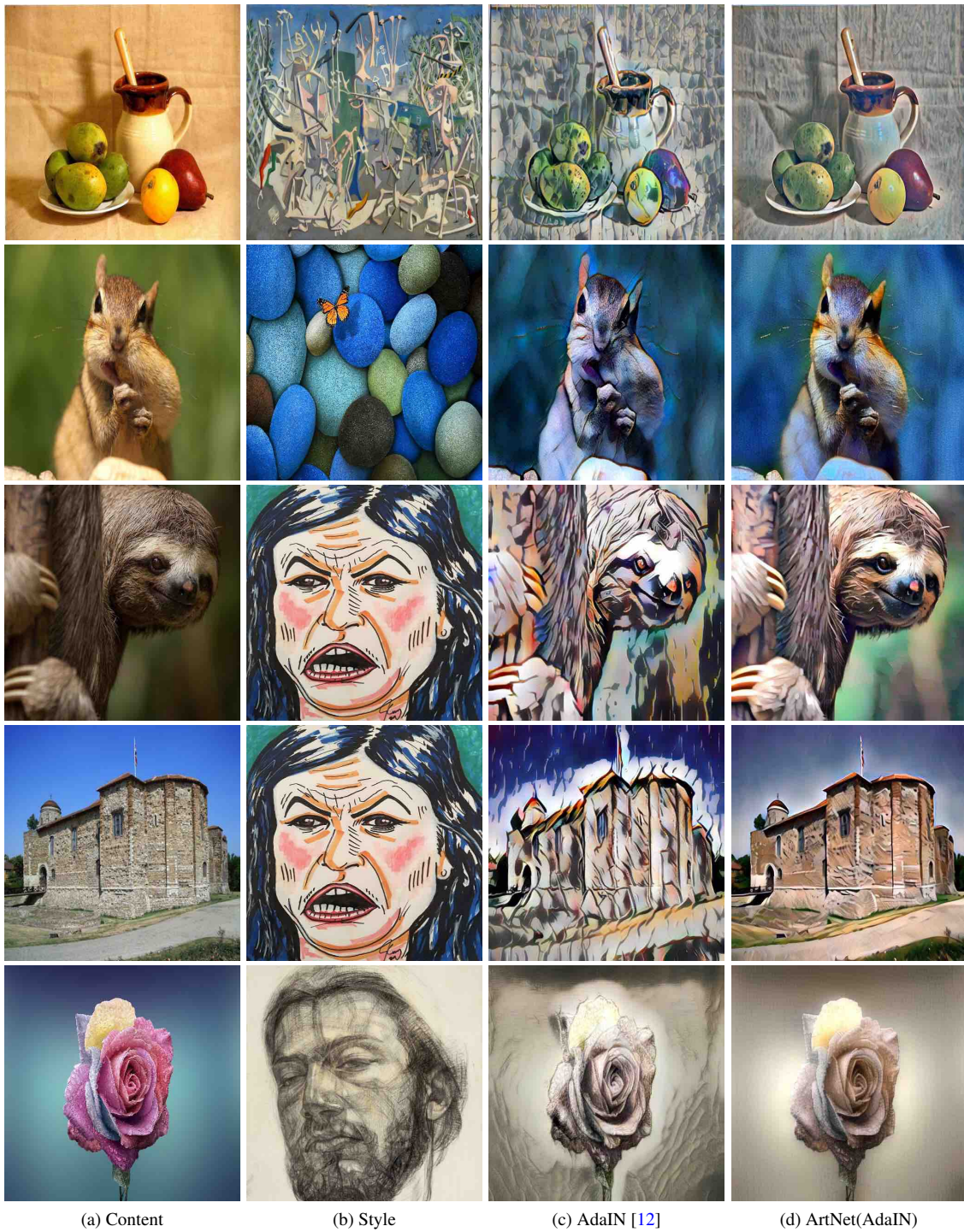


Figure 14: Artistic style transfer comparison between the ArtNet(AdaIN) and AdaIN [12].



