

Image Neural Style Transfer with Global and Local Optimization Fusion

HUI-HUANG ZHAO^{1,2}, PAUL L. ROSIN³, YU-KUN LAI³, MU-GANG LIN^{1,2}, AND QIN-YUN LIU^{1,2}

¹College of Computer Science and Technology, Hengyang Normal University, Hengyang China

²Hunan Provincial Key Laboratory of Intelligent Information Processing and Application, Hunan, China

³School of Computer Science and Informatics, Cardiff University, Cardiff, UK

Corresponding author: Hui-Huang Zhao (e-mail: happyday.huihuang@gmail.com).

This work was supported by National Natural Science Foundation of China (61503128,61772179), Science and Technology Plan Project of Hunan Province (2016TP1020), Scientific Research Fund of Hunan Provincial Education Department (16C0226,17C0223,18A333), Hengyang guided science and technology projects and Application-oriented Special Disciplines (Hengkefa [2018]60-31), Double First-Class University Project of Hunan Province (Xiangjiaotong [2018]469), Hunan Province Special Funds of Central Government for Guiding Local Science and Technology Development (2018CT5001) and Subject Group Construction Project of Hengyang Normal University (18XKQ02). We would like to thank NVIDIA for the GPU donation.

ABSTRACT This paper presents a new image synthesis method for image style transfer. For some common methods, the textures and colors in the style image are sometimes applied inappropriately to the content image, which generates artifacts. In order to improve the results, we propose a novel method based on a new strategy that combines both local and global style losses. On the one hand, a style loss function based on a local approach is used to keep the style details. On the other hand, another style loss function based on global measures is used to capture the more global structural information. Results on various images show that the proposed method reduces artifacts while faithfully transferring the style image's characteristics and preserving the structure and color of the content image.

INDEX TERMS Deep Neural Networks, Style Transfer, Markov Random Field, Gram Matrix, Local Patch.

I. INTRODUCTION

The motivation of deep learning is to build a multi-layer neural network to analyze data, with the aim of interpreting data such as images [22], [31], sounds, internet of things [23] and texts by simulating the mechanism of human brain [17], [21], [29]. Since 2016, deep learning has been applied to a new field, imitating artists' painting style [11], achieving so-called Neural Style Transfer (NST). By inputting a style image and a content image into a trained neural network model, a new image is synthesized. The newly generated image not only has the structure and content features of the original content image, but also has the style or textural features of the original style image.

Such NST has become an active research topic in the field of artificial intelligence. Its basic principle is to transfer style from the "style image" to the "content image" by using a neural network model with these two known images [13]. The purpose is to generate new images with different styles from the same content image according to the guidance of different style images. Nowadays, Neural Style Transfer is widely used to solve many problems, such as video stylization [28], texture synthesis [9], head style transfer [10], [18], [30], super-resolution [8], [16] and font style transfer [2]. In this paper,

we specifically consider the problem of image style transfer which is guided by different style loss function strategies. To achieve this, we propose a new loss function which combines a global measure and a local patch based approach. The global style loss helps avoid patch transfer errors, such as mouth, eyes and moustache being transferred into wrong places. An example is shown in figure 1. The local style loss helps better retain detailed styles. By combining both, our method better transfers styles and reduces artifacts.

II. RELATED WORK

Because neural style transfer can produce impressive results, and can generalize better across different styles than traditional non-photorealistic rendering methods [27], it has become one of the active research topics in academia and industry in recent years. Many research institutes and laboratories have conducted extensive and in-depth research on style transfer. Among them, Stanford University's research group led by Li Feifei achieved real-time style transfer of images by pre-training the network model of style images, which greatly improves the speed of style transfer and the resolution of image generation [16]. Scholars from the University of Science and Technology of China and Microsoft Research proposed a style



FIGURE 1: Comparison of style transfer results with global style loss and local loss.

library called Style Bank [5] which can be used for image style transfer. StyleBank consists of several convolutional filter banks, each of which explicitly expresses a style. The Durham University team used image style transfer to propose a new method of real-time monocular depth estimation for adaptive synthetic data [1]. The team from Princeton University, Adobe and UC Berkeley proposed a style transfer algorithm called Paired CycleGAN, which can automatically enhance and remove makeup [4]. Scholars from University of Science and Technology of China, Peking University and Microsoft Research reprocessed the depth features of stylistic images (i.e. arranging the spatial positions of feature maps) to achieve the style transfer of arbitrary images [14]. Researchers from Shanghai Jiaotong University and Microsoft proposed a generalized style transmission network model consisting of a style coder, a content coder, a mixer and a decoder to generate images with target style and content [38]. In order to meet the needs of 3D movies and AR/VR, scholars from the University of Science and Technology of China and Microsoft Research have studied the method of stereo neural style transfer and achieved satisfactory results [6]. Researchers at Tsinghua University and Cardiff University have proposed a Cartoon GAN style transfer algorithm [7]. It can generate arbitrary cartoon images using real scenes as source images where the style is learned from unpaired images from cartoon movies.

