# Arbitrary Style Transfer with Deep Feature Reshuffle

Shuyang Gu[1*]     Congliang Chen[2*]     Jing Liao[3]     Lu Yuan[3]

[1]University of Science and Technology of China     [2]Peking University     [3]Microsoft Research

gsy777@mail.ustc.edu.cn     chcoliang@pku.edu.cn     {jliao,luyuan}@microsoft.com

## Abstract

*This paper introduces a novel method by reshuffling deep features (*i.e.*, permuting the spacial locations of a feature map) of the style image for arbitrary style transfer. We theoretically prove that our new style loss based on reshuffle connects both global and local style losses respectively used by most parametric and non-parametric neural style transfer methods. This simple idea can effectively address the challenging issues in existing style transfer methods. On one hand, it can avoid distortions in local style patterns, and allow semantic-level transfer, compared with neural parametric methods. On the other hand, it can preserve globally similar appearance to the style image, and avoid wash-out artifacts, compared with neural non-parametric methods. Based on the proposed loss, we also present a progressive feature-domain optimization approach. The experiments show that our method is widely applicable to various styles, and produces better quality than existing methods.*

## 1. Introduction

This process of rendering a content image in the style of another image is referred to as Style Transfer. The problem of style transfer has its origin from non-photo-realistic rendering [26], and is closely related to texture synthesis and transfer [11, 12, 13]. These methods typically rely on low-level statistics and often fail to capture semantic structures. Recently, the work of Gatys et al. [16] opened up a new field called Neural Style Transfer, which uses Convolutional Neural Network (CNN) [25] to change the style of an image while preserving its content. It is flexible enough to combine content and style of arbitrary images.

Style transfer is receiving increasing attention from computer vision researchers because it involves two interesting topics: image representation and image synthesis.

Some early representations, like multi-resolution [9], pyramid [18], wavelet [35], used in traditional texture synthesis and transfer, are mainly for statistics matching. The recent work [16] showed that the representations of image content and style were separable by variant CNN convolutional layers. Moreover, the representation provides the possibility for image decoupling and recombining.

Image synthesis methods, whether traditional or neural, can be broadly categorized as *parametric* and *non-parametric*. Specifically, for neural methods, the parametric methods match the global statistics of deep features, like Gram matrix [16, 23] and its approximates [30, 33], mean and variance [20, 10], histogram [41]; while the non-parametric methods [28, 29, 8, 32] directly find neural patches similar to the given example. However, to the best of our knowledge, there are no work connecting these neural methods to form a complementary solution.

The *neural parametric* models (*e.g.*, [16]) yields results, preserving the content of image and the overall looking of the artwork. However, the models will distort local style patterns (shown in the first row of Fig. 1) or cannot obtain locally semantic-level transfer (*e.g.*, eye-to-eye in the second row of Fig. 1). The *neural non-parametric* models (*e.g.*, [28]) can address these issues well, but their example matching uses a greedy optimization, causing the decreasing in the richness of the style patterns (shown in the first row of Fig. 2), and introducing wash-out artifacts [21] (see the second row of Fig. 2). It suggests that such *neural non-parametric* models should consider global constraint, borrowing from *neural parametric models*.

In this paper, we propose a novel neural style transfer algorithm which owns the advantages of both neural parametric and non-parametric methods. This is achieved by *deep feature reshuffle*, which refers to spatially rearranging the position of neural activations. We reshuffle the features of the style image according to the content image for style transfer. On one hand, the feature reshuffle enforces the distribution of style patterns, between the transferring result and the style image, to be globally consistent. It can be

theoretically proved that reshuffling features of style image equals to the optimization of Gram matrices, a commonly used statistics in *neural parametric* models. On the other hand, a certain type of reshuffling features can help achieve locally semantic matching between images, as well as *neural non-parametric* models.

We reformulate the objective function of neural style transfer in the fashion of reshuffle, and then connects both kinds of methods. To avoid exhaustively optimizing the energy function in image domain as similar as [16], we propose a novel optimization method in feature domain. We can progressively recover features from high-level layers to low-level ones, and train a decoder to convert recovered features back to the image. This way is more efficient.

Our experiments show that this method can effectively accomplish the transfer for arbitrary styles, yield results with global similarity to the style and local plausibility. We summarize main contributions as follows:

- We provide a new understanding of *neural parametric* models and *neural non-parametric* models. Both can be integrated by the idea of *deep feature reshuffle*.

- We define a new energy function based on *deep feature reshuffle*, which is simple, flexible, and better than either *neural parametric* or *non-parametric* methods.

- We train a new level-wise decoder to allow us efficiently optimize our feature-domain energy function in a pyramid manner.

## 2. Related Work

The problem of style transfer involves two sub-problems: representation and synthesis. Inspired by the success of CNN in style transfer, we also use neural representation for image decoupling, and better matching. In this paper, we focus on synthesis problem, which can be categorized as parametric and non-parametric.

In fact, parametric and non-parameter synthesis methods early occur at texture synthesis and transfer. Parametric methods [9, 18, 35] start from random noise, and then iteratively update it until the desired global statistics is satisfied. However, it is challenging to find a proper statistical model for representation and fine matching. By contrast, non-parametric models [11, 12, 40] use a simple patch representation (*e.g.*, color [11, 19, 2], curvilinear features [42], edge [24] and its orientation [27]), and find the most similar patches by nearest neighbor search. All above methods only use the low-level features for synthesis, which limits to capture semantic structures.