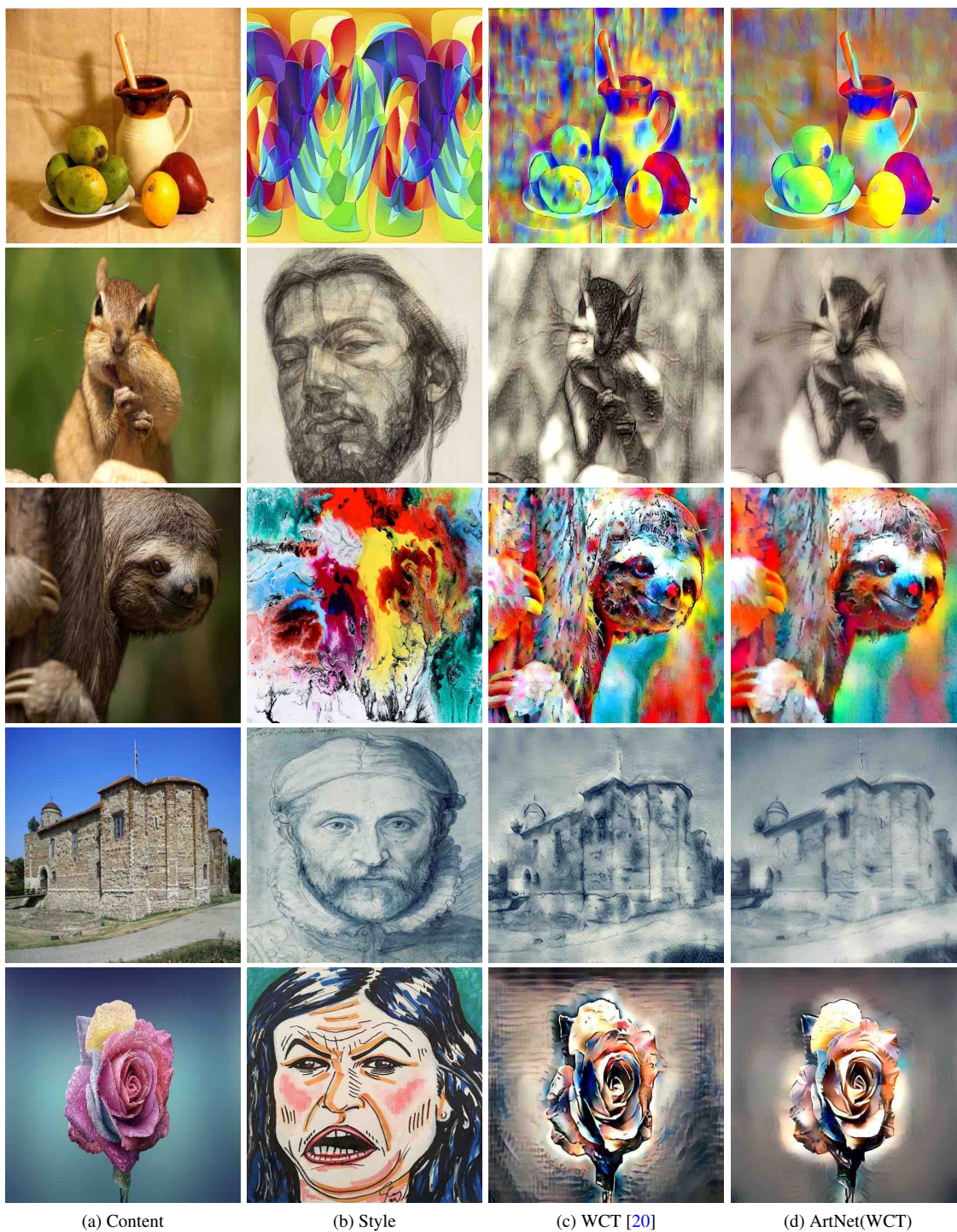


Figure 15: Artistic style transfer comparison between the ArtNet(WCT) and WCT [20].



