

Stroke Controllable Fast Style Transfer with Adaptive Receptive Fields*

Yongcheng Jing[†] Yang Liu[†] Yezhou Yang[‡] Zunlei Feng[†] Yizhou Yu[†] Mingli Song^{††}

[†] Microsoft Visual Perception Laboratory, College of Computer Science and Technology,
Zhejiang University

[‡] School of Computing, Informatics, and Decision Systems Engineering,
Arizona State University

{ycjing, lyng_95, zunleifeng, brooksong}@zju.edu.cn yz.yang@asu.edu yizhouy@acm.org

Abstract

Recently, in the community of Neural Style Transfer, several algorithms are proposed to transfer an artistic style in real-time, which is known as Fast Style Transfer. However, controlling the stroke size in stylized results still remains an open challenge. To achieve controllable stroke sizes, several attempts were made including training multiple models and resizing the input image in a variety of scales, respectively. However, their results are not promising regarding the efficiency and quality. In this paper, we present a stroke controllable style transfer network that incorporates different stroke sizes into one single model. Firstly, by analyzing the factors that influence the stroke size, we adopt the idea that both the receptive field and the style image scale should be taken into consideration for most cases. Then we propose a StrokePyramid module to endow the network with adaptive receptive fields, and two training strategies to achieve faster convergence and augment new stroke sizes upon a trained model respectively. Finally, by combining the proposed runtime control techniques, our network can produce distinct stroke sizes in different output images or different spatial regions within the same output image. The experimental results demonstrate that with almost the same number of parameters as the previous Fast Style Transfer algorithm, our network can transfer an artistic style in a stroke controllable manner.

1. Introduction

Rendering a photograph with a given artwork style has been a long-standing research topic [12, 30, 28, 13]. Conventionally, the task of style transfer is usually studied as

*Another version of the paper with high-resolution results can be found in the project website: <http://yongchengjing.com/StrokeControllable>

^{††}Mingli Song is the corresponding author. Yongcheng Jing and Yang Liu contributed equally.

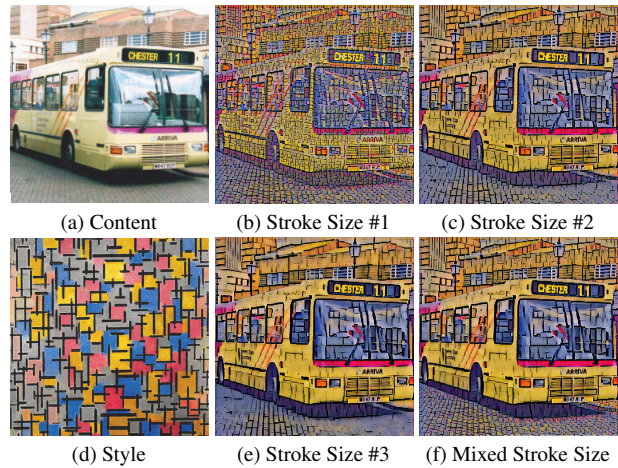


Figure 1. Stylized results with different stroke sizes. All these results are produced by one single model in real-time using our proposed algorithm. We did not do any image pre- or post-processing before and after forwarding.

a generalization of texture synthesis [6, 8, 7]. Based on the recent progress in visual texture modelling [9], Gatys *et al.* firstly propose an algorithm that exploits Convolutional Neural Network (CNN) to recombine the content of a given photograph and the style of an artwork, and reconstruct a visually plausible stylized image, which is known as the process of Neural Style Transfer [10]. Since the seminal work of Gatys *et al.*, Neural Style Transfer has been attracting wide attention from both academia and industry [20, 24, 27]. However, the algorithm of Gatys *et al.* is based on iterative image optimizations and it requires a slow optimization process for each pair of content and style. To tackle this issue, several algorithms are proposed to speed up the style transfer process, which is called Fast Style Transfer in the literature [11, 26].

Currently, there are three categories of Fast Style Transfer [15], namely Per-Style-Per-Model (PSPM) [31, 32, 16, 21], Multiple-Style-Per-Model (MSPM) [5, 36, 22, 2] and

Arbitrary-Style-Per-Model (ASPM) [14, 23]. The gist of the PSPM is to train a feed-forward style-specific generator and produces a corresponding stylized result with a forward pass. The MSPM improves the efficiency by further incorporating multiple styles into one single generator. The ASPM tries to transfer an arbitrary style through only one single model.

For these Fast Style Transfer algorithms, there is a trade-off between efficiency and quality [14, 23]. In terms of the quality, the PSPM is usually regarded to produce more appealing stylized results [32, 14]. However, the PSPM is not flexible in terms of controlling perceptual factors (*e.g.*, style-content tradeoff, color control, spatial control). Among these perceptual factors, strokes are one of the most important geometric primitives to characterize an artwork, as is shown in Figure 1. In reality, for the same texture, different artists have their own way to place different sizes of strokes as a reflection of their unique “styles” (*e.g.*, Monet and Pollock). To achieve different stroke sizes with the PSPM, one possible solution is to train multiple models, which is time and space consuming. Another solution is to resize the input image to different scales, which will inevitably hurt the quality of stylization.

In this paper, we propose a stroke controllable Fast Style Transfer algorithm that can incorporate multiple stroke sizes into one single model. Firstly, by analyzing the factors that influence the stroke size in stylized results, we adopt the idea that both the receptive field and the style image scale should be considered for most cases. Based on this idea, we propose a *StrokePyramid* module to endow the network with adaptive receptive fields and different stroke sizes are learned with different receptive fields. Then a progressive training strategy is introduced to make the network converge faster, and an incremental training strategy is presented to learn new stroke sizes upon a trained model. Finally, by combining two proposed runtime control techniques which are stroke interpolation and spatial stroke control, our network can produce distinct stroke sizes in different outputs or different spatial regions within the same output image.

