# Geometric Style Transfer



Figure 1: Left to right: content image, style image, texture only transfer using Gatys *et al*. [14], geometric and texture transfer using our method. Our method can not only capture texture features of the style image, but also deform the content image to match geometric structures of the style image.

Xiao-Chang Liu     Xuan-Yi Li    Ming-Ming Cheng     Peter Hall

University of Bath      Nankai University      University of Bath

## Abstract

*Neural style transfer (NST), where an input image is rendered in the style of another image, has been a topic of considerable progress in recent years. Research over that time has been dominated by transferring aspects of color and texture, yet these factors are only one component of style. Other factors of style include composition, the projection system used, and the way in which artists warp and bend objects. Our contribution is to introduce a neural architecture that supports transfer of geometric style. Unlike recent work in this area, we are unique in being general in that we are not restricted by semantic content. This new architecture runs prior to a network that transfers texture style, enabling us to transfer texture to a warped image. This form of network supports a second novelty: we extend the NST input paradigm. Users can input content/style pair as is common, or they can chose to input a content/texture-style/geometry-style triple. This three image input paradigm divides style into two parts and so provides significantly greater versatility to the output we can produce. We provide user studies that show the quality of our output, and quantify the importance of geometric style transfer to style recognition by humans.*

## 1. Introduction

Neural style transfer (NST) is an active area of research with the aim of synthesising artistic images. The most common paradigm is to input two images, one provides the content that the output should contain, the other input indicates the "style" in which the provided content is to be rendered. To date, NST has been dominated by the transfer of texture, but artistic style is not characterised by texture alone. Style includes changes of shape of objects, rules of composition, the projection model used, and many other factors.

Our contribution is to step closer to artistic styles by including *geometric style transfer* (GST) that changes the shape of the content image to better match shapes in the style image. Figure 1 demonstrates that transfer of geometric style yields an output that is a closer match to the style image that can be achieved by texture alone. This paper appeals to a little of the Art History literature to argue for the importance of GST, explains how it can be achieved in a general setting, and provides empirical evidence that GST yields output closer to the target style than texture transfer alone.

The style of an artist, or a school of artists, is only partly characterized by the way they make marks on the image surface. More than sixty years ago, Art Theorist Rudolf Arnheim [2] argued that art style should be described by attributes such as color, shape, and composition. Just over forty years later another Art Theorist, John Willats, was also concerned with art style and coined the terms *projection style* and *denotation style*. Denotation refers to the way in which an artist makes marks on the image surface. Denotation is impacted by the substrate (paper, canvas, *etc*.), the media (paint, pencil, *etc*.) and the application method. Projection refers to the spatial organisation of parts [52], it includes both standard cameras and orthogonal cameras, but more generally refers to the spatial organization of an objects parts. Willats shows that projection variety is at least if not more important than denotation variety in characterizing style. For example, ancient Egyptian art is characterised by the way people are unnaturally posed; Byzantine artists routinely used inverse perspective; Chinese artists traditionally used orthogonal pro-

jection; Cubism exhibits a multitude of views in a single image. Denotation is important, but ancient Egyptian art (say) is recognizable in paintings, in bas-relief, in sculpture – all of which differ in denotation but share the same "projection" variety.

It is worth noting that the geometric changes humans introduce tend not to be arbitrary. Some artists deform shape to bring emphasis to some aspect of the subject being depicted. For example, Stubbs would deliberately paint bulls to be larger and stronger than any real bull could be – he did so to please the landowners whose animals he was depicting. Other artists, such as Modigliani and El Greco, distort faces and humans as a matter of personal style. We can imagine that there is some underlying photograph that has been somehow warped, and then painted over.

NST has been dominated by the transfer of texture, which approximates denotation. There are some examples of *geometric style* transfer in NST, but these require strong models of faces (*e.g.* [55, 49]) or of text [54] so have a limited content domain. Our contribution is to introduce *geometric style* transfer as a general case to sit alongside texture transfer. Our method (see Section 3) conforms to the standard paradigm of providing a content and a style image, but in our case changes and distortions of geometry are transferred in addition to texture. We can unique do style transfer with three inputs (content source, style source and geometric source). Processing the geometric changes requires an additional processing path that is parallel to an otherwise standard NST architecture. The additional path is used to compute a geometric mapping that warps the content image before it texture transfer takes place. Section 7 shows that our results not only tend to be preferred to the output of alternative algorithms, but also then to be regarded as more similar to the target output.

## 2. Related Work and Background

Image stylization is the process of mapping an input image in one style to an output image in a new style, with content preserved. The problem has been a significant field of research within Visual Computing for over two decades, beginning with non-photorealistic rendering (NPR) and more recently continuing with neural style transfer (NST).