In addition, many world-renowned commercial enterprises have joined the style transfer research and its application. For example, Tencent AI along with scholars of Tsinghua University proposed a method of video style transfer using a feedforward network, and adopted a new two-frame cooperative training mechanism to achieve video style transfer [15]. Researchers at Adobe and Cornell University have proposed a method to generate realistic style transfer in various scenarios, which can achieve style transfer for images including daytime, weather, season and art [24]. Adobe, in conjunction with researchers from University of California, proposed a multimodal convolution neural network, which uses separate representations of the color and brightness channels to study hierarchical style transfer with multiple scaling losses [37]. 360 AI Lab, in conjunction with researchers from Peking

University and National University of Singapore, proposed a new meta-network model for image style transfer [33]. The meta model is only 449 KB in size and can run the image style transfer program in real time on mobile devices. SenseTime and researchers from the Chinese University of Hong Kong proposed a multi-scale zero-shot style transfer method based on feature decoration [34]. A style decorator was designed to make use of semantic alignment style features from arbitrary style images to form content features. This not only matches their feature distribution on the whole, but also retains the detailed style patterns in decorative features. The partition algorithm is a significant digital image processing technique for image denoising and image reconstruction. [32] proposed a new image denoising method based on Hard-partition Weighted Sum filters. In order to solve the multi-person pose estimation problem, [25] proposed a novel Pose Partition Network (PPN) which has a good performance in low complexity and high accuracy of joint detection and partition. [3] proposed an Adaptive Triangular Partition Algorithm named IATP for digital images. The method considers the grayscale distribution of the image and removes the shared edges between the adjacent triangles in the partitioned mesh.

Two main types of methods for representing elements of an image are used in deep learning based style transfer: global approaches based on the Gram matrix [11] or other global measures (e.g. histograms [26]) and local approaches based on patch matching [18], [19]. Compared to the global methods, methods based on patch matching are more flexible and better cope with cases in which the visual styles or elements vary across the image. However, they could also produce visible artifacts when there are local matching errors. Compared to the methods based on local approaches, the structure and color of the content image can be preserved better with global approaches, although detailed styles may not be fully captured. An example is shown in figure 2, in which artifacts produced by existing methods are clearly evident.

III. ARCHITECTURE

We now discuss our style transfer DCNN (Deep Convolutional Neural Network) architecture which is based on the VGG 19-

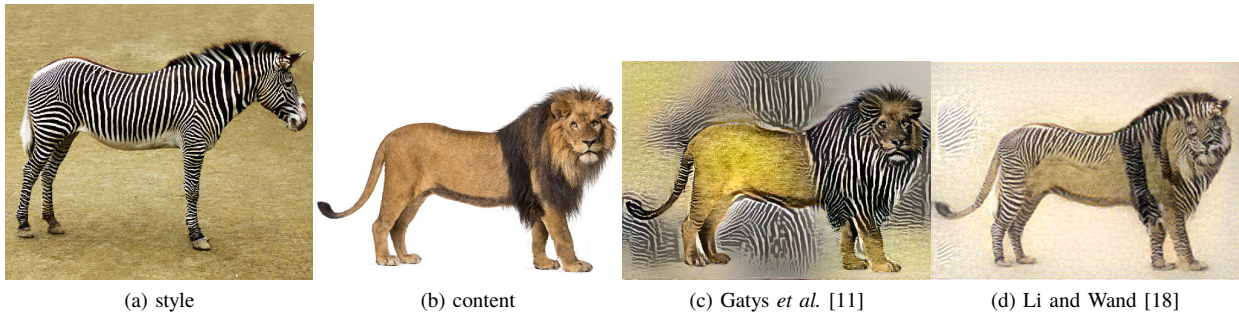


FIGURE 2: Comparison of style transfer results with different methods.

layer network [35]. It takes as input a content image and a style image, both of which are fed into the VGG 19-layer network. The DCNN architecture combines pooling and convolution layers with 3×3 filters. Our style transfer network consist of two DCNN models which are extracted from specific layers from the VGG 19-layer network. The style transfer network is shown in figure 3.

In the proposed architecture in figure 3, the global style loss network and local style loss network consist of different layers of the VGG 19-layer network. We set them to have equal contribution in the overall loss. Following the patch-based approach of [18], we synthesize at multiple increasing resolutions, and randomly initialize the optimization.