In summary, the contributions of this work are threefold:

- We analyze the factors that influence the stroke size in stylized results. Based on our analysis, we propose that both the receptive field and the style image scale should be considered for stroke size control in most cases.
- We propose a stroke controllable style transfer network and two corresponding training strategies in order to achieve faster convergence and augment new stroke sizes upon a trained model respectively. To the best of our knowledge, this is the first style transfer network that can transfer an artistic style in a stroke controllable manner.

- We present two runtime control techniques to empower the network with the ability of producing more diverse stroke sizes in different output images and distinct stroke sizes in different spatial regions within the same output image.

2. Related Work

Controlling perceptual factors in Fast Style Transfer.

Stroke size control belongs to the domain of controlling perceptual factors during stylization. In this field, several significant works are recently presented. However, there are few efforts devoted to controlling stroke size during Fast Style Transfer. In [11], Gatys *et al.* study the color control and spatial control for Fast Style Transfer. Lu *et al.* further extend Gatys *et al.*’s work to meaningful spatial control by incorporating semantic content, achieving the so-called Fast Semantic Style Transfer [26]. Another related work is Wang *et al.*’s algorithm which aims to learn large brush strokes for high-resolution images [33]. They find that current Fast Style Transfer algorithms fail to paint large strokes in high-resolution images and propose a coarse-to-fine architecture to solve this problem. Note that the work in [33] is intrinsically different from this paper as one single pre-trained model in [33] still produces one stroke size for the same input image.

Regulating receptive field in neural network.

The receptive field is one of the basic concepts in convolutional neural networks, which refers to a region of the input image that one neuron is responsive to. It can affect the performance of the networks and becomes a critical issue in many tasks (*e.g.*, semantic segmentation, image parsing). To regulate the receptive field, [35] proposes the operation of dilated convolution (also called atrous convolution in [3]), which supports the expansion of receptive field by setting different dilation values. Another work in [4] further proposes a deformable convolution which augments the sampling locations in regular convolution with additional offsets. Furthermore, Wei *et al.* [34] propose a learning-based receptive field regulating method which is to inflate or shrink feature maps automatically on the basis of the learned knowledge.

3. Pre-analysis

Before our pre-analysis, we model the concept of the stroke size firstly. Consider an image in style transfer as a composition of a series of small stroke textons, which are referred as the fundamental geometric micro-structures in images [17, 37]. The stroke size of an image can be defined as the average scale of the composed stroke textons.

In the deep neural network based Fast Style Transfer, three factors are found to influence the stroke size, namely

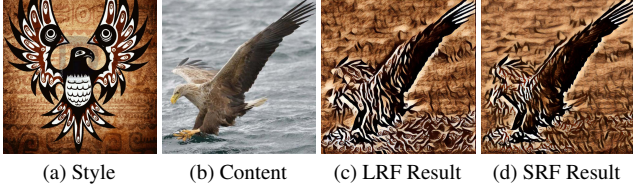


Figure 2. Results of learning the same size of large strokes with large and small receptive fields, respectively. LRF represents the result produced with a large receptive field and SRF represents the result produced with a small receptive field. The content image is credited to flickr user *Kevin Robson*.

the scale of the style image [33], the receptive field in the loss network [11], and the receptive field in the generative network.

The objective style is usually learned by matching the style image’s gram-based statistics [10] in style transfer algorithms, which are computed over the feature maps from the pre-trained VGG network [29]. These gram-based statistics are scale-sensitive, *i.e.*, they contain the scale information of the given style image. One reason for this characteristic is that the VGG features vary with the image scale. We also find that for other style statistics (*e.g.*, BN-based statistics in [24]), it reaches the same conclusion. Therefore, given the same content image, generative networks trained with different scales of the style image can produce different stroke sizes.

Although the stroke in stylized results generally becomes larger with the increase of the style image scale, this is infeasible when the style image is scaled to a high resolution (*e.g.*, 3000×3000 pixels [11]). The reason for this problem is that a neuron in pre-trained VGG loss network can only affect a region with the receptive field size in the input image. When the stroke texton is much larger than the fixed receptive field in VGG loss network, there is no visual difference between a large and larger stroke texton in a relatively small region.

Apart from these above two factors, we further find that the receptive field size in the generative network also has influence on the stroke size. In Figure 2, we change the receptive field size in the generative network and other factors remain the same. It is noticeable that a larger stroke size is produced with a larger receptive field for some styles. To explain this result, we interpret the training process of the generative network as teaching the convolutional kernels to paint a pre-defined size of stroke textons in each region with the size of receptive field. Therefore, given two different sizes of input images, the kernels of a trained network paint almost the same size of stroke textons in the same size of regions, as shown in Figure 3. In particular, when the receptive field in a generative network is smaller than the stroke texton, the kernels can only learn to paint a part of the whole stroke texton in each region, which influences the

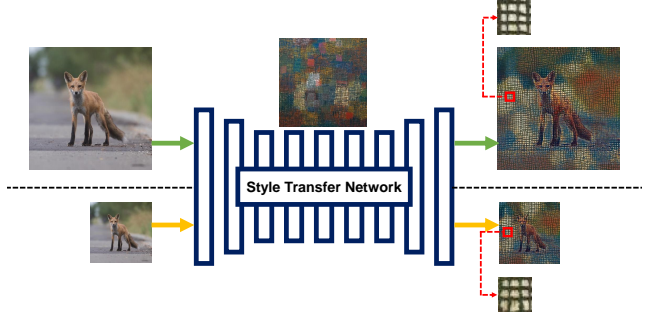


Figure 3. The feed-forward process of Fast Style Transfer. For the same size of regions in the outputs of both small and large input images respectively, their stroke sizes are almost the same. The content image is credited to flickr user *BillChenSF*.

stroke sizes.