Gatys et al. [16] pioneer the neural texture synthesis and style transfer by successfully applying CNN (pre-trained VGG networks [38]) to this problem. The spirit of their synthesis method is parametric, which statistically matches
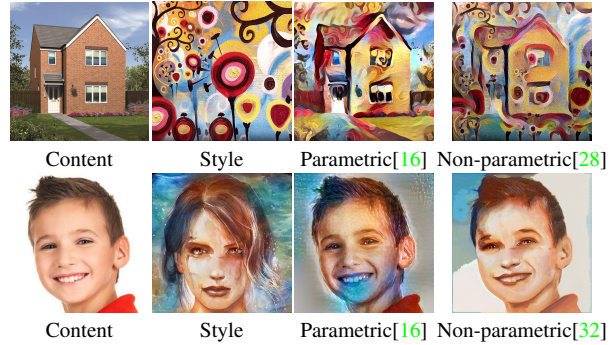


Figure 1. Parametric style transfer fails to preserve some local texture patterns of the style, e.g. circles in the upper row, and it have weak spatial constraint, like rendering background colors in the face region (lower row).
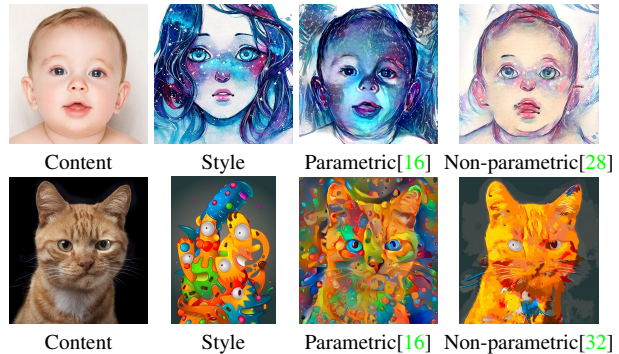


Figure 2. Non-parametric style transfer is sometimes globally less similar to the style (upper row) and repetitively uses the same patches to cause wash-out artifacts (lower row).

both the content and style features by their Gram matrices. The solution is general to varies of artistic styles, and begins to considering semantic structures. To improve the quality, some complimentary information is incorporated into the statistical model, including spatial correlation [37], face guidance [36], user controls [4, 17], and segmentation masks [34]. To accelerate, a feed-forward generative network [23, 39, 10, 6] is directly learnt instead, but they are still limited to a fixed number of pre-trained styles. More recently, some work [30, 20, 33] further allow arbitrary style transfer in feedforward networks. The backend idea is to match the statistics of content features at intermediate layers to that of the style features, and then train a decoder to turn features to the image. And [7, 5] further extend this kind of method to video and stereoscopic 3D style transfer.

Non-parametric neural style transfer method is firstly proposed by Li et al. [28]. They reformulate the style transfer based on Markov Random Field (MRF): searching local neural patches from the style image to satisfy the local structure prior of content image. Compared with neural parametric methods, this method can reproduce local textures more faithfully. They also train a Generative Adversarial Networks (GAN) to accelerate the optimization [29]. Chen et al. [8] use de-VGG networks and the patch-based method for fast arbitrary style transfer. Another represen-

2

tative work is called Deep Analogy [32], which proposes accurate semantic-level patch match algorithm by considering bidirectional constraint and pyramids refinement. Our method is based on non-parametric model, sharing all the advantages, like local similarity to style and semantic-level transfer. Moreover, it allows some global constraints to avoid washout artifacts, and obtains global consistency, as similar as parametric methods.

## 3. Understanding Neural Style Transfer

In this section, we explore the relationship between neural parametric method (*e.g.*, [16]) and neural non-parametric method (*e.g.*, [28]). Then, we realize that the *feature reshuffle* can theoretically be a complementary solution for both methods.

For the task of style transfer, we want to generate a stylization result $I_o$, given the content image $I_c$ and the style image $I_s$. For simplicity, we suppose $I_o$, $I_c$ and $I_s$ are with the same size[1], and consider the single layer feature of these image, which are denoted as $F_o$, $F_c$ and $F_s$ respectively. $F \in \Omega^{c \times h \times w}$ is indeed 3D tensor, where $c, h, w$ denote channel number, height and width respectively.

**Neural Parametric.** In [16], the energy function consists of a content term $\mathcal{L}_{cont}$ and a style term $\mathcal{L}_{sty}$ [2]:

$$\mathcal{L}_{total} = \alpha\mathcal{L}_{cont} + (1-\alpha)\mathcal{L}_{sty}, \qquad (1)$$

where $\alpha$ is the tradeoff to balance content and style.

The content loss $\mathcal{L}_{cont}$ is defined by the feature difference between the content image $I_c$ and the (yet unknown) stylized image $I_o$:

$$\mathcal{L}_{cont} = ||F_o - F_c||_F^2, \qquad (2)$$

where $|| \cdot ||_F^2$ denotes Frobenius norm.

For the style loss $\mathcal{L}_{sty}$, Gatys et al. [16] uses Gram matrix $G(i, j)$ to obtain correlations between filter responses. It was used to measure texture correlation in texture synthesis algorithm [15]. Gram matrix $G(i, j)$ is defined as the inner product between the $i$_th and $j$_th feature channels:

$$G(i, j) = \sum_p F(i, p)F(j, p), \qquad (3)$$

where $p$ denotes the 2D spatial location in feature map and $G \in \Omega^{C \times C}$. Then, the style loss $\mathcal{L}_{sty}$ (also called global style loss later) is defined by the difference between Gram matrices of $F_o$ and $F_s$:

$$\mathcal{L}_{sty} = ||G_o - G_s||_F^2 {}^3 \qquad (4)$$

---

[1] For different size, the principle still holds true, when these features and terms are normalized according to feature size.

[2] We ignore the image regularization term here, because it is common, and only takes effects to subtle noise. Its effect to our analysis can be negligible.

[3] It should be $\mathcal{L}_{sty} = \frac{1}{\mathcal{Z}}||G_o - G_s||_F^2$, where $\mathcal{Z} = 4 \times c^2 \times h^2 \times w^2$. We leave out $\mathcal{Z}$ for simplicity, and it won't effect our analysis.

**Neural Non-parametric.** In [28], the energy function is also defined with two terms:

$$\mathcal{L}_{total} = \alpha\mathcal{L}_{cont} + (1-\alpha)\mathcal{L}_{match}. \qquad (5)$$

The content loss $\mathcal{L}_{cont}$ is identical to Eq. (2); while the style loss $\mathcal{L}_{match}$ measures the neural patch-based similarity. Let $\Psi(F)$ denote the list of all local patches extracted from feature map $F$. Each *neural patch* centered at the location $p$ of feature $F$ is indexed as $\Psi_p(F)$. The loss $\mathcal{L}_{match}$ (also called local style loss later) is defined as:

$$\mathcal{L}_{match} = \sum_p ||\Psi_p(F_o) - \Psi_{\text{NN}(p)}(F_s)||_F^2, \qquad (6)$$

where $\text{NN}(p)$ is the index of the patch in $\Psi(F_s)$ which is the most similar to $\Psi_p(F_o)$. The best matching index is calculated using normalized cross-correlation over all local patches in the style feature $F_s$:

$$\text{NN}(p) = \underset{p'=1,2,...,\Theta}{\arg\max} \frac{\Psi_p(F_o) \cdot \Psi_{p'}(F_s)}{||\Psi_p(F_o)||_F^2 \cdot ||\Psi_{p'}(F_s)||_F^2}, \qquad (7)$$

where the operator $\cdot$ denotes inner product, and $\Theta = h \times w$.