IV. STYLE TRANSFER OPTIMIZATION FUNCTION

Next, we introduce our style transfer model. We use a fusion loss function which is based on a patch-based approach [18] and a Gram matrix method [11] for style transfer, using optimization to minimize content reconstruction loss $L_{content}$, and style losses L_{global} and L_{local} . We are given a style image $x_s \in \mathbb{R}^{3 \times w_s \times h_s}$, a content image $x_c \in \mathbb{R}^{3 \times w_c \times h_c}$, where w_s and h_s are the width and height of the style image, and w_c and h_c are the width and height of the content image. The style transfer result image is denoted by $x \in \mathbb{R}^{3 \times w_c \times h_c}$. We define a loss function as follows, and seek x that minimizes it:

$$L(x) = \alpha_1 L_{global}(F(x), F(x_s)) + \alpha_2 L_{local}(F(x), F(x_s)) + \alpha_3 L_{content}(F(x), F(x_c)) \quad (1)$$

where L_{global} , L_{local} and $L_{content}$ are defined as global and local style loss functions and content loss function respectively, where $F(x)$ is x 's feature map (activation) that the network outputs in some layer, $F(x_s)$ is the feature map of the style image x_s in the same layer. For our method, L_{global} and L_{local} aim to penalize inconsistencies in neural activations between x and x_s from global and local perspectives. $L_{content}$ computes the squared distance between the feature map of the synthesized image and that of the content source image x_c .

A. GLOBAL STYLE LOSS

Gatys et al. found correlations between feature maps can be used to match textures between images, and incorporated this

in their image style transfer methods [11], [12]. A Gram matrix between a feature map $F_i^l(x)$ and $F_j^l(x)$ at layer l is defined as

$$G_{i,j}^l(x) = \langle F_i^l(x), F_j^l(x) \rangle \quad (2)$$

So the global style loss function is defined as

$$L_{global}(x, x_s) = \sum \| G(x) - G(x_s) \|^2 \quad (3)$$

Global feature and local feature are two main categories. The main difference between Global Style Loss and Local Style Loss is that the Global Style Loss methods use the feature maps in the neural network model as the optimization goal, while the Local Style Loss methods divide each feature map into blocks for optimization.

B. LOCAL STYLE LOSS FUNCTION:

Li and Wand [18] proposed a method combining MRF (Markov random fields) and DCNN for image synthesis. Our local style loss function similarly combines an MRF and the VGG 19-layer network model. We extract all the local patches from $F(x)$, denoted as $\Psi(F(x))$. For a given layer, assuming N is the number of channels, each patch in $\Psi_i(F(x))$ has size $t \times t \times N$, where t is the width and height of the patch. The local style loss function E_s is defined as

$$L_{local}(F(x), F(x_s)) = \sum_{i=1}^P \| \Psi_i(F(x)) - \Psi_{NN(i)}(F(x_s)) \|^2 \quad (4)$$

where P is the number of patches in the synthesized image. For each patch $\Psi_i(F(x))$, we find its best matching patch $\Psi_{NN(i)}(F(x_s))$ using a normalized cross-correlation over all P_s example patches in $\Psi(F(x_s))$:

$$NN(i) := \arg \max_{j=1, \dots, P_s} \frac{\Psi_i(F(x)) \cdot \Psi_j(F(x_s))}{|\Psi_i(F(x))| \cdot |\Psi_j(F(x_s))|}, \quad (5)$$

where $\Psi_i(\Phi(x))$ is the concatenation of neural activation for the j^{th} patch of the style image. The nearest patch thus takes style similarity of the style and the output image into account.