To sum up, both the scale of the style image and the receptive field in the generative network should generally be considered for stroke size control. As the style image is not high-resolution in most cases, the influence of the receptive field in the loss network is not considered in this work.

4. Proposed Approach

4.1. Problem Formulation

Assume that $\mathcal{T}_i \in \mathbb{T}$ denotes the stroke size of an image, \mathbb{T} denotes the set of all stroke sizes, and $I^{\mathcal{T}_i}$ represents an image I with the stroke size \mathcal{T}_i . The problem studied in this paper is to incorporate different stroke sizes $\mathcal{T}_i \in \mathbb{T}$ into the feed-forward fast neural style transfer model. Firstly, we formulate the feed-forward stylization process as:

$$g(I_c) = I_o, \quad I_o \sim p(I_o | I_c, I_s), \quad (1)$$

where g is the trained generator. And the target statistic $p(I_o)$ of the output image I_o is characterized by two components, which are the semantic content statistics derived from the input image I_c , and the visual style statistics derived from the style image I_s .

Our feed-forward style transfer process for producing multiple stroke sizes can then be modeled as:

$$g'(I_c, \mathcal{T}_i) = I_o^{\mathcal{T}_i}, \quad I_o^{\mathcal{T}_i} \sim p(I_o^{\mathcal{T}_i} | I_c, I_s, \mathcal{T}_i) \quad (\mathcal{T}_i \in \mathbb{T}). \quad (2)$$

We aim to enable one single generator g' to produce stylized results with multiple stroke sizes $\mathcal{T}_i \in \mathbb{T}$ for the same content image I_c .

4.2. Network Architecture

Based on the analysis in Section 3, to incorporate different stroke sizes into one single model, we propose to design a network with adaptive receptive fields and each receptive field is used to learn a corresponding size of stroke. The

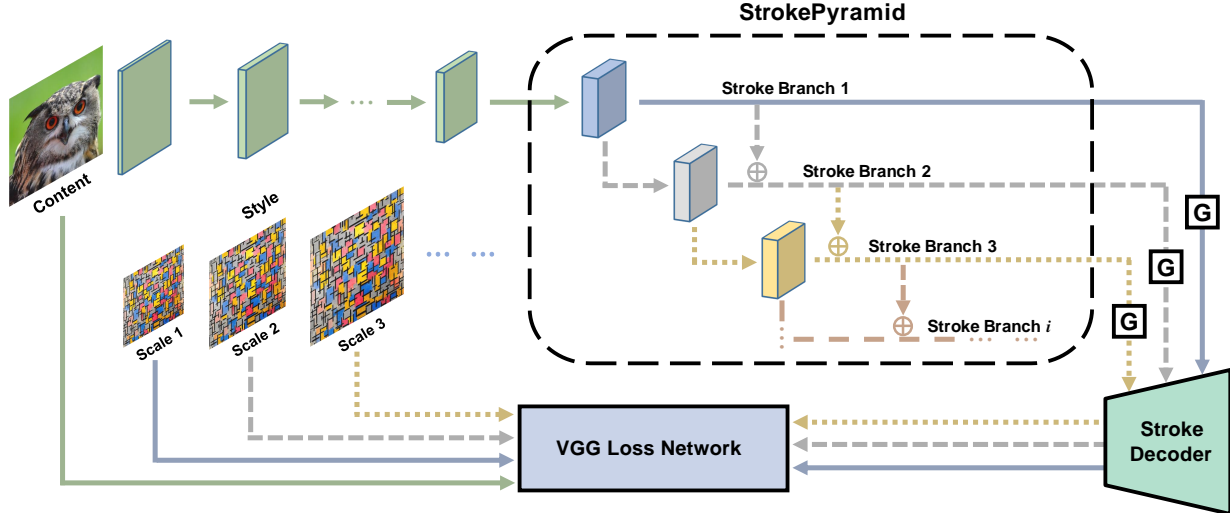


Figure 4. An overview of our network architecture with the *StrokePyramid*. It consists of several stroke branches with gating functions. Each stroke branch corresponds to a specific stroke size.

network architecture of our proposed approach is depicted in Figure 4.

Our network consists of three components. At the core of our network, a *StrokePyramid* module is proposed to decompose the network into several stroke branches. Each branch has a larger receptive field than the previous branch through progressively growing convolutional filters. By handling different stroke branches, the *StrokePyramid* can regulate the receptive field in the generative network. With different receptive fields, the network learns to paint strokes with different sizes. In particular, to better preserve the desired size of strokes, larger strokes are learned with larger receptive fields, as is explained in Section 3. During the testing phase, given a signal which indicates the desired stroke size, the *StrokePyramid* automatically adapts the receptive field in the network and the stylized result with a corresponding stroke size can be produced.