Both methods share the same content loss, but have different style loss terms. The global style loss (Eq. (4)) measures global statistics, but ignores the spatial layout of features [31]. On the contrast, the local style loss (Eq. (6)) encourages to find optimal feature layout for each local patch individually (Eq. (7)), but without global constraint.

**Neural Feature Reshuffle.** Can an optimal feature $F_o$ be achieved to satisfy both global and local style terms simultaneously? Yes, we find that *feature reshuffle* can theoretically be an ideal condition, which can make both global and local style terms be zero.

Specifically, *feature reshuffle* means that we permute the spatial location of feature map $F_s$ to reconstruct a new feature map $F_o$. Let $\text{SF}(p)$ be the permuted location of the feature corresponding to the original location $p$. By feature reshuffle, the reconstructed feature map $F_o$ is denoted as:

$$F_o(p) = F_s(\text{SF}(p)), \text{ and } F_s(p) = F_o(\widetilde{\text{SF}}(p)), \qquad (8)$$

where $\widetilde{\text{SF}}(p)$ is the inverse reshuffle function. In other words, $\text{SF}(p)$ is a bijection: for each location in $F_o$ there exists exactly one location in $F_s$ corresponding to it.

We can deduce the global style loss (in Eq. (4)) with the reshuffle solution (see in Eq. (8)) as:

$$\mathcal{L}_{sty} = \sum_{i,j} ||\sum_p F_o(i,p)F_o(j,p) - \sum_p F_s(i,p)F_s(j,p)||^2$$
$$= \sum_{i,j} ||\sum_p F_s(i, \text{SF}(p))F_s(j, \text{SF}(p)) - \sum_p F_s(i,p)F_s(j,p)||^2$$
$$= 0,$$

3

which indicates reshuffling features of style does not effect Gram matrices of style features, leading the global style loss to be zero.

According to Eq. (8), the feature $F_o(p)$ at each point $p$ is reconstructed by $F_s(SF(p))$. When the patch size is $1 \times 1$, $\Psi_p(F_o)$ can be denoted as $F_o(p)$. Then, the nearest neighbor of $F_s(SF(p))$ to be found in the feature map $F_s$ is just itself, *i.e.*, $F_s(SF(p)) = F_s(NN(p)) = \Psi_{NN(p)}(F_s)$. Based on these deductions, we can rewrite Eq. (6) as

$$\mathcal{L}_{match} = \sum_i || \sum_p (F_o(i,p) - F_s(i, \text{SF}(p)))||^2 = 0.$$

As summary, the above theoretical derivation demonstrates feature reshuffled from the style image simultaneously minimizes both global style loss and local style loss (when patch size is $1 \times 1$).

## 4. Method

Based on the idea of *deep feature reshuffle*, we propose a novel neural style transfer algorithm, which integrates global and local style losses in the whole objective function. We first present a new style loss called *reshuffle loss*, which would be combined with content loss as well. Then, we show the optimization in a single feature layer. Such an optimization can be done in image domain, similar to the manner of [16], by iteratively forward and backward passing the networks. For acceleration, we propose two efficient ways: 1) the optimization can be done in feature domain, and does not need to back propagate to the image at every time; 2) we progressively optimize the features across multiple layers, and the exhaustive patch match in the fine layer can be guided by the matching result in the coarse layer.

### 4.1. Reshuffle Loss Function

We define a new *reshuffle loss* for the style loss, which only slightly modifies the local style loss term (see Eq. (6)):

$$\mathcal{L}_{shuf} = \sum_p ||\Psi_p(F_o) - \Psi_{\text{NNC}(p)}(F_s)||_{\text{F}}^2, \qquad (9)$$

where the original nearest neighbor (NN) search $NN(p)$ is replaced by a new function $NNC(p)$. It is also the NN search, but constrained by the times of patch usage. For strict reshuffle, we require each patch in the source to be only mapped once, as shown in Fig. 3.

Indeed, sometimes such one-usage constrain is too harsh. For example, the content image has two faces but the style image only has one. We relax the constraint in our energy term to allow more times of patch usage. Although it will sacrifice the global term to some extent, it greatly improve the robustness. A hard cutoff of usage is difficult to be found for every case. Instead, our constrained NN search function softly encourages the uniform usage of patch, inspired by



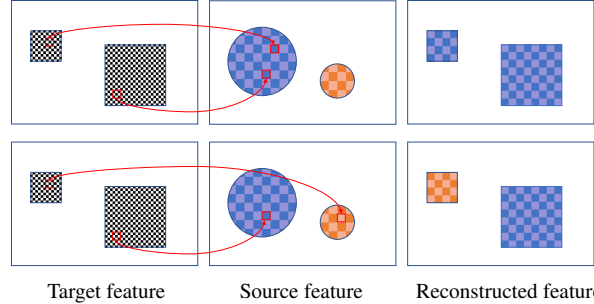| Target feature | Source feature | Reconstructed feature |

Figure 3. The comparison between normal NN search (upper row) and our NN search constrained by reshuffle (lower row). Ours strictly requires each patch in the source to be only mapped once. Notice that the bottom only shows a possible solution of reshuffle.