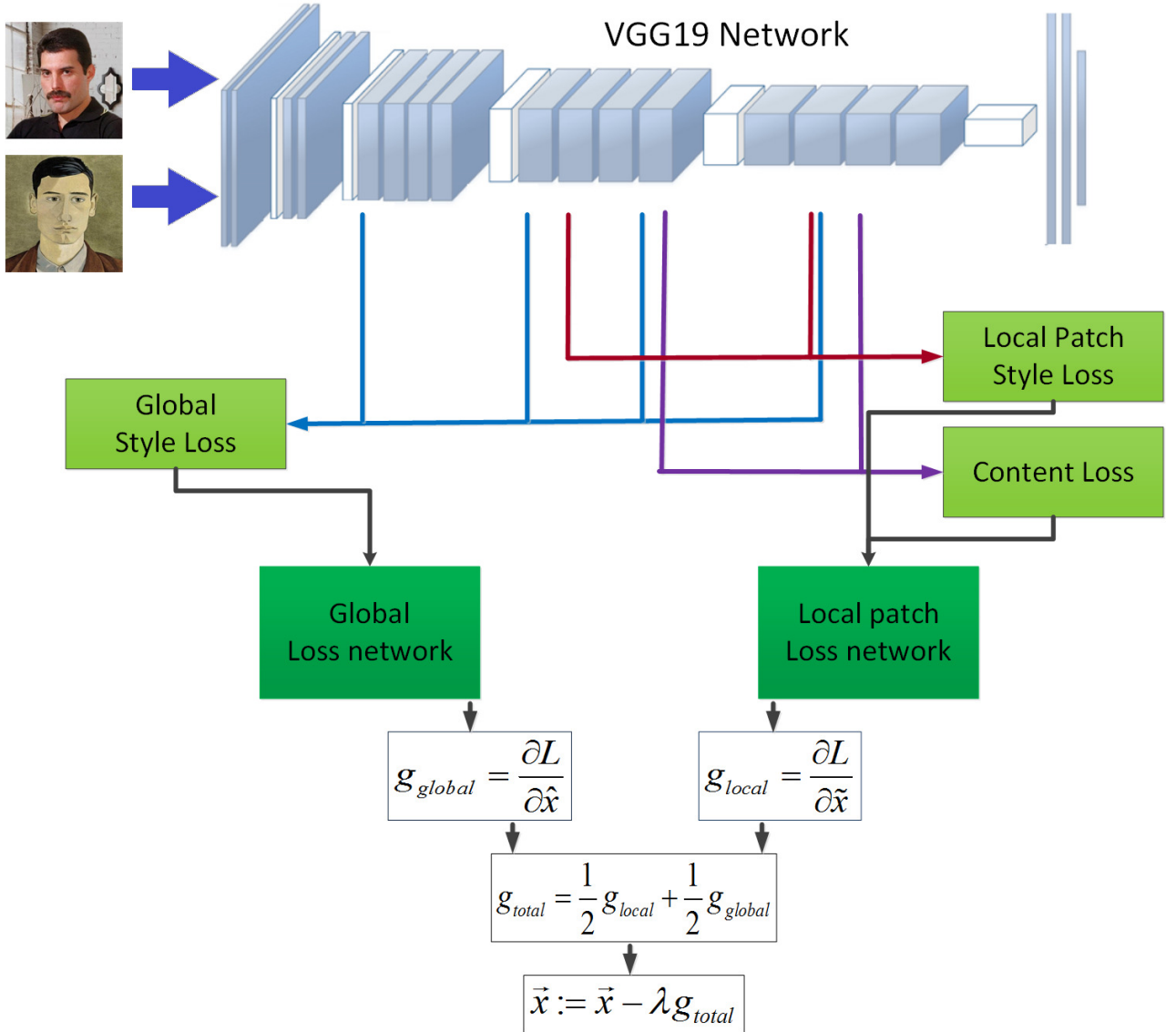


FIGURE 3: Our neural style transfer framework.

C. CONTENT LOSS FUNCTION:

In order to control the content of the synthesized image, we define $L_{content}$ as the squared Euclidean distance between $F(x)$ and $F(x_c)$ at layer l :

$$L_{content}(F(x), F(x_c)) = \sum_{i=1}^{N_l} \| (F(x) - F(x_c)) \|^2. \quad (6)$$

Like most methods [11], [18], we also minimize Equation 1 using backpropagation with L-BFGS. In Equation 1, α_1 , α_2 and α_3 are weights for the global style, local patch style and the content constraints, respectively. During experiments, we set $\alpha_1 = 20$, $\alpha_2 = 10^{-4}$ and $\alpha_3 = 20$, and these values can be fine-tuned to interpolate between the content and the global and local style preservation.

V. EXPERIMENTS AND DISCUSSION

In common with previous research, we use the pre-trained VGG 19-layer network to generate feature maps. During the global loss part, layers *relu1_2*, *relu2_2*, *relu3_3* and *relu4_2* are selected. During the local loss part, layer *relu3_1* is selected. We use 3×3 patches, and we set the stride to one. Layer *relu3_1* is selected in the content loss part. On a GTX Titan with 12GB of GPU RAM, synthesis takes from 5 to 10 minutes depending on the desired output quality and resolution.

We will now compare the proposed method with several popular methods: [11], [18], [20], [36] which are representative global and local neural style transfer methods. Some examples of style transfer for male and female portraits are shown in figure 4.

One can see from figure 4 that the proposed method not only achieves better results in style details and avoids mistakes