In addition to the *StrokePyramid*, there are two more components in the network, namely the pre-encoder and the stroke decoder. The pre-encoder module refers to the first few layers in the network and is shared among different stroke branches to learn both the semantic content of a content image and the basic appearances of a style. The stroke decoder module takes the feature maps from the *StrokePyramid* as input and decodes the stroke feature into the stylized result with a corresponding stroke size. To determine which stroke feature to decode, we augment a gating function G in each stroke branch. Assume that the selected feature map to be decoded is from the branch \mathcal{B}_{s_k} . The gating function G is then defined as

$$G(\mathcal{F}^{\mathcal{B}_{s_i}}) = \begin{cases} 0 \cdot \mathcal{F}^{\mathcal{B}_{s_i}}, & i \neq k \\ 1 \cdot \mathcal{F}^{\mathcal{B}_{s_i}}, & i = k \end{cases}, \quad (3)$$

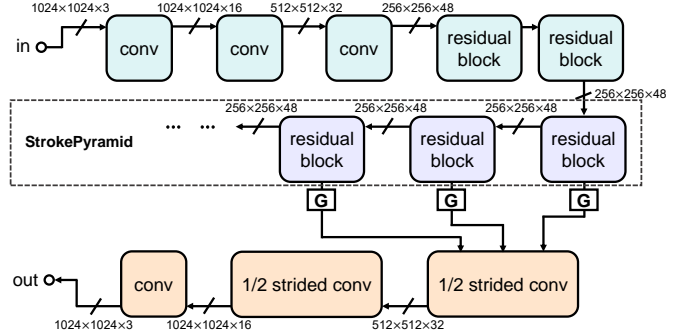


Figure 5. Details of our architecture. An image with size 1024×1024 is taken as our example input.

where $\mathcal{F}^{\mathcal{B}_{s_i}}$ is the output feature map of the branch \mathcal{B}_{s_i} in the *StrokePyramid*, which corresponds to the stroke size \mathcal{T}_i . All the stroke features from the *StrokePyramid* need to go through the gating function and then be fed into the stroke decoder Dec to be decoded into the output result $I_o^{\mathcal{T}_k}$ with the desired stroke size:

$$Dec\left(\sum_i G(\mathcal{F}^{\mathcal{B}_{s_i}})\right) = I_o^{\mathcal{T}_k}. \quad (4)$$

We also demonstrate the detailed settings of our network in Figure 5. Compare with [16], our network uses fewer channels to further reduce the computational complexity.

4.3. Loss Function

Semantic loss. The semantic loss is defined to preserve the semantic information in the content image, which is formulated as the Euclidean distance between the content image I_c and the output stylized image I_o in the feature space of the VGG network [10].

Assume that $\mathcal{F}^l(I) \in \mathbb{R}^{C \times H \times W}$ represents the feature map at layer l in VGG network with a given image I , where C , H and W denote the number of channels, the height and width of the feature map respectively. The semantic content loss is then defined as:

$$\mathcal{L}_c = \sum_{l \in \{l_c\}} \|\mathcal{F}^l(I_c) - \mathcal{F}^l(I_o)\|^2. \quad (5)$$

where $\{l_c\}$ represents the set of VGG layers used to compute the content loss.

Stroke loss. The visual style statistics can be well represented by the correlations between filter responses of the style image I_s in different layers of pre-trained VGG network. These feature correlations can be obtained by computing the Gram matrix over the feature map at a certain layer in VGG network. As the gram-based statistic is scale-sensitive, representations of different stroke sizes can be obtained by simply resizing the given style image.

By reshaping $\mathcal{F}^l(I)$ into $\mathcal{F}^l(I)' \in \mathbb{R}^{C \times (H \times W)}$, the Gram matrix $\mathcal{G}(\mathcal{F}^l(I)') \in \mathbb{R}^{C \times C}$ over feature map $\mathcal{F}^l(I)'$ can be computed as:

$$\mathcal{G}(\mathcal{F}^l(I_s)') = [\mathcal{F}^l(I_s)'] [\mathcal{F}^l(I_s)']^T. \quad (6)$$

The stroke loss for size \mathcal{T}_k can be therefore defined as:

$$\mathcal{L}_{\mathcal{T}_k} = \sum_{l \in \{l_s\}} \|\mathcal{G}(\mathcal{F}^l(\mathcal{R}(I_s, \mathcal{T}_k))') - \mathcal{G}(\mathcal{F}^l(I_o^{\mathcal{B}_{s_k}})')\|^2, \quad (7)$$

where \mathcal{R} represents the function that resizes the style image to an appropriate scale according to the desired stroke size \mathcal{T}_k , and $I_o^{\mathcal{B}_{s_k}}$ represents the output of the k -th stroke branch. $\{l_s\}$ is the set of VGG layers used to calculate the style loss.

The total loss for stroke branch \mathcal{B}_{s_k} can then be written as:

$$\mathcal{L}_{\mathcal{B}_{s_k}} = \alpha \mathcal{L}_c + \beta_k \mathcal{L}_{\mathcal{T}_k} + \gamma \mathcal{L}_{tv}, \quad (8)$$

where α , β and γ are balancing factors. \mathcal{L}_{tv} is a total variation regularization loss to encourage the smoothness in the generated images.

4.4. Training Strategies

Progressive training. To train different stroke branches in one single network, we propose a progressive training strategy. This training strategy stems from the intuition that the training of the latter stroke branch benefits from the knowledge of the previously learned branches. Taken this into consideration, the network learns different stroke sizes with different stroke branches progressively. Assume that the number of the stroke sizes to be learned is K . For every K iterations, the network firstly updates the first stroke branch in order to learn the smallest size of stroke. Then,

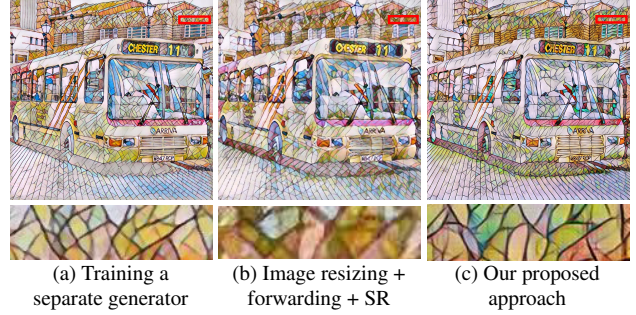


Figure 6. Quality comparison of our proposed algorithm and aforementioned two possible stroke control solutions in Section 4. SR represents the image super-resolution technique [18]. The second line represents the zoom region in the red frame in order to compare the fine texture.

based on the learned knowledge of the first branch, the network uses the second stroke branch to learn the second stroke size with a corresponding scale of the style image. In particular, since the second stroke branch grows the convolutional filters on the basis of the first stroke branch, the updated components in the previous iteration are also adjusted. Similarly, the following stroke branches are updated with the same progressive process. In the next K iterations, the network repeats the above progressive process, since we need to ensure that the network preserves the previously learned stroke sizes.