[24, 14], which considered similar constraints in image domain, and is defined as:

$$\text{NNC}(p) = \underset{p'=1,2,...,\Theta}{\arg\max} \; \left( \frac{\Psi_p(F_o) \cdot \Psi_{p'}(F_s)}{||\Psi_p(F_o)||_{\text{F}}^2 \cdot ||\Psi_{p'}(F_s)||_{\text{F}}^2} \right. \qquad (10)$$
$$\left. -\lambda \frac{\Gamma(\Psi_{p'}(F_s))}{R \times R} \right),$$

where $R$ is the patch size. $\Gamma(p)$ keeps track how many times of each pixel has been used in all patches covering it, and $\Gamma(\Psi_p)$ refers to the total usages of a patch normalized by its area $R \times R$, namely,

$$\Gamma(\Psi_p) = \sum_{p \in \Psi_p} \frac{\Gamma(p)}{R \times R}. \qquad (11)$$

This term requires each pixel to be used only once, as possible as it can. $\lambda$ controls the relative contribution of uniformity enforcement. The $NNC(p)$ can be optimized with the EM-like algorithm described in [24, 14], which extends PatchMatch algorithm [3] to keep track of the usage as well.

### 4.2. Single Layer Optimization

The objective function for single layer is defined as:

$$\mathcal{L}_{total} = \alpha \mathcal{L}_{cont} + (1-\alpha)\mathcal{L}_{shuf}. \qquad (12)$$

The most direct solution is to optimize it in image domain. Similar to [16, 28], we suppose the output image $I_o$ to be either the content image or random noise initially. Then we pass it to the VGG19 network (pre-trained on ImageNet for object recognition), and get the feature maps $F_o^l$ at $relul\_1$ layer. We optimize NNC with the feature map $F_o^l$ and $F_s^l$, shown in Eq. (10). The gradient of the energy function $\mathcal{L}_{total}$ (in Eq. (12)) is then computed and back propagated to update $I_o$. Such processing always needs hundreds of iterations to converge by gradient decent method (*e.g.*, L-BFGS [43]). Hundreds times of foward-backward passes of networks and NN search make it prohibitively slow.

An alternatively fast solution is to directly optimize $F_o^l$ in feature domain, and then reverse it to the output image
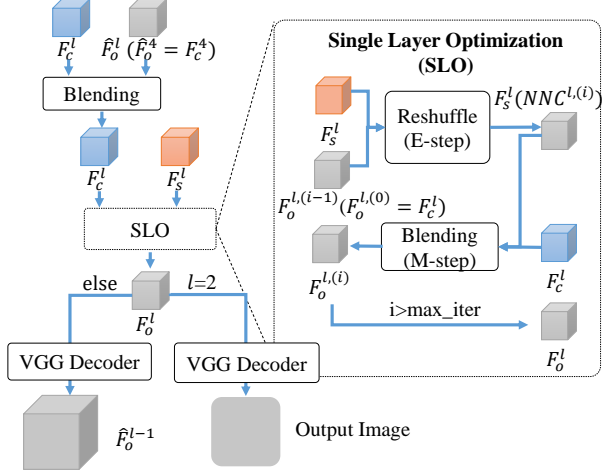
Figure 4. System Pipeline

by training VGG-like image decoder. Here, we first discuss how to get $F_o^l$ by minimizing Eq. (12). Considering Eq. (10) is not directly differentiable, the objective function (Eq. (12)) is hard to be optimized by a standard solver like SGD. To solve this issue, we adopt an iterative EM-like algorithm, which was used in [24] with good convergency. Specifically, $F_o^l$ is initialized to the content features, $F_o^{l,(0)} = F_c^l$. In the E-step of each iteration $i$, the constrained NN field $NNC^{l,(i)}$ is computed by matching $F_o^{l,(i-1)}$ and $F_s^l$ (see Eq. (10)). We then get the feature map $F_s^l(NNC^{l,(i)})$, by warping $F_s^l$ with $NNC^{l,(i)}$ and average-voting overlapping neighbor patches at each location. In the M-step, we obtain the blended feature map $F_o^{l,(i)}$ by a linear combination: $F_o^{l,(i)} = \alpha F_c^l + (1-\alpha)F_s^l(NNC^{l,(i)})$, according to Eq. (12). After several iterations (less than 10 normally), $F_o^{l,(i)}$ will converge to the optimal feature $F_o^l$.

Once $F_o^l$ is achieved, we will get the output image $I_o$ by decoding $F_o^l$ to image domain. The decoder can be pre-trained for efficiency [30, 20, 33]. By contrast, we adopt an different decoder learning from theirs. The details of decoder training are discussed in Section 4.4.

### 4.3. Multi-layer Progressive Optimization

We actually adopt a multi-layer progressive optimization instead of independent single layer optimization described in the above section. This way provides three advantages. First, leveraging mutil-layer features can generate richly textured results. Second, multi-resolution processing can help avoid getting stuck in worse local minima. Last, with the NNF guidance from coarse layers, we may accordingly decrease the search range for more efficient matching.

Fig. 4 shows our algorithm pipeline. The multi-layer optimization considers such three layers $l = 2, 3, 4$. In initialization, we obtain feature maps $\{F_c^l\}_{l=4,3,2}$ and $\{F_s^l\}_{l=4,3,2}$ of $I_c$ and $I_s$ respectively by feeding them to VGG-19.

We start from the coarsest layer $l = 4$, and perform the single layer optimization (see Section 4.2) at this layer, which uses EM-like constrained patch match algorithm. After it, we get the updated feature $F_o^4$, which is then decoded to the next layer $l - 1$ by our trained decoder network, denoted as $\hat{F}_o^{l-1}$. To leverage the result from the coarse layer, we update the content feature $F_c^{l-1}$ at layer $l - 1$ by linearly combining it with $\hat{F}_o^{l-1}$, namely, $F_c^{l-1} \leftarrow \beta \hat{F}_o^{l-1} + (1-\beta)F_c^{l-1}$. Next, we use the updated content feature $F_c^{l-1}$ for single layer optimization in layer $l - 1$. The blending can inherit information from the coarser layers. We iterate the process from the coarsest layer $l = 4$ to the finest layer $l = 2$. Finally, we decode the optimal feature $F_o^2$ to obtain the output image $I_o$. The pseudo code of our implementation is listed in Algorithm 1. Code has been made available at: https://github.com/msracver/Style-Feature-Reshuffle

### 4.4. Decoder Training

Li et al. [30] proposed a universal decoder for fast style transfer. However, their method is not very economic, since $L$ various decoders are needed to respectively decode features from every different layer. These decoders do not share weights in training. In this paper, we propose a new training strategy which provides only a single decoder for features may from different layers. The comparisons between the two training strategies are shown in Fig. 5.