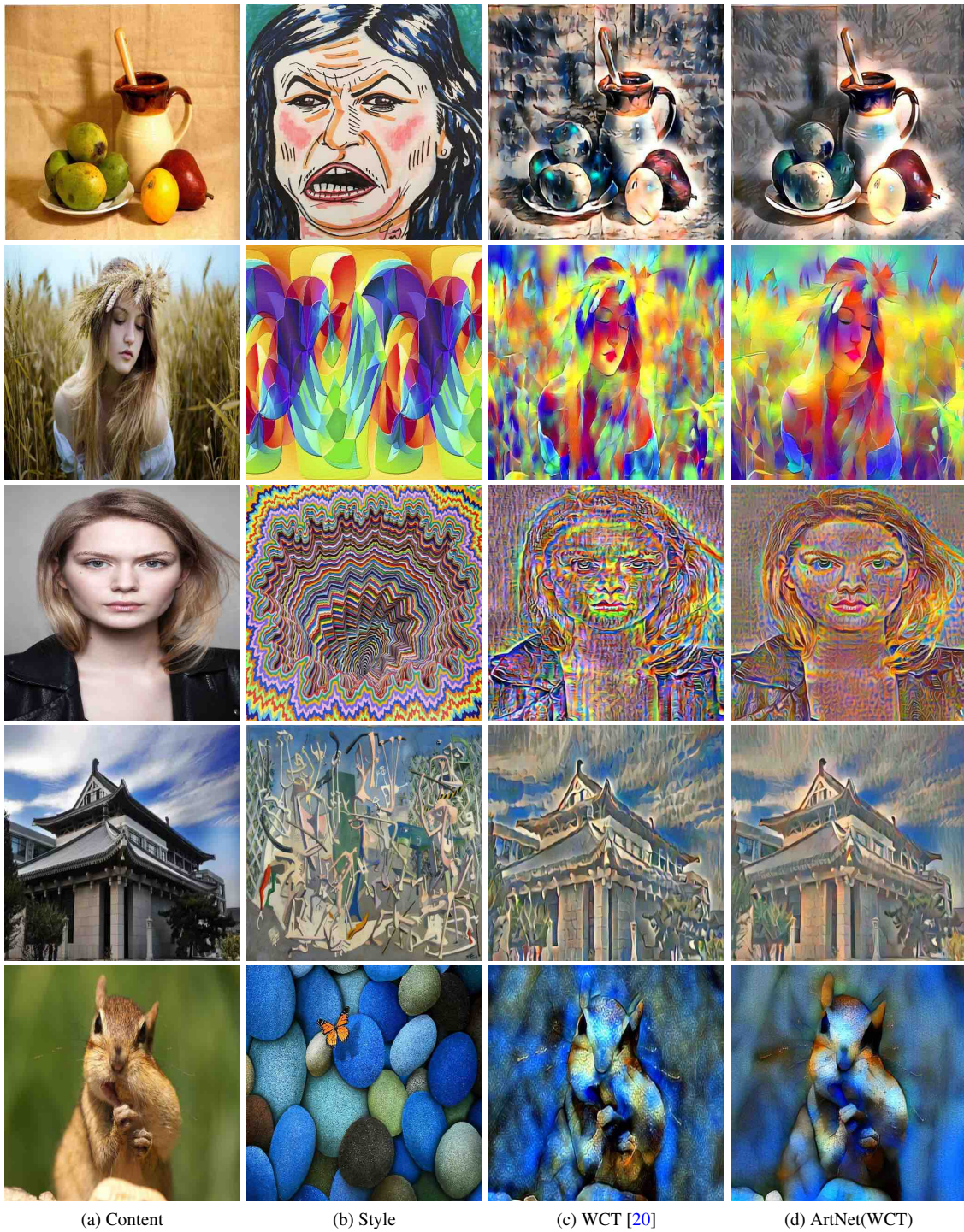


Figure 16: Artistic style transfer comparison between the ArtNet(WCT) and WCT [20].