NPR has been the subject of research for many years and is too large to give a comprehensive overview. It includes image synthesis from 3D models, the emulation of substrate and media, and user interaction, but we confine ourselves to NPR from images. Early algorithms mapped pixel patches into "blobs" [17]. Later, the mapping became more sophisticated, targeting salient regions *e.g.* [11]; blobs became brush strokes [22]. These few examples are indicative of a much larger body of work in which image stylization is seen as a sophisticated filtering process. Higher forms of abstraction are far less common, but have been tackled using *ad-hoc* approaches emulating movements such as Cubism [10] and artists such as Archimboldo [25]. Projection style in the sense Willats intends has been wholly neglected in NST, though there are examples in NPR such as [58, 19].

Some of these early algorithms contained stochastic elements but were all prescriptive in the sense that the style of the output was predefined. Examples of learning style appeared as early as 1998 [18], and later in 2001 [23].

Image stylization moved firmly in the direction of learning, in about 2015, when Gatys *et al.* [13] introduced neural style transfer. The key idea was to adjust a variable image $X$ so that it matched some image $A$ for content and another image $B$ for style. The definition for content loss and style loss were both premised on features extracted from a network pre-trained for recognition (VGG-16 was used, [50]), the loss for content being the L2-norm between response vectors, and the style loss being the L2-norm between Gram matrices comprising feature correlations.

The ability to learn style transforms is useful, but slow optimization motivated work towards fast transfer [27, 29, 26, 36]. A second problem is the need to retrain the network for each new style, which encouraged work to learn styles more generally, including but not limited to [38, 53, 7]. The loss functions have received attention, Huang *et al.* [26] provide one example in which the loss function is based on the statistical distribution of features; Li *et al.* [38] are another – they use a whitening and coloring transforms to better map feature vectors. Kotovenko *et al.* [33] introduce two additional losses that learn subtle variations within one style and ensure stylization is not conditioned on the input photograph.

To date, NST has been extended to do many different tasks [28], such as portrait painting style transfer [48, 55]; visual attribute transfer [40, 32, 56, 34]; semantic style transfer in natural images [45, 9, 5]; video style transfer [47, 24, 16, 6, 36]; 3D style transfer [8, 30], and photorealistic style transfer [42, 44, 39, 57].

Nearly all NST is limited in the sense that there is no explicit attempt to change the geometry or shape of objects in the picture. For clarity, many NST algorithms can output images with a different geometry to the input, but any such changes are accidents of the texture transfer process and are typically confined to the boundaries of objects. There is usually no effort to deliberately transfer any geometric distortions introduced by an artist. This limits the ability of NST algorithms to emulate style, because geometric changes are an integral part of style.

The need to transfer *geometric* style has been recognised within NST, albeit in specialised domains. Facial caricature is relatively popular [59, 35, 4, 49], while Yaniv *et al.* [55] consider artistic portraiture more generally; all of that work is limited to faces. Yang *et al.* [54] explicitly control the shape deformation of artistic text. Our work is unique by being the first to provide a general approach to geometric style transfer.

In summary, NST research has largely followed the trajectory of NPR in that work began on texture and later moved towards high forms of abstraction. If the history of NPR is a teacher, then NST will continue to develop away from texture transfer.

## 3. Texture and Geometric Style Transfer

Our system transfers both texture and geometric styles, the latter being our contribution. As is usual for neural style transfer, our system can take two inputs: a content image $I^c$ and a style image $I^s$. Our novelty is this: rather than transfer texture directly, the content image is warped first. The warp carries the content image onto the style target, so the output image $I^o$ is the same size as $I^s$. The warp transfers geometric style. Figure 9 illustrates by showing photo-textured warped images and the textured results that come from them.

Our system is novel too in being able to accept three inputs: one content image as before, one geometric style image, and one texture style image. This paper is written assuming two inputs for familiarity and simplicity of explanation – the extension to three inputs simply use the geometric style image through the geometric warping network, and the texture style image through the texture transfer network.

The whole procedure consists of three major steps: feature extraction, geometric warping, and texture transfer. The first step is to extract features from a network, as explained in Section 4. As is common in style transfer, we need "content features" to preserve content, "denotation" features that will be used for denotation (texture) transfer, and uniquely we need "geometric features" for geometry transfer. The second step sends these geometric features into a CNN architecture to compute a mapping $\Re^2 \mapsto \Re^2$; the mapping is used to warp the content image, see Section 5. The third step uses the content and texture features in an Image-Optimization-Based online style transfer method with a multi-scale strategy to transfer texture while preserving content to generate the final result, as described in Section 6.