The architecture of our decoder uses a symmetrical structure to that of VGG-19 encoder network. The training strategy is bottom-up, starting from layer 1 to layer $L$. When we train the $l$_th layer of decoder, image $I$ first feeds to the encoder sub-net $\mathcal{E}_1^{l-1}$ involving encoder layer 1 to
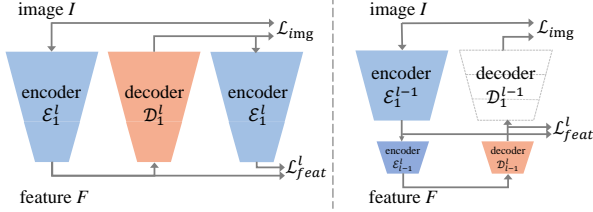
Figure 5. Comparison of our decoder (right) and Li et al.'s [30] (left). To train a decoder in layer $l$, ours fixes and shares the pretrained part from layer 1 to $l-1$, while theirs retrains all parts.



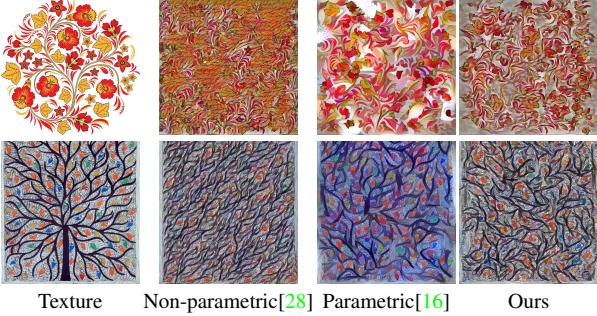<div style="text-align:center">Texture    Non-parametric[28]   Parametric[16]    Ours</div>

Figure 6. Texture synthesis results by minimizing different style losses on layer $4$

$l-1$, and achieves feature map at layer $l-1$, $F^{l-1} = \mathcal{E}_1^{l-1}(I)$. It will further feed to the encoder $\mathcal{E}_{l-1}^l$ involving encoder layer $l$, and get feature $F^l = \mathcal{E}_{l-1}^l(F^{l-1})$. Here, our loss includes two folds. On one hand, feature $F^l$ would directly feed to the decoder for image reconstruction, namely $\hat{I} = \mathcal{D}_1^{l-1}(\hat{F}^{l-1})$, where $\hat{F}^{l-1} = \mathcal{D}_{l-1}^l(F^l)$. We use $L_2$-norm based image reconstruction loss, $\mathcal{L}_{img} = ||\hat{I} - I||^2$. On the other hand, we measure feature loss by $\mathcal{L}_{feat} = ||\hat{F}^{l-1} - F^{l-1}||^2$. In summary, the decoder sub-net $\mathcal{D}_{l-1}^l$ only involving decoder layer $l$, can be achieved by $\min(\mathcal{L}_{img} + \mathcal{L}_{feat})$, which can be rewritten as: $\min(||I - \mathcal{D}_1^{l-1}(\mathcal{D}_{l-1}^l(F^l))||^2 + ||F^{l-1} - \mathcal{D}_{l-1}^l(F^l)||^2)$, where the $\mathcal{D}_1^{l-1}$ is fixed when $\mathcal{D}_{l-1}^l$ is computed in our training. By contrast, both $\mathcal{D}_1^{l-1}$ and $\mathcal{D}_{l-1}^l$ are always retrained for every decoder layer [30]. We train our decoder on the ImageNet dataset [25].

## 5. Ablation Study

### 5.1. Style Loss Analysis

We study the effect of different style loss terms, including global style loss (in Eq. (4)), local style loss (in Eq. (6)), and reshuffle style loss (in Eq. (9)), by neglecting the common content loss. Thus, we only evaluate them on texture synthesis. We collect 60 image pairs from existing papers, and start from random noise to respectively minimize the three losses in layer $4$. The optimization is conducted by L-BFGS method, and stopped at 500 iterations, where results have no visible changes with further iterations.

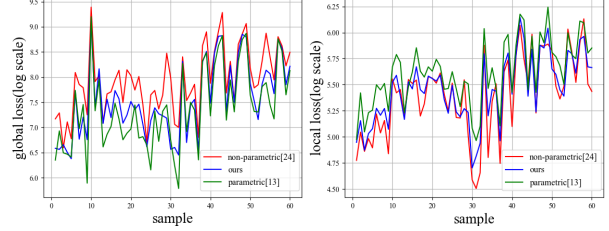Some results are shown in Fig. 6. As we can see, the



Figure 7. Global style loss (Eq. (4)) (left) and local style loss (Eq. (4)) (right) measured on every result with various methods.
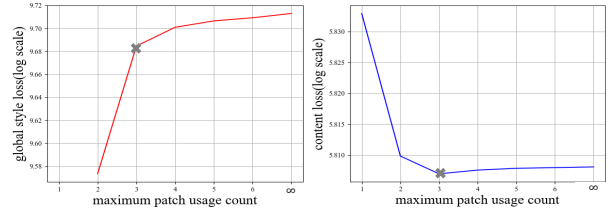


Figure 8. The global style loss (left) increases while content loss (right) decreases with the increase of maximum patch usage count.

results by minimizing the global style loss (e.g., [15]) better reproduce the overall feeling of the style, while the results by minimizing the local style loss (e.g., [28]) are more faithful to the local shapes of example texture. The results using our reshuffle loss own both merits: global similarity and local plausibility. The quantitative comparison also demonstrate the same point, as shown in Fig. 7. We can see that our loss function achieves lower global loss (Eq. (4)) than nonparametric method; while obtains lower local loss (Eq. (6)) than parametric method, in all the cases.