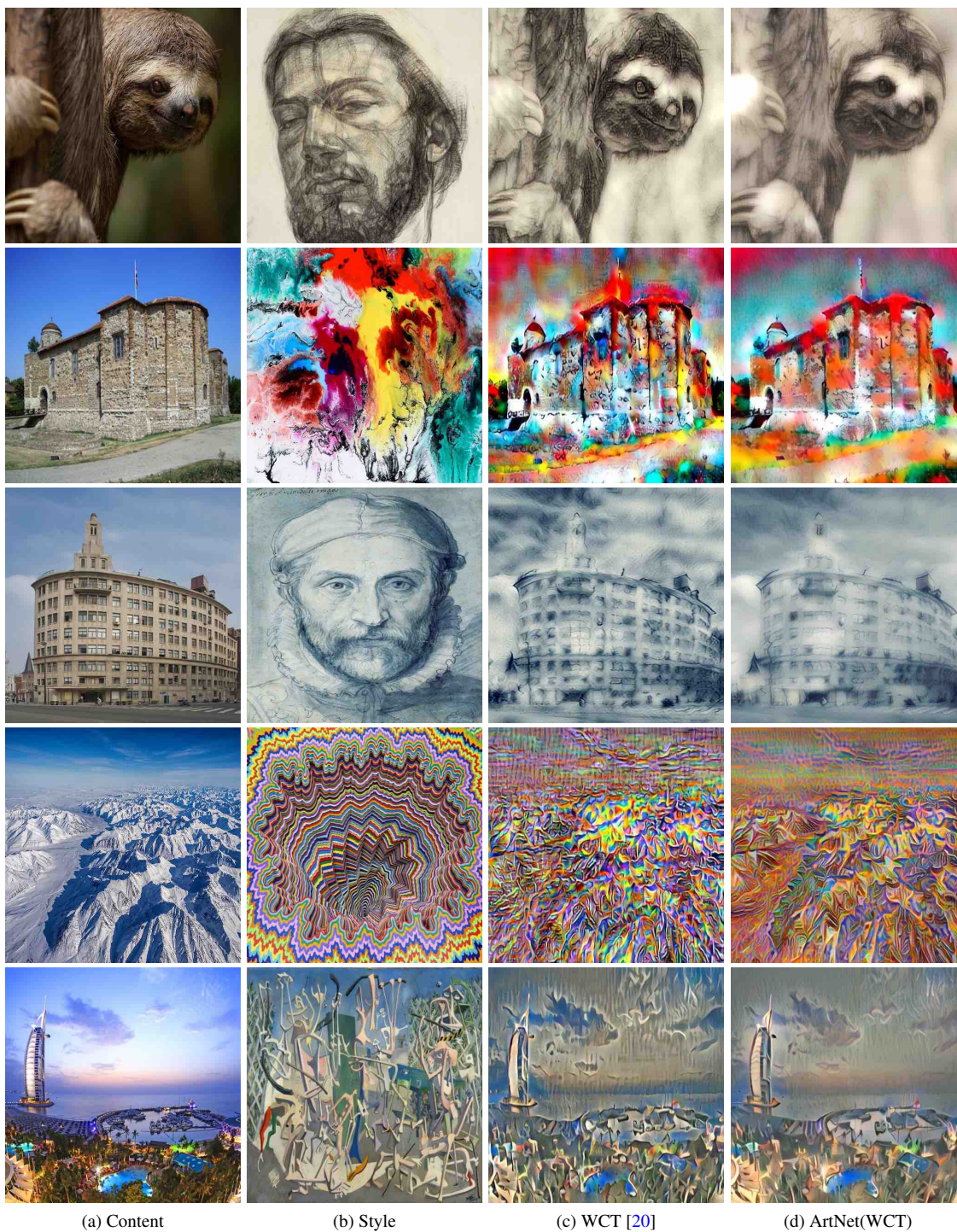


Figure 17: Artistic style transfer comparison between the ArtNet(WCT) and WCT [20].