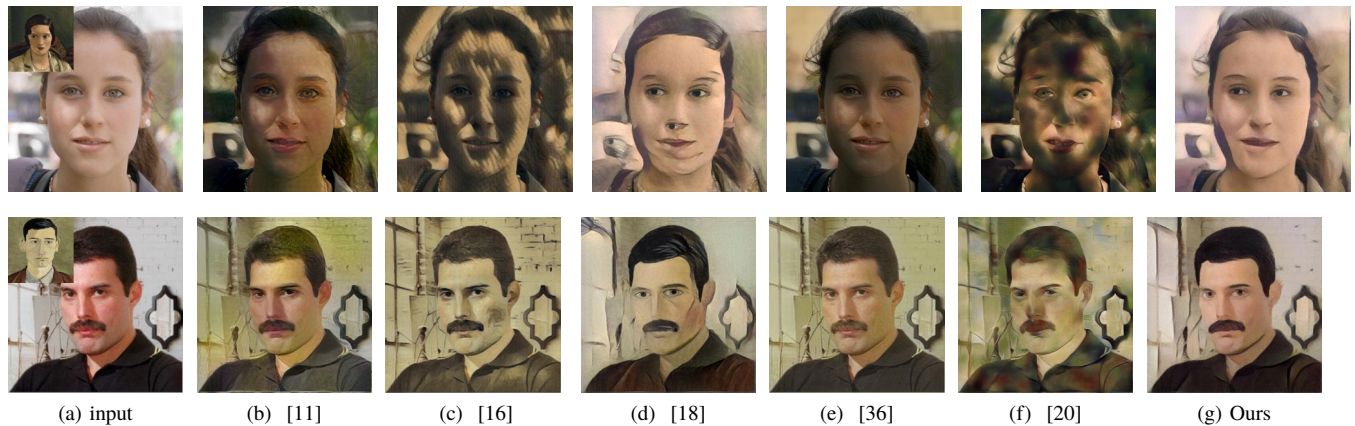


FIGURE 4: Style transfer comparison.

in transfer, but also keeps the background clean without being ruined.

A. STYLE TRANSFER OF DIFFERENT OBJECT TYPES.

More examples of style transfer for objects like tiger, lion and cat are shown in figure 5 – figure 7. The styles are zebras with different poses. One can see from that the proposed method can transfer more zebra stripe style to the tiger, lion and cat while keeping the structure of the content images.

B. MODIFYING THE CONTENT AND STYLE WEIGHT.

Further experiments are carried out, in which α_2 is modified while $\alpha_1 = 20$ and $\alpha_3 = 20$ are fixed. Figure 8 demonstrates the effect of modifying α_2 from $\alpha_2 = 10^{-1}$ to $\alpha_2 = 10^{-6}$.

When α_2 is too small, the result does not have sufficient style information. On the other hand, setting α_2 too large may result in matched patches having poor content consistency. According to our experiments, $\alpha_2 \in [10^{-3}, 10^{-4}]$ and both α_1 and $\alpha_3 \in [20, 40]$ achieve the best results.

VI. CONCLUSIONS

Our paper demonstrates the benefits of fusing global and local losses in image style transfer. A fusion architecture is designed. On the one hand, a style loss function based on a local approach defined in several layers is used to keep the detailed styles. On the other hand, a style loss function based on global measures defined in several layers is used to capture the more global information. Results on various images show the proposed method can preserve the structure and color of the content image while having the style transferred, reducing artifacts.

REFERENCES

- [1] Atapour-Abarghouei, A., Breckon, T.P.: Real-time monocular depth estimation using synthetic data with domain adaptation via image style transfer. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
- [2] Azadi, S., Fisher, M., Kim, V., Wang, Z., Shechtman, E., Darrell, T.: Multi-content GAN for few-shot font style transfer (2017)
- [3] Cai, X.Y.Z.: An adaptive triangular partition algorithm for digital images. IEEE Transactions on Multimedia **21**(6), 1372 – 1383 (2018)
- [4] Chang, H., Lu, J., Yu, F., Finkelstein, A.: PairedCycleGAN: Asymmetric style transfer for applying and removing makeup. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
- [5] Chen, D., Yuan, L., Liao, J., Yu, N., Hua, G.: StyleBank: an explicit representation for neural image style transfer. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
- [6] Chen, D., Yuan, L., Liao, J., Yu, N., Hua, G.: Stereoscopic neural style transfer. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
- [7] Chen, Y., Lai, Y.K., Liu, Y.J.: CartoonGAN: Generative adversarial networks for photo cartoonization. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
- [8] Deng, X.: Enhancing image quality via style transfer for single image super-resolution. IEEE Signal Processing Letters (2018)
- [9] Efros, A., Freeman, W.: Image quilting for texture synthesis. In: ACM SIGGRAPH (2001)
- [10] Fišer, J., Jamriška, O., Simons, D., Shechtman, E., Lu, J., Asente, P., Lukáč, M., Sýkora, D.: Example-based synthesis of stylized facial animations. ACM Transactions on Graphics **36**(4) (2017)
- [11] Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2414–2423 (2016)
- [12] Gatys, L.A., Ecker, A.S., Bethge, M., Hertzmann, A., Shechtman, E.: Controlling perceptual factors in neural style transfer. arXiv preprint arXiv:1611.07865 (2016)
- [13] Gatys, L.A., Ecker, A.S., Bethge, M., Hertzmann, A., Shechtman, E.: Controlling perceptual factors in neural style transfer. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
- [14] Gu, S., Chen, C., Liao, J., Yuan, L.: Arbitrary style transfer with deep feature reshuffle. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
- [15] Huang, H., Wang, H., Luo, W., Ma, L., Jiang, W., Zhu, X., Li, Z., Liu, W.: Real-time neural style transfer for videos. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
- [16] Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: European Conference on Computer Vision, pp. 694–711. Springer (2016)
- [17] Lecun, Y., Bengio, Y., Hinton, G.: Deep learning. Nature **521**(7553), 436 (2015)
- [18] Li, C., Wand, M.: Combining Markov random fields and convolutional neural networks for image synthesis. In: Proc. Conf. Computer Vision and Pattern Recognition, pp. 2479–2486 (2016)
- [19] Li, C., Wand, M.: Precomputed real-time texture synthesis with Markovian generative adversarial networks. In: European Conference on Computer Vision (2016)
- [20] Li, Y., Fang, C., Yang, J., Wang, Z., Lu, X., Yang, M.H.: Universal style transfer via feature transforms. In: Advances in Neural Information Processing Systems, pp. 386–396 (2017)
- [21] Lu, H., Li, Y., Chen, M., Kim, H., Serikawa, S.: Brain intelligence: go beyond artificial intelligence. Mobile Networks and Applications **23**(2), 368–375 (2018)