### 5.2. Count of Patch Usage

We will examine the effect of patch usage count. The relaxation of maximum usage count only reduces global style loss, but the local style loss always remains 0. We still consider the above 60 examples, and use our single layer optimization (in Section 4.2, only layer $4$) to compute the content loss (Eq. (2)) and the global style loss (Eq. (4)). Here, we try varied maximum counts of patch usage, and shows their corresponding content losses and global style losses in Fig. 8. With the increase of maximum usage count, the global style loss increases while the content loss decreases. It is not hard to understand that allowing more usage times will provide more choices to match content. Moreover, we find the good upper bound for patch usage is 3, with minimum style loss given the best preservation to content, which helps infer a tradeoff weight $\lambda = 0.05$ in Eq. (10). A visual comparison of different usage constraints is shown in Fig. 9.

### 5.3. Patch Size Selection

Another hyperparameter in our method is the patch size $R$. Increasing patch-size will sacrifice the global style loss (in Eq. (4)) to some extend. However, in some scenarios, large patch size is needed to preserve spatial coherence. As
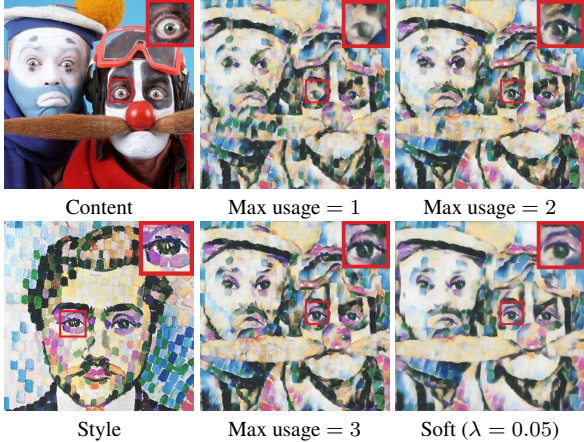
Figure 9. A comparison of results with different usage count constraints. Please note the eye in the red rectangles are miss-matched when max usage = 1, but fixed when max usage increases.
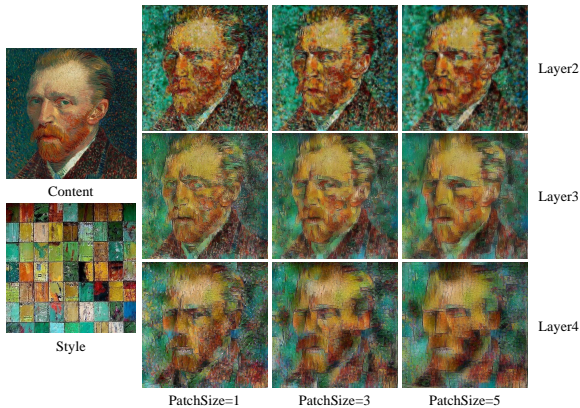


Figure 10. An example of style transfer results with different patch sizes and at different layers.

shown in Fig. 10, lager patch size we use, better local structure of style patterns can be preserved. Another factor is that patches in coarser layers may cover larger reception fields, making good matching difficult. So we choose empirically $3 \times 3$ patch in layer $4$ and $5 \times 5$ patch in layers $2, 3$.

# 6. Results

## 6.1. Implementation Details

All results are produced by multi-layer aggressive optimization in feature domain (in Section 4.3). On each layer, we do 5 EM iterations. The used patch sizes are $\{R^l\}_{l=2,3,4} = \{5, 5, 3\}$ and the patch usage parameter is set to be $\lambda = 0.05$ respectively according to the experiment in Section 5. As to the weight for content and style balancing, we set $\alpha = 0.5$ and $\beta = 1.0$ to make our stylization level similar to previous works.

## 6.2. Comparisons

We compare our result with other neural style transfer methods including representative parametric methods [16, 20, 30] and non-parametric methods [28, 32]. For fair comparison, all our results are generated with fixed parameters as described in Section 6.1. And theirs are obtained by running author-released code with default settings. We have tested on more than 100 content and style pairs collected from previous papers and Ostagram website [1]. Fig. 11 shows some representative results. More results can be found in our supplemental material.