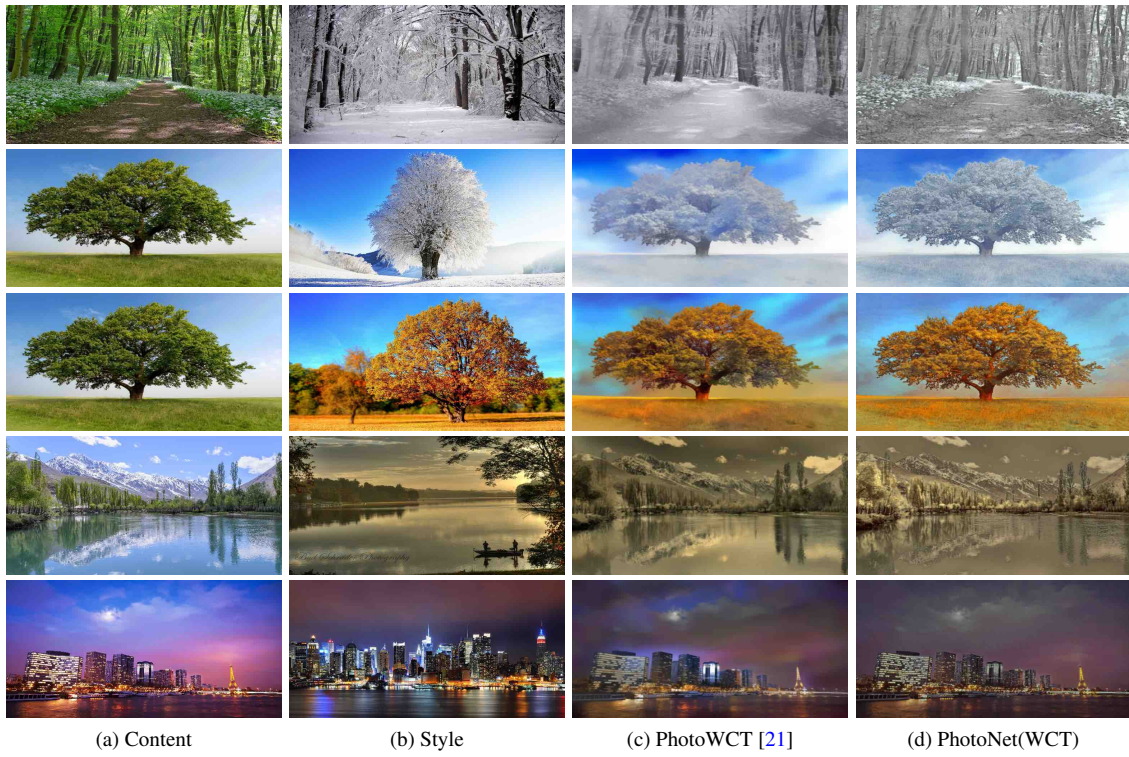


Figure 18: Photorealistic style transfer comparison between the PhotoNet(WCT) and PhotoWCT [21].