Incremental training. We also propose a flexible incremental training strategy to efficiently augment new stroke sizes upon a trained model. Given a new desired stroke size, instead of learning from scratch, our algorithm incrementally learns the new stroke size by adding one single layer as a new stroke branch in the *StrokePyramid*. The position of the augmented layer depends on the previously learned stroke sizes and their corresponding receptive fields. By fixing other network components and only updating the augmented layer, the network learns to paint a new size of strokes on the basis of the previously learned stroke features and thus can reach convergence quickly.

5. Experiment

5.1. Implementation Details

Our proposed network is trained on MS-COCO dataset [25]. All the images are cropped and resized to 512×512 pixels before training. We adopt the Adam optimizer [19] during training. The pre-trained VGG-19 network [29] is selected as the loss network and $\{relu1_1, relu2_1, relu3_1, relu4_1, relu5_1\}$ are used as the style layers and $relu4_2$ is used as the content layer. The number of initially learned stroke sizes are set to 3 by

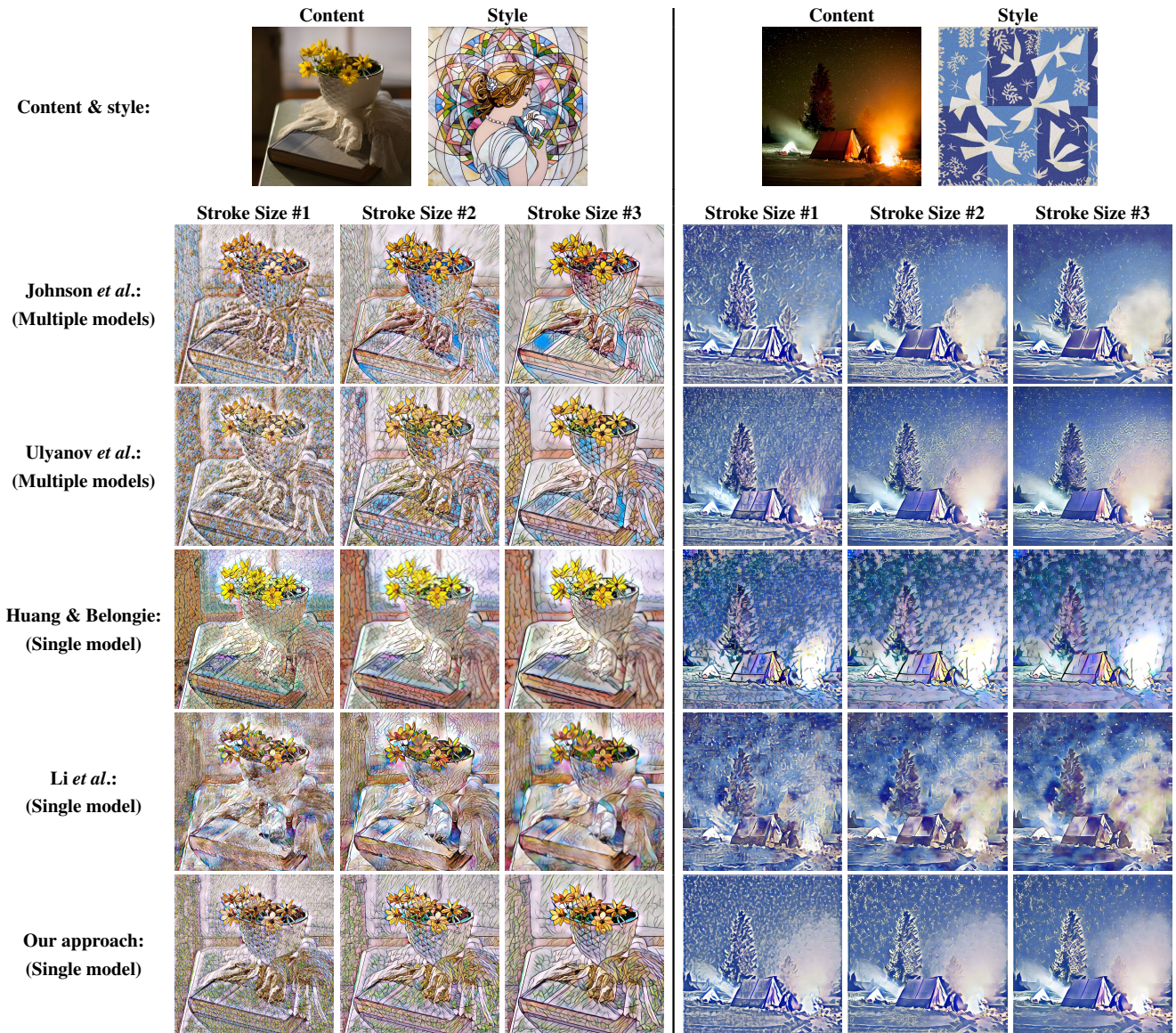


Figure 7. Some example results of different stroke sizes produced by our algorithm and other Fast Style Transfer algorithms [16, 32, 14, 22]. Note that results of each size of strokes with Johnson *et al.* [16] and Ulyanov *et al.* [32] are produced by a separate model. Results of Huang and Belongie’s algorithm [14], Li *et al.*’s algorithm [22], and our algorithm are produced by one single model. All the results with different stroke sizes are all produced by forwarding content images with a fixed size 1024×1024 pixels. All the test content images are never seen during training. Content images are credited to flickr users *Julie Jablonski* and *b togol*.

default. For a fair comparison, the parameters of the existing algorithms are set to be the default values according to their published literature. We implement our algorithm based on Tensorflow [1].