As shown in Fig. 11, our method shares advantages from both kinds of methods. First, as a non-parametric method, our results preserve the local texture patterns better, more faithful to the style, compared with parametric methods [16, 20, 30], as shown in row 1&2 of Fig. 11. Unfortunately, in their results, the local textures are distorted, and some new patterns (not belongs to the style) appears. Second, our method seeks for best matches for each local patch in the content, so it can better achieve semantic-level transfer (*e.g.*, eye-to-eye, mouth-to-mouth) than parametric methods which only mimic the global statistics [16, 20, 30], as shown in row 3&4 of Fig. 11. Third, our method can own good global properties, making it different from other non-parametric methods [28, 32]. On one hand, our result can better preserve overall feels of the exemplar style, as shown in the row 5&6 of Fig. 11. Non-balanced neural patch sampling [28, 32] makes their results different from the global distribution of the style patterns; while ours are globally more faithful. On the other hand, our method can successfully avoid excessively repetitive use of the same sample, which will cause the washout effect [21]. We can clearly see these undesired effects in row 2&7 of Liao et al.'s results [32]. The property benefits from our reshuffle constraint.
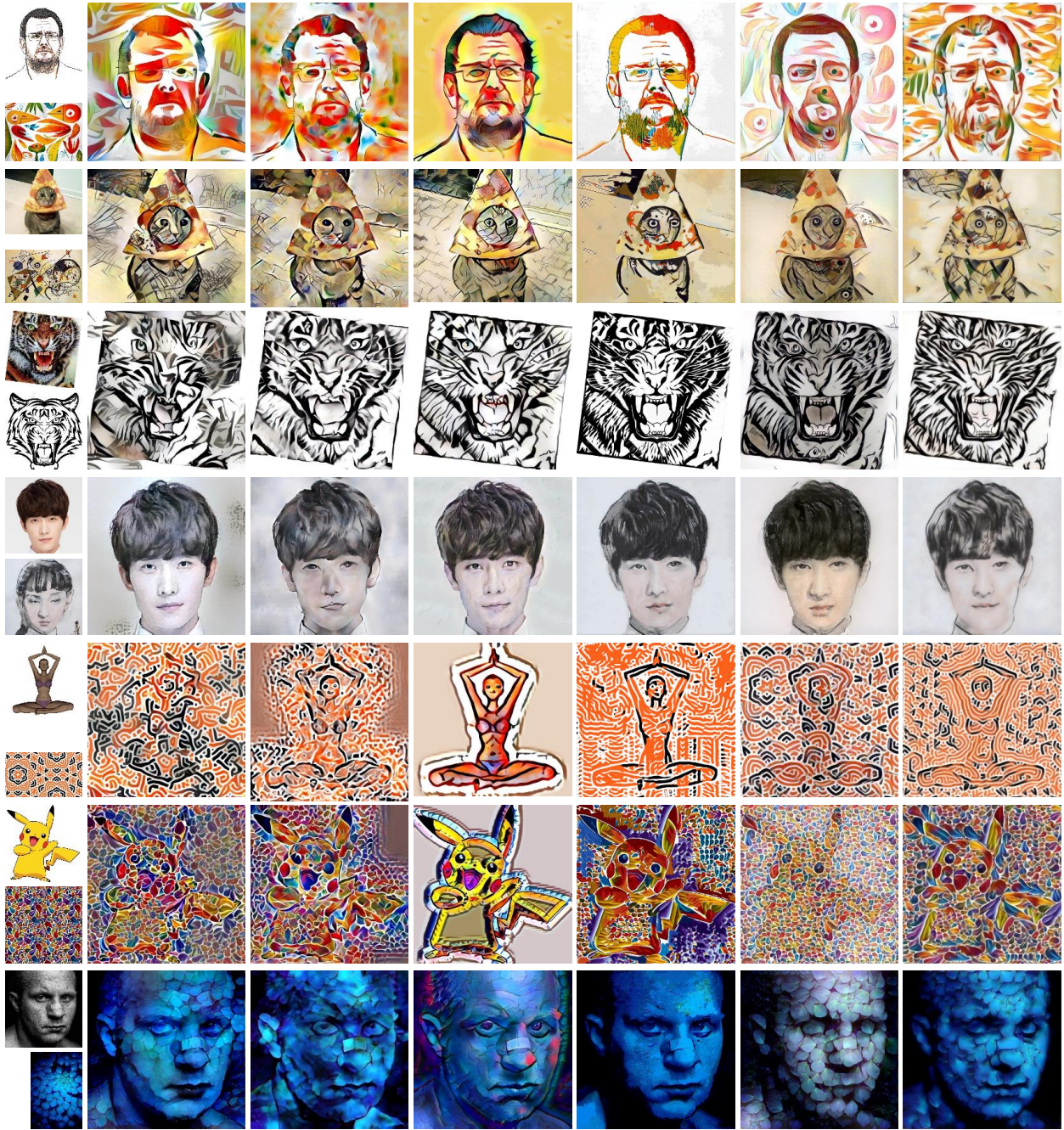
We also compare the time cost of all these methods. Tabel below Fig. 11 gives the average running time of each method on $512 \times 512$ image pairs. All the methods are tested on a PC with an Intel E5 2.6GHz CPU and an NVIDIA Tesla K40C GPU. Our method is slower than [20, 30], comparable to [32] and faster than [16, 28]. The bottleneck is the constrained NN field search step.

## 6.3. Perceptual Study

We conduct a user study similar to [22]. In the study, we use 150 groups of images shown in the supplemental material. Each group contains two inputs and six outputs (involving 5 results from [16, 30, 20, 32, 28], and ours). All six results in each group are presented side-by-side and in a random order to participants. Participants are given unlimited time to rank the score from 1 to 6 (1 is the best, 6 is the worst) according to preference. We show the average ranking scores over 15 participants in Table 1. Overall, subjects prefer our result more than others.
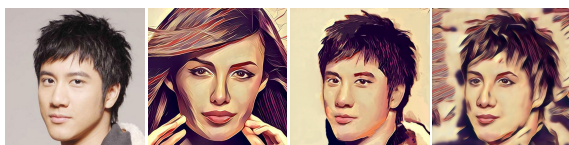
Table 1. Average stylization rank scores of six algorithms

| Method | [16] | [30] | [20] | [32] | [28] | Ours |
|---|---|---|---|---|---|---|
| Average rank | 3.08 | 3.55 | 3.92 | 3.42 | 4.2 | 2.88 |

| Method | Gatys et al.[16] | Li et al.[30] | Huang et al.[20] | Liao et al.[32] | Li et al.[28] | Ours |
|--------|------------------|---------------|------------------|-----------------|---------------|------|
| Time(s) | 370 | 8.46 | 0.556 | 114 | 195 | 114 |

Figure 11. Comparison with previous neural style transfer methods.



| Content | Style | Liao et al. [32] | Ours |

Figure 12. A failure case.

## 7. Discussion and Conclusion

Despite the success of neural style transfer, the relationship between different methods was far from clear. In this paper, we give a new perspective to connect them. We then propose a novel and simple idea, called deep feature reshuffle, which is the first to unify both commonly-used global

and local style losses. Based on this idea, we propose a new and efficient neural style transfer algorithm by progressively optimizing the new loss in feature domain. The results have shown that our approach is widely applicable to various inputs, and produces better quality than existing methods.

However, the method still suffers from some limits. Constraining the usages of neural patch for the sake of style, will cause less accurate matching and thus damage the content structure, as shown in Fig. 12. It can be solved by fine-tuning the usage parameter $\lambda$. How to automatically determine the optimal parameter for each input will become a vital and practical problem to be explored in future work.

# References

[1] Ostgram. https://ostagram.ru. 7

[2] N. Ashikhmin. Fast texture transfer. *IEEE Computer Graphics and Applications*, 23(4):38–43, 2003. 2

[3] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.*, 28(3):24–1, 2009. 4

[4] A. J. Champandard. Semantic style transfer and turning two-bit doodles into fine artworks. *arXiv preprint arXiv:1603.01768*, 2016. 2

[5] D. Chen, J. Liao, L. Yuan, N. Yu, and G. Hua. Coherent online video style transfer. *arXiv preprint arXiv:1703.09211*, 2017. 2

[6] D. Chen, L. Yuan, J. Liao, N. Yu, and G. Hua. Stylebank: An explicit representation for neural image style transfer. In *Proc. CVPR*, 2017. 2

[7] D. Chen, L. Yuan, J. Liao, N. Yu, and G. Hua. Stereoscopic neural style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 10, 2018. 2