## 4. Feature Extraction

We extract all our features from VGG-19, which is trained on more than a million images from the ImageNet dataset [12] and can classify images into a thousand or more object categories. A given input image is encoded in each layer by the filter responses to that image, and layer $l$ with $N_l$ filters has $N_l$ distinct feature maps each of size $W_l \times H_l$. The extraction procedure is summarised in Figure 2.

**Content Features** are used to preserve content during transfer. We draw them from the intermediate layer of the network, because such layers contain mid to high level image representations. Specifically, we use the $N_4$-element feature vectors from a $W_4 \times H_4$ array at at layer $conv4\_2$ to obtain feature maps $F^c$ for the content image $I^c$ and $F^o$ for the output image. These maps are each of size $W_4 \times H_4 \times N_4$,

**Texture Features** are used to transfer denotation style (which prior work refers to as "style", with no modifying adjective). Low-level statistics tend to characterize denotation, but these benefit from context. Following [14, 29], we use early $conv1\_1$, $conv2\_1$, and later layer $conv3\_1$, $conv4\_1$ and $conv5\_1$, and compute feature correlations as Gram matrices
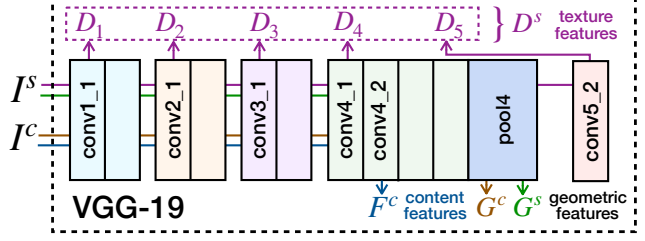


Figure 2: Feature extraction.

$D_l \in \Re^{N_l \times N_l}$ for each layer $l$. Taken over all layers we have a set $D = \{D_l : \forall l \in [1, 2, 3, 4, 5]\}$. This set of correlations characterizes the texture style of an image; $D^s$ for the style image $I^s$ and $D^o$ for the output image.

**Geometric Features** are used to compute a spatial mapping $\mathcal{M} : \Re^2 \mapsto \Re^2$. The features should be reasonably robust to local local spatial variability in the form of rotation and translation, and aggregate content. We use the feature map at $pool4$ layer, followed by L2-normalization of each feature channel, and get geometric features $G^c$ and $G^s$ for $I^c$ and $I^s$ respectively.

## 5. Geometric Style Transfer

Geometric style transfer warps the content image onto the style image before texture is transferred. GST comprises of two parts, as illustrated in Figure 3. First, feature correlation measures the degree to which geometric features (see Section 4) in the content and style image correlate. Second, the correlations are input to a trained network that provides a first approximation to a geometric mapping $\mathcal{M} : \Re^2 \mapsto \Re^2$; this mapping is governed by geometric style. We use parametric spatial mappings; we've found either affine transforms or quadratic thin-plate to be sufficient but there is no reason in principle not to use other mappings. Each part is described in detail below.

### 5.1. Feature Correlation

Feature correlation is between the geometric feature sets, $G^c$ and $G^s$ defined in Section 4. These are each arrays of size $W_4 \times H_4$ and contain vector elements of length $N_4$, Elements in these sets are indexed by their location in the sample array. Then correlation function is a four-dimensional function:

$$C_{i,j,k,l} = \hat{f}^c_{i,j} \odot \hat{f}^s_{k,l}, \tag{1}$$

with $\odot$ being the inner product between vectors, and $\hat{\ }$ indicating $L_2$ normalization. $C$ is then postprocessed to zero out negative values. A visualization of this procedure is illustrated in Figure 3.

### 5.2. Learning a Spatial Mapping

To transfer geometric style, we require a spatial transform $\mathcal{T} : \Re^2 \mapsto \Re^2$ to warp the content image to the style image. We assume a parametric transform, and therefore use a regression CNN to determine parameter values, $\Theta$. At this stage we
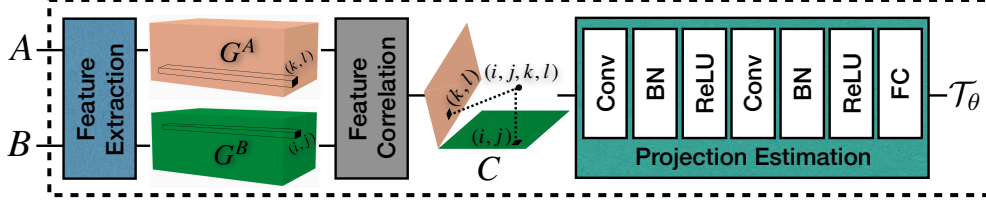
Figure 3: Geometric warping Network; it runs parallel to the texture transfer network.