5.2. Qualitative Evaluation

Sample results of our algorithm and two aforementioned possible solutions are shown in Figure 6 (the generator for Figure 6(a) is trained using [16]). Our algorithm achieves comparative results with the first possible solution in Figure 6(a) regarding the quality while preserving the flexi-

bility of the second possible solution in Figure 6(b). Figure 7 shows sample results of our algorithm and other Fast Style Transfer algorithms. More results can be found in the supplementary material¹. For results from the algorithm of Johnson *et al.* [16] and Ulyanov *et al.* [32], training a separate scale-specific generator for each stroke size is adopted, which is the first aforementioned possible solution of stroke size control. Totally three generators are trained using [16] and [32] for each style in Figure 7. It can be noticed that our

¹<https://youtu.be/UNG38tdMSMg>

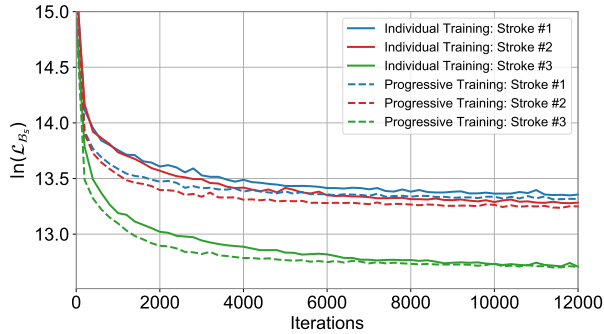


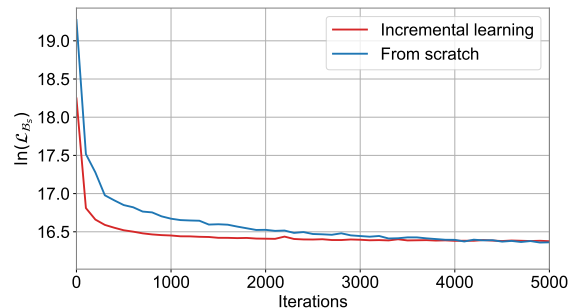
Figure 8. Training curves for three stroke sizes with (dashed curves) and without (solid curves) progressive training strategy. Our progressive training strategy is faster than the individual training in terms of the converging speed, except for the learning of stroke #3 where progressive training achieves comparable convergence speed.

algorithm achieves competitive results against [16, 32] but exploits only one single pre-trained model. To achieve visually plausible results with K stroke sizes, [16, 32] need to train K corresponding stroke-specific models, which brings additional time cost. In contrast, our algorithm needs less training time due to our progressive training strategy. The algorithm of Huang and Belongie [14] and the algorithm of Li *et al.* [23] belong to the category of ASPM algorithms and are able to transfer an arbitrary style through one single model. Therefore, these two algorithms do not need to train a style-specific generator in advance and can control the stroke size during stylization by just feeding different scales of style images. However, [14] is not effective at producing some fine textures in some styles (Figure 7, the fourth row). Although [23] captures finer textures, the details are not well preserved in some cases (Figure 7, the fifth row).

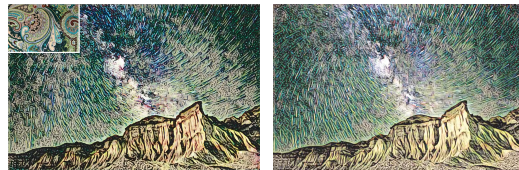
5.3. Quantitative Evaluation

In terms of the quantitative evaluation, we focus on three evaluation metrics in this section, which are: training curves during jointly training and incremental training; average content and style loss for test content images; training time for our single model and corresponding generating time for results with different stroke sizes.

Training curve analysis. To demonstrate the effectiveness of our progressive training strategy, we record the stroke losses when learning several sizes of strokes progressively and learning different strokes individually. The result is shown in Figure 8. The reported loss values were averaged over 15 randomly selected batches of content images. It can be observed that the network which progressively learns multiple stroke sizes converges relatively faster than the one which learns only one single stroke size individ-



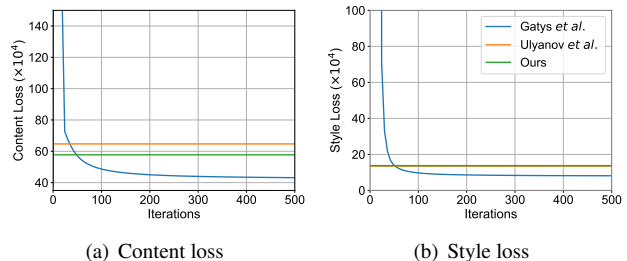
(a) Training curves of learning incrementally and from scratch.



(b) Incremental learning.

(c) Learning from scratch.

Figure 9. Comparisons between incremental learning and learning from scratch in terms of training curves and stylization quality.



(a) Content loss

(b) Style loss

Figure 10. Comparisons of the average content and style loss of our algorithm with state-of-the-art Neural Style Transfer algorithms.

ually. The result indicates that during progressive training, the latter stroke branch benefits from the learned knowledge of the previous branches, and can even improve the training of previous branches through a shared network component in turn. To validate our stroke incremental training strategy, we present both the training curves of the incremental learning and learning from scratch in Figure 9. While achieving comparable stylization quality, incrementally learning a stroke can significantly speed up the training process compared to learning from scratch.