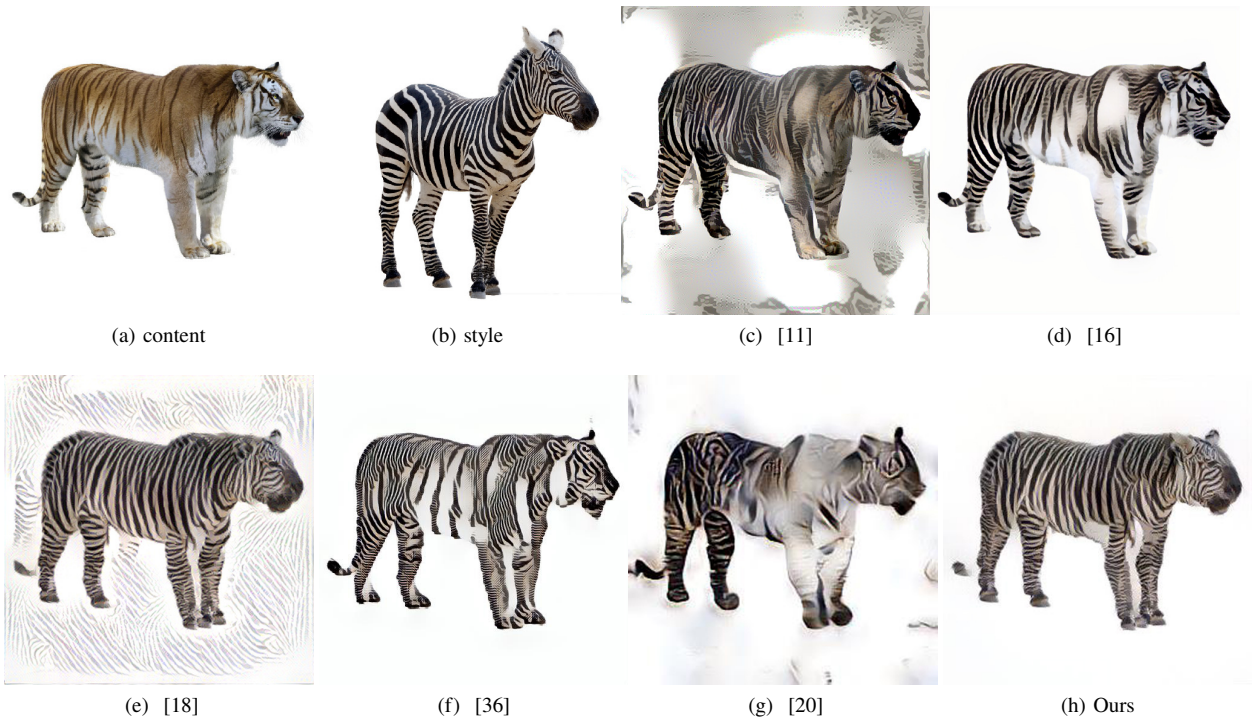


FIGURE 5: More objects style transfer; zebra to tiger.