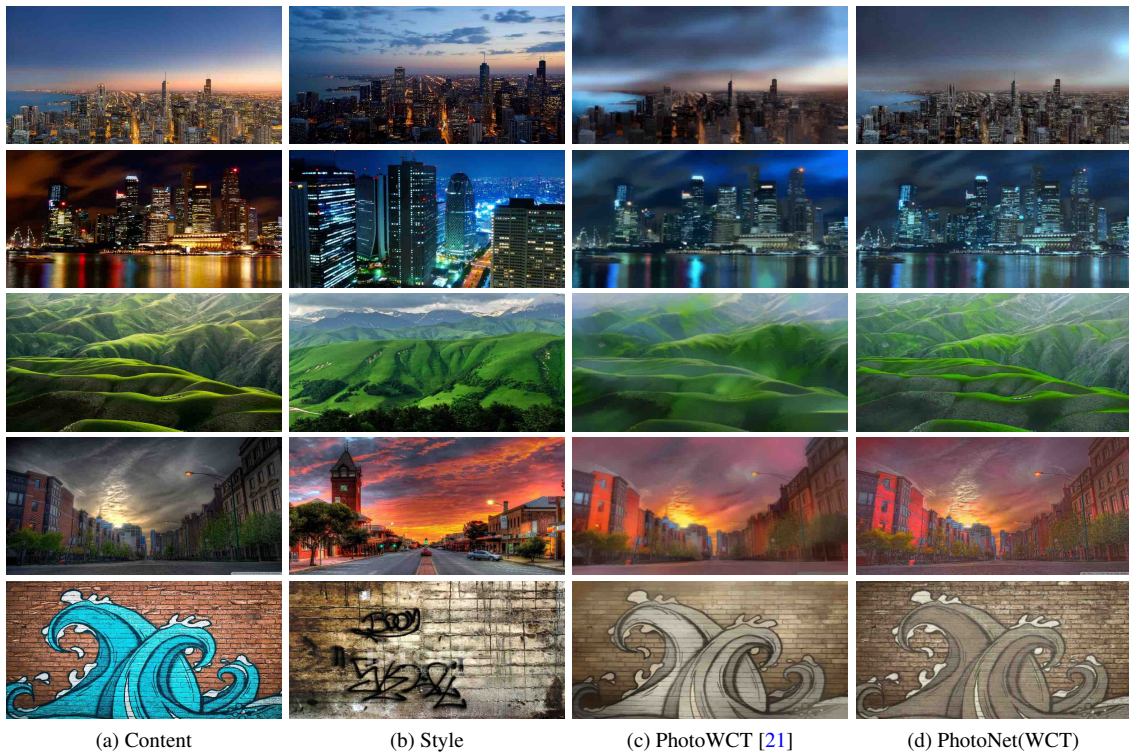


Figure 19: Photorealistic style transfer comparison between the PhotoNet(WCT) and PhotoWCT [21].



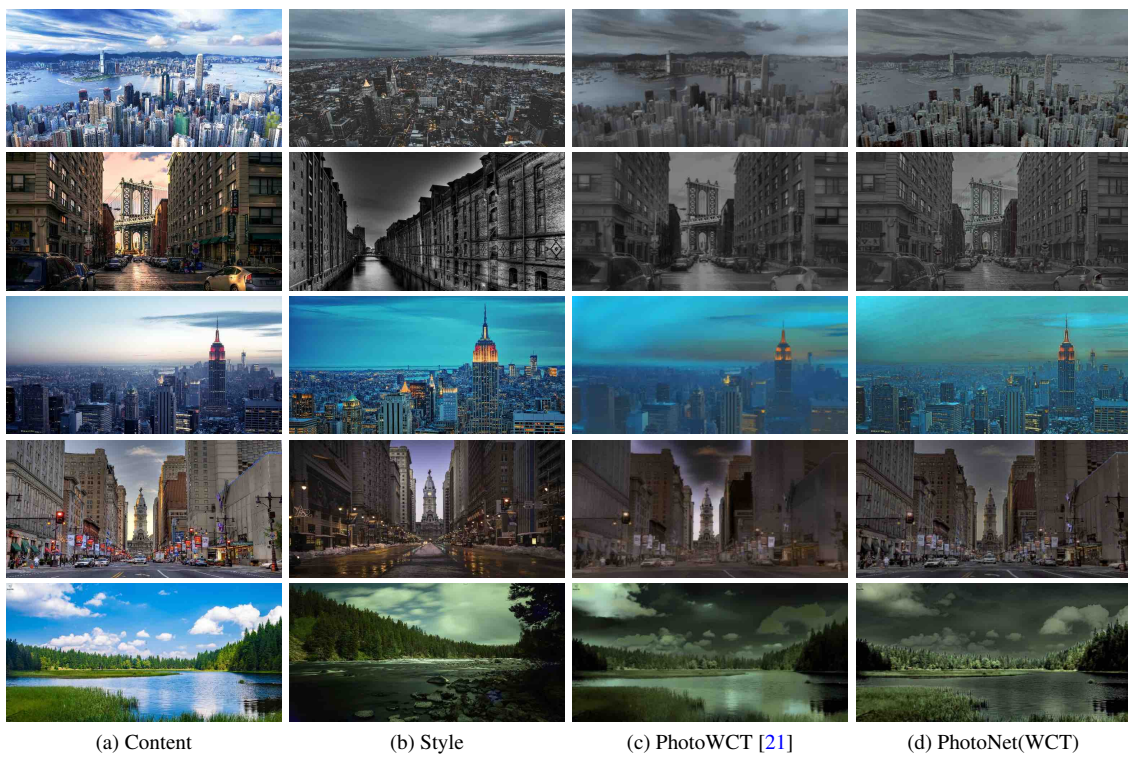


Figure 20: **Photorealistic style transfer comparison between the PhotoNet(WCT) and PhotoWCT [21].**