Average loss analysis. To measure how well the loss function is minimized, we compare the average content and style loss of our algorithm with other style transfer methods. For a fair comparison, the loss functions of different algorithms are kept the same. The recorded values are averaged over 100 content images and 5 style images. For each style, we calculate the average loss of the three stroke sizes. As shown in Figure 10, the average style loss of our

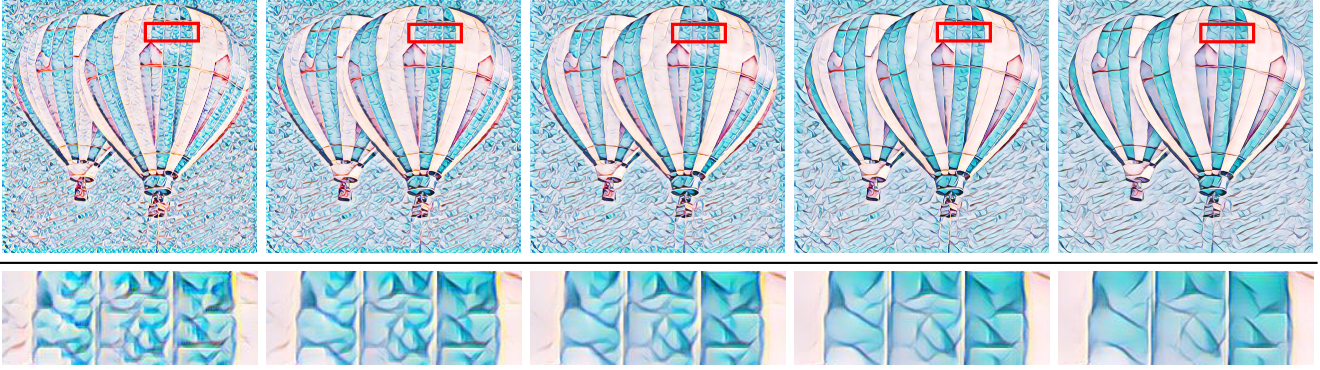


Figure 11. Results of stroke interpolation. We zoom in on the same region (red frame) to observe the variations of stroke sizes.

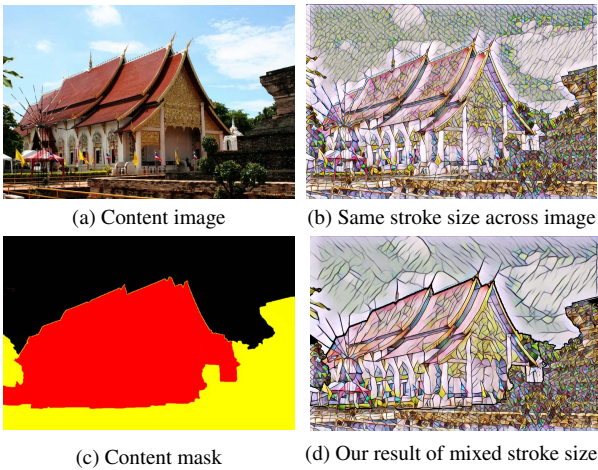


Figure 12. Our algorithm allows flexible spatial stroke size control during stylization. The result produced by our single model can have mixed stroke sizes, which is more consistent with an artist’s artwork in reality.

algorithm is similar to [32], and our average content loss is slightly lower than [32]. This indicates that our algorithm is comparable to [32] regarding the ability to minimize the loss function.

Speed analysis. Fully training one single model with three stroke sizes takes about 2 hours on a single NVIDIA Quadro M6000. For generating time, it takes averagely 0.09 seconds to stylize an image with size 1024×1024 on the same GPU using our algorithm. Since our network architecture is similar with [16, 32] but with a shorter path for some stroke sizes, our algorithm can be on average faster than [16, 32] and further faster than Huang and Belongie’s algorithm and Li *et al.*’s algorithm according to the speed analysis in [14, 23].

5.4. Runtime Controls

Stroke Interpolation. By interpolating between the output feature maps in the *StrokePyramid*, our algorithm can achieve arbitrary intermediate stroke sizes. Given a content image I_c , we assume that $\mathcal{F}^{\mathcal{B}_{s_m}}$ and $\mathcal{F}^{\mathcal{B}_{s_n}}$ are two output feature maps in the *StrokePyramid*, which can be decoded into the stylized results with two stroke sizes $I_o^{\mathcal{T}_m}$ and $I_o^{\mathcal{T}_n}$ respectively. The interpolated feature $\mathcal{F}^{\mathcal{B}_{\bar{s}}}$ can then be calculated as:

$$\mathcal{F}^{\mathcal{B}_{\bar{s}}} = a\mathcal{F}^{\mathcal{B}_{s_m}} + (1 - a)\mathcal{F}^{\mathcal{B}_{s_n}} \quad (9)$$

By changing the value of a and feeding the obtained $\mathcal{F}^{\mathcal{B}_{\bar{s}}}$ into the stroke decoder module, stylized results with arbitrary intermediate stroke sizes $I_o^{\tilde{\mathcal{T}}}$ can be produced, as shown in Figure 11.

Spatial Stroke Control. Previously, in the community of Fast Style Transfer, stylized results usually have almost the same stroke size across the whole image, which is impractical in the real case. As is shown in Figure 12, our algorithm supports mixed stroke sizes in different spatial regions and in this way, the contrast information in the content image can be enhanced. Our spatial stroke control is achieved by feeding masked content image through different corresponding stroke branches and then combining these stylized results. The mask can be obtained either by manual labelling or forwarding the content image through a pre-trained semantic segmentation network.