[8] T. Q. Chen and M. Schmidt. Fast patch-based style transfer of arbitrary style. In *Proc. of NIPS*, 2016. 1, 2

[9] J. S. De Bonet. Multiresolution sampling procedure for analysis and synthesis of texture images. In *Proc. ACM SIGGRAPH*, pages 361–368. ACM Press/Addison-Wesley Publishing Co., 1997. 1, 2

[10] V. Dumoulin, J. Shlens, and M. Kudlur. A learned representation for artistic style. *arXiv preprint arXiv:1610.07629*, 2016. 1, 2

[11] A. A. Efros and W. T. Freeman. Image quilting for texture synthesis and transfer. In *Proc. ACM SIGGRAPH*, pages 341–346. ACM, 2001. 1, 2

[12] A. A. Efros and T. K. Leung. Texture synthesis by non-parametric sampling. In *Proc. ICCV*, volume 2, pages 1033–1038. IEEE, 1999. 1, 2

[13] M. Elad and P. Milanfar. Style-transfer via texture-synthesis. *arXiv preprint arXiv:1609.03057*, 2016. 1

[14] J. Fišer, O. Jamriška, M. Lukáč, E. Shechtman, P. Asente, J. Lu, and D. Sỳkora. Stylit: illumination-guided example-based stylization of 3d renderings. *ACM Transactions on Graphics (TOG)*, 35(4):92, 2016. 4

[15] L. Gatys, A. S. Ecker, and M. Bethge. Texture synthesis using convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 262–270, 2015. 3, 6

[16] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015. 1, 2, 3, 4, 6, 7, 8

[17] L. A. Gatys, A. S. Ecker, M. Bethge, A. Hertzmann, and E. Shechtman. Controlling perceptual factors in neural style transfer. *arXiv preprint arXiv:1611.07865*, 2016. 2

[18] D. J. Heeger and J. R. Bergen. Pyramid-based texture analysis/synthesis. In *Proc. ACM SIGGRAPH*, pages 229–238. ACM, 1995. 1, 2

[19] A. Hertzmann, C. E. Jacobs, N. Oliver, B. Curless, and D. H. Salesin. Image analogies. In *Proc. ACM SIGGRAPH*, pages 327–340. ACM, 2001. 2

[20] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. 2017. 1, 2, 5, 7, 8

[21] O. Jamriška, J. Fišer, P. Asente, J. Lu, E. Shechtman, and D. Sỳkora. Lazyfluids: Appearance transfer for fluid animations. *ACM Trans. Graph. (Proc. of SIGGRAPH)*, 34(4):92, 2015. 1, 7

[22] Y. Jing, Y. Yang, Z. Feng, J. Ye, and M. Song. Neural style transfer: A review. *arXiv preprint arXiv:1705.04058*, 2017. 7

[23] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. *arXiv preprint arXiv:1603.08155*, 2016. 1, 2

[24] A. Kaspar, B. Neubert, D. Lischinski, M. Pauly, and J. Kopf. Self tuning texture optimization. In *Computer Graphics Forum*, volume 34, pages 349–359. Wiley Online Library, 2015. 2, 4, 5

[25] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012. 1, 6

[26] J. E. Kyprianidis, J. Collomosse, T. Wang, and T. Isenberg. State of the" art: A taxonomy of artistic stylization techniques for images and video. *IEEE transactions on visualization and computer graphics*, 19(5):866–885, 2013. 1

[27] H. Lee, S. Seo, S. Ryoo, and K. Yoon. Directional texture transfer. In *Proc. of NPAR*, pages 43–48. ACM, 2010. 2

[28] C. Li and M. Wand. Combining markov random fields and convolutional neural networks for image synthesis. *arXiv preprint arXiv:1601.04589*, 2016. 1, 2, 3, 4, 6, 7, 8

[29] C. Li and M. Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. *arXiv preprint arXiv:1604.04382*, 2016. 1, 2

[30] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Universal style transfer via feature transforms. 2017. 1, 2, 5, 6, 7, 8

[31] Y. Li, N. Wang, J. Liu, and X. Hou. Demystifying neural style transfer. *arXiv preprint arXiv:1701.01036*, 2017. 3

[32] J. Liao, Y. Yao, L. Yuan, G. Hua, and S. B. Kang. Visual attribute transfer through deep image analogy. *ACM Trans. Graph. (Proc. of SIGGRAPH)*, 2017. 1, 2, 3, 7, 8

[33] M. Lu, H. Zhao, A. Yao, F. Xu, Y. Chen, and L. Zhang. De-coder network over lightweight reconstructed feature for fast semantic style transfer. In *Proc. ICCV*, 2017. 1, 2, 5

[34] F. Luan, S. Paris, E. Shechtman, and K. Bala. Deep photo style transfer. *CoRR, abs/1703.07511*, 2017. 2

[35] J. Portilla and E. P. Simoncelli. A parametric texture model based on joint statistics of complex wavelet coefficients. *International Journal of Computer Vision*, 40(1):49–70, 2000. 1, 2

[36] A. Selim, M. Elgharib, and L. Doyle. Painting style transfer for head portraits using convolutional neural networks. *ACM Trans. Graph. (Proc. of SIGGRAPH)*, 35(4):129, 2016. 2

[37] O. Sendik and D. Cohen-Or. Deep correlations for texture synthesis. *ACM Trans. Graph.*, 36(5):161, 2017. 2

[38] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2

[39] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. *arXiv preprint arXiv:1603.03417*, 2016. 2

[40] L.-Y. Wei and M. Levoy. Fast texture synthesis using tree-structured vector quantization. In *Proc. ACM SIGGRAPH*, pages 479–488. ACM, 2000. 2

[41] P. Wilmot, E. Risser, and C. Barnes. Stable and controllable neural texture synthesis and style transfer using histogram losses. *arXiv preprint arXiv:1701.08893*, 2017. 1

[42] Q. Wu and Y. Yu. Feature matching and deformation for texture synthesis. *ACM Trans. Graph.*, 23(3):364–367, 2004. 2

[43] C. Zhu, R. H. Byrd, P. Lu, and J. Nocedal. Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale bound-constrained optimization. *ACM Transactions on Mathematical Software (TOMS)*, 23(4):550–560, 1997. 4