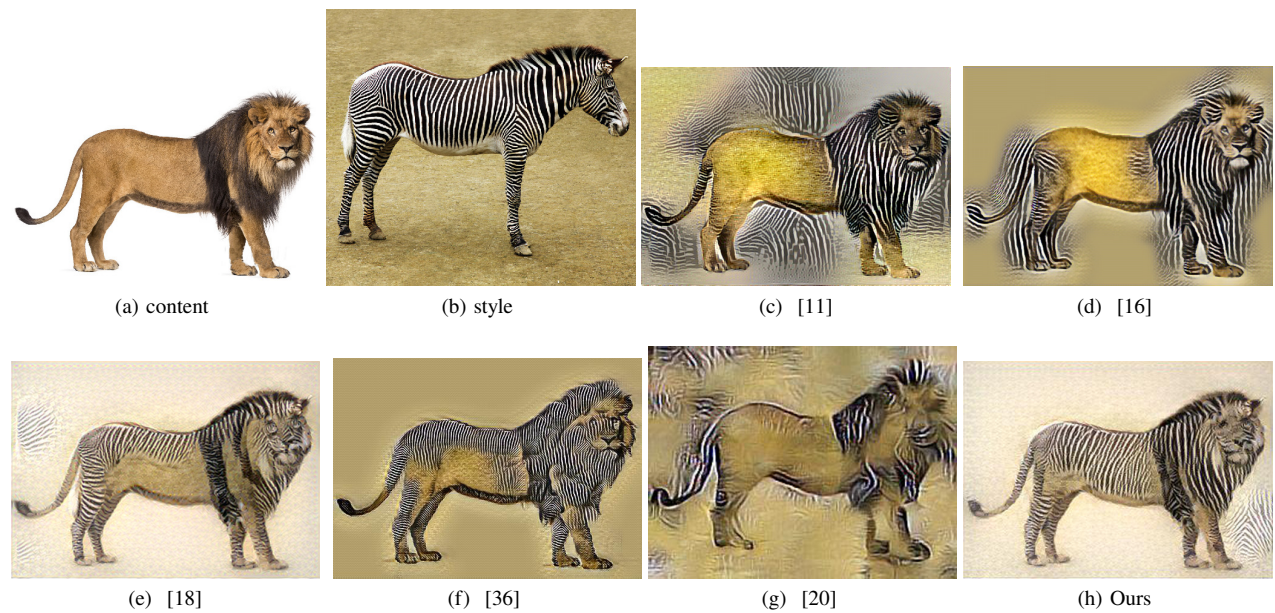


FIGURE 6: More objects style transfer; zebra to lion.

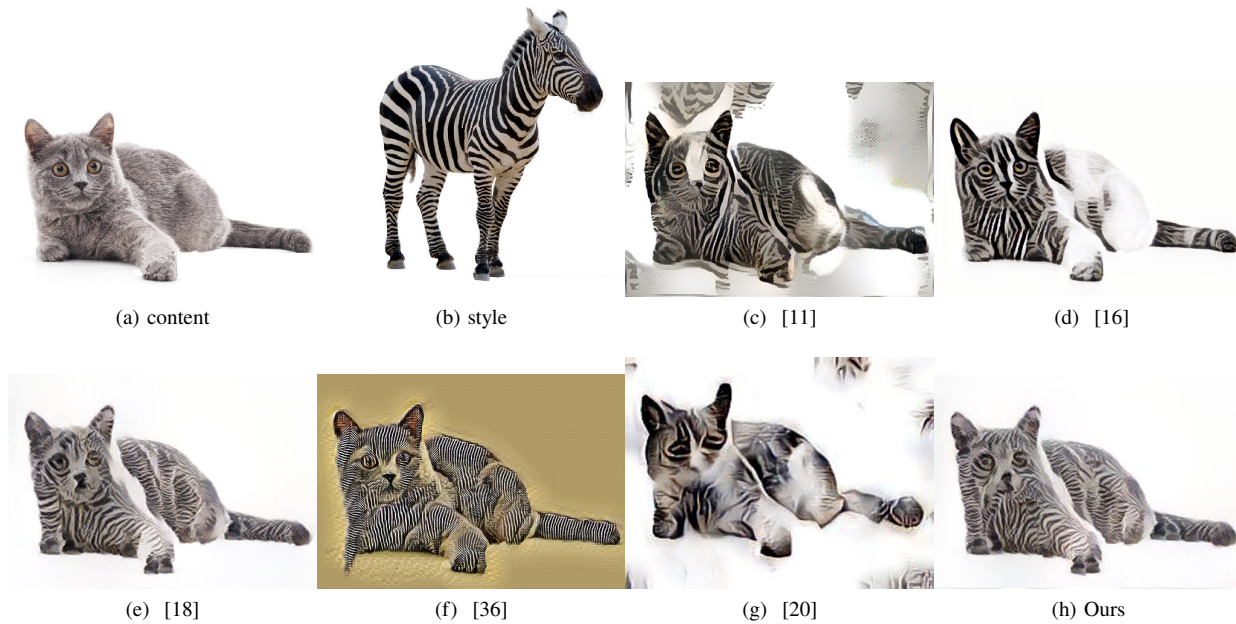
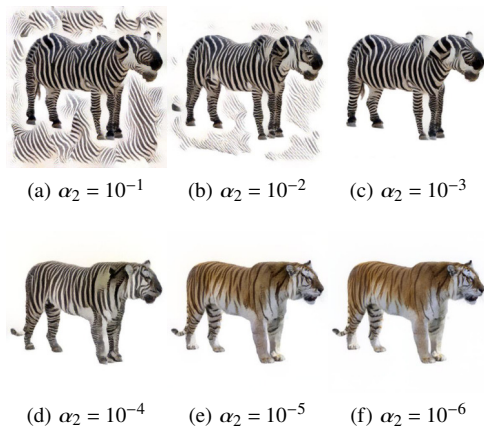


FIGURE 7: More objects style transfer; zebra to cat.