6. Conclusions

In this paper, we present a Fast Style Transfer deep network that allows a flexible control for the stroke size during stylization. By using almost the same number of parameters as the previous Fast Style Transfer algorithm, our network is capable of learning multiple stroke sizes. The main idea behind our technique is proposing a *StrokePyramid* module to endow the network with adaptive receptive

fields and the network can learn to paint different stroke sizes with the corresponding size of the receptive field. By adopting the proposed progressive training strategy, our network achieves faster convergence and through incremental training strategy, new stroke sizes can be augmented in a trained mode. Finally, our network can produce distinct stroke sizes in different output images or different spatial regions within the same output image. Embedded with other existing perceptual factor controlling strategies, our work takes a step towards breaking the tradeoff between the flexibility and efficiency in Fast Style Transfer. Our future work will focus on exploring the influence of other factors on the stroke size. Another research direction is to apply the idea of the *StrokePyramid* into the MSPM so as to efficiently learn multiple stroke sizes for multiple styles.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016. [6](#)
- [2] D. Chen, L. Yuan, J. Liao, N. Yu, and G. Hua. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. [1](#)
- [3] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets. *Atrous Convolution, and Fully Connected CRFs*. *arXiv*, 2016. [2](#)
- [4] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei. Deformable convolutional networks. *arXiv preprint arXiv:1703.06211*, 2017. [2](#)
- [5] V. Dumoulin, J. Shlens, and M. Kudlur. A learned representation for artistic style. In *International Conference on Learning Representations*, 2017. [1](#)
- [6] A. A. Efros and W. 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. ACM, 2001. [1](#)
- [7] M. Elad and P. Milanfar. Style transfer via texture synthesis. *IEEE Transactions on Image Processing*, 26(5):2338–2351, 2017. [1](#)
- [8] O. Frigo, N. Sabater, J. Delon, and P. Hellier. Split and match: Example-based adaptive patch sampling for unsupervised style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 553–561, 2016. [1](#)
- [9] L. A. Gatys, A. S. Ecker, and M. Bethge. Texture synthesis using convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 262–270, 2015. [1](#)
- [10] L. A. Gatys, A. S. Ecker, and M. 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](#), [3](#), [4](#)
- [11] L. A. Gatys, A. S. Ecker, M. Bethge, A. Hertzmann, and E. Shechtman. Controlling perceptual factors in neural style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. [1](#), [2](#), [3](#)
- [12] B. Gooch and A. Gooch. *Non-photorealistic rendering*. A. K. Peters, Ltd., Natick, MA, USA, 2001. [1](#)
- [13] A. Hertzmann, C. E. Jacobs, N. Oliver, B. Curless, and D. H. Salesin. Image analogies. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 327–340. ACM, 2001. [1](#)
- [14] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. [2](#), [6](#), [7](#), [8](#)
- [15] Y. Jing, Y. Yang, Z. Feng, J. Ye, Y. Yu, and M. Song. Neural style transfer: A review. *arXiv preprint arXiv:1705.04058*, 2017. [1](#)
- [16] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711, 2016. [1](#), [4](#), [6](#), [7](#), [8](#)
- [17] B. Julesz et al. Textons, the elements of texture perception, and their interactions. *Nature*, 290(5802):91–97, 1981. [2](#)
- [18] J. Kim, J. Kwon Lee, and K. Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1646–1654, 2016. [5](#)
- [19] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. [5](#)
- [20] C. Li and M. Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2479–2486, 2016. [1](#)
- [21] C. Li and M. Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *European Conference on Computer Vision*, pages 702–716, 2016. [1](#)
- [22] Y. Li, F. Chen, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Diversified texture synthesis with feed-forward networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. [1](#), [6](#)
- [23] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems*, 2017. [2](#), [7](#), [8](#)
- [24] Y. Li, N. Wang, J. Liu, and X. Hou. Demystifying neural style transfer. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 2230–2236, 2017. [1](#), [3](#)
- [25] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. [5](#)
- [26] M. Lu, H. Zhao, A. Yao, F. Xu, Y. Chen, and L. Zhang. Decoder network over lightweight reconstructed feature for fast semantic style transfer. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. [1](#), [2](#)

- [27] I. Prisma Labs. Prisma: Turn memories into art using artificial intelligence, 2016. [1](#)
- [28] P. Rosin and J. Collomosse. *Image and video-based artistic stylisation*, volume 42. Springer Science & Business Media, 2012. [1](#)
- [29] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [3](#), [5](#)
- [30] T. Strothotte and S. Schlechtweg. *Non-photorealistic computer graphics: modeling, rendering, and animation*. Morgan Kaufmann, 2002. [1](#)
- [31] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *International Conference on Machine Learning*, pages 1349–1357, 2016. [1](#)
- [32] D. Ulyanov, A. Vedaldi, and V. 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*, 2017. [1](#), [2](#), [6](#), [7](#), [8](#)
- [33] X. Wang, G. Oxholm, D. Zhang, and Y.-F. Wang. Multimodal transfer: A hierarchical deep convolutional neural network for fast artistic style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. [2](#)
- [34] Z. Wei, Y. Sun, J. Wang, H. Lai, and S. Liu. Learning adaptive receptive fields for deep image parsing network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2434–2442, 2017. [2](#)
- [35] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *International Conference on Learning Representations*, 2016. [2](#)
- [36] H. Zhang and K. Dana. Multi-style generative network for real-time transfer. *ArXiv e-prints*, Mar. 2017. [1](#)
- [37] S.-C. Zhu, C.-E. Guo, Y. Wang, and Z. Xu. What are textures? *International Journal of Computer Vision*, 62(1):121–143, 2005. [2](#)