FIGURE 8: Result showing the effects of varying parameter α_2 .

- [22] Lu, H., Li, Y., Uemura, T., Kim, H., Serikawa, S.: Low illumination underwater light field images reconstruction using deep convolutional neural networks. *Future Generation Computer Systems* **82**, 142–148 (2018)
- [23] Lu, H., Wang, D., Li, Y., Li, J., Li, X., Kim, H., Serikawa, S., Humar, I.: CONet: a cognitive ocean network. *arXiv preprint arXiv:1901.06253* (2019)
- [24] Luan, F., Paris, S., Shechtman, E., Bala, K.: Deep photo style transfer. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017)
- [25] Nie, X., Feng, J., Xing, J., Yan, S.: Pose partition networks for multi-person pose estimation. In: *The European Conference on Computer Vision (ECCV)* (2018)
- [26] Risser, E., Wilmot, P., Barnes, C.: Stable and controllable neural texture synthesis and style transfer using histogram losses. *arXiv preprint arXiv:1701.08893* (2017)
- [27] Rosin, P.L., Collomosse, J.P. (eds.): *Image and Video-Based Artistic Stylisation*. Springer (2013)

- [28] Ruder, M., Dosovitskiy, A., Brox, T.: Artistic style transfer for videos and spherical images. *International Journal of Computer Vision* **126**(11), 1199–1219 (2018)
- [29] Schmidhuber, J.: Deep learning in neural networks: An overview. *Neural Networks* **61**, 85–117 (2015)
- [30] Selim, A., Elgharib, M., Doyle, L.: Painting style transfer for head portraits using convolutional neural networks. *ACM Transactions on Graphics* **35**(4), 1–18 (2016)
- [31] Serikawa, S., Lu, H.: Underwater image dehazing using joint trilateral filter. *Computers & Electrical Engineering* **40**(1), 41–50 (2014)
- [32] Shao M, B.K.E.: Optimization of partition-based weighted sum filters and their application to image denoising. *IEEE Transactions on Image Processing* **15**(7), 1900–1915 (2006)
- [33] Shen, F., Yan, S., Zeng, G.: Neural style transfer via meta networks. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)
- [34] Sheng, L., Lin, Z., Shao, J., Wang, X.: Avatar-Net: multi-scale zero-shot style transfer by feature decoration. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)
- [35] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
- [36] Ulyanov, D., Vedaldi, A., Lempitsky, V.: 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*, p. 6 (2017)
- [37] Wang, X., Oxholm, G., Zhang, D., Wang, Y.F.: Multimodal transfer: A hierarchical deep convolutional neural network for fast artistic style transfer. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017)
- [38] Zhang, Y., Zhang, Y., Cai, W.: Separating style and content for generalized style transfer. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)



HUI-HUANG ZHAO received his Ph.D. degree in 2010 from XiDian University. He was a Sponsored Researcher in the School of Computer Science and Informatics, Cardiff University. Now he is an Associate Professor in the College of Computer Science and Technology, Hengyang Normal University. His main research interests include Solder Joint Inspection, Compressive Sensing, Machine Learning, and Image Processing.



PAUL L. ROSIN received the B.Sc. degree in Computer Science and Microprocessor Systems in 1984 from Strathclyde University, Glasgow, and the Ph.D. degree in Information Engineering from City University, London in 1988. He is a full professor in the School of Computer Science and Informatics, Cardiff University. His main research interests include Non-Photorealistic Rendering, Mesh Processing, and Computer Vision.



YU-KUN LAI received his bachelor's and Ph.D. degrees in computer science from Tsinghua University in 2003 and 2008, respectively. He is currently a Reader of Visual Computing in the School of Computer Science and Informatics, Cardiff University. His research interests include computer graphics, geometry processing, image processing and computer vision. He is on the editorial board of *The Visual Computer*.



MU-GANG LIN received the M.S. degree in computer application technology from Changsha University of Science & Technology, P.R. China in 2005. He is currently pursuing the Ph.D. degree in computer science with the School of Computer Science and Engineering, Central South University, Changsha, P.R. China. He is working as an associate professor in the College of Computer Science and Technology, Hengyang Normal University, Hengyang, P.R. China. His current research interests include Computer Algorithms and Parameterized Algorithms.



QING-YUN LIU received the M.S. degree in computer application technology from Jiangsu University, P.R. China in 2014. He worked in National Defense University of Science and Technology for 2 years as a supercomputer systems engineer. He is working as an assistant professor in the College of Computer Science and Technology, Hengyang Normal University, Hengyang, P.R. China. His current research interests include Image Style Transfer and Deep Learning.

...