are computing a first approximation to the spatial transform and have found an affine transform; in the next step will will upgrade to a thin-plate spline. Whether we use these or some other transform, we wish to compute the mapping parameters using the correlation matrix, $C$, so that $\Theta = \mu(C)$. To do this, we use a regression CNN to learn the mapping $\mu : \Re^{W \times H \times W \times H} \mapsto \Re^p$ where $p$ is the number of parameters, and $W,H$ control the grid of sample sites.

The regression CNN is trained by iterating over many identical trials. At the $i^{th}$ trial we randomly sample transform parameters $\Theta_i$, sampling details are given below. The parameters specify a spatial transform we will call $\mathcal{T}_i$. Each ground truth transform is used to warp training images $A$ to make $B = \mathcal{T}_i(A)$, and also to move sample locations for geometric features, $(x,y)_{wh}$, to get sample locations in the warped image at $\mathcal{T}_i(x,y)_{wh}$. The image $B$ is then processed using Gatys *et al.* [14] to create a artistic-texture-augmented copy. We can now use the original image $A$ and the warped and artistic-texture augmented image $B$ to compute a correlation matrix as described by Equation 1; we will call this correlation $C_i$. This process explicitly connects the known parameters $\Theta_i$ to a known correlation matrix $C_i$. The quality of the (current) mapping $\mu$ is measured using the L2-norm between the sample locations mapped under the known transform, $\mathcal{T}_i$, and the transform constructed from the parameters $\mu(C_i)$, which we write at $\mathcal{T}_{\mu|C_i}$. Thus, the regression net has the loss function:

$$\mathcal{L} = \sum_i ||\mathcal{T}_{\mu|C_i}(x_{jk}), \mathcal{T}_i(x_{jk})||_2. \tag{2}$$

Once trained, the network will compute spatial transform parameters, given a correlation matrix $C_{ijkl}$.

The learning process above will work in principle for any transform. Even so, finding the optimal thin-plate spline (TPS) transformation [3] is not easy. We have found it useful in practice to learn such higher-order mappings by first estimating an affine mapping using the above, and using this affine mapping to warp the image $A$ before then warping it a second time using the higher-order mapping that is being learned. Note that this requires two copies of the geometric warping network: the first outputs the 6 parameters of an affine mapping, the second outputs the 18 parameters of a quadratic thin-plate spline. The two networks are distinct.

**Comments on Geometric Warping:** The examples in this

paper all use either affine or TPS warping. However, there is nothing about the architecture that limits it to that pair of warping families. We could have used bicubic warps, or a homography, for example. The geometric warping network could sit in parallel to many existing texture transfer architectures. The reader is free to implement our network alongside their own, and to explore different families of geometric warp.

## 6. Texture Transfer

Texture style transfer is used to approximate denotational style to input images. In our case these will be a content image that has been warped by the geometry style network of the previous section. Let $I^o$ be the generated stylized result, and $I^s$ and $I^c$ are the style and content image respectively; the content image has been warped to match geometric style. As stated in Section 4, the style of an image is represented by the texture feature $D$. The texture style reconstruction loss is a weighted sum of L2-norms:

$$\mathcal{L}_{texture}(I^s, I^o) = \frac{1}{2} \sum_l \omega_l \| D_l^s - D_l^o \|_2, \tag{3}$$

where $l \in \{1,2,3,4,5\}$, stands for the selected layers, and $\omega_l$ is the weighting factor for each layer.

The content of an image is represented by the content feature $F$ (see Section 4), and the content reconstruction loss is the L2-norm of two features:

$$\mathcal{L}_{content}(I^c, I^o) = \frac{1}{2} \| F^c - F^o \|_2. \tag{4}$$

The loss function we minimise is:

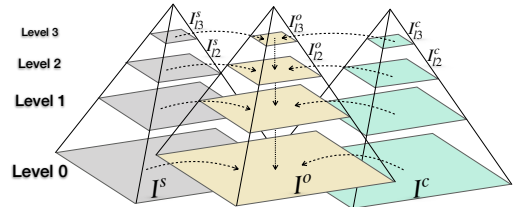$$\mathcal{L}_{total} = \alpha \mathcal{L}_{texture}(I^s, I^o) + \beta \mathcal{L}_{content}(I^c, I^o), \tag{5}$$



Figure 4: Multi-scale strategy used in style texture transfer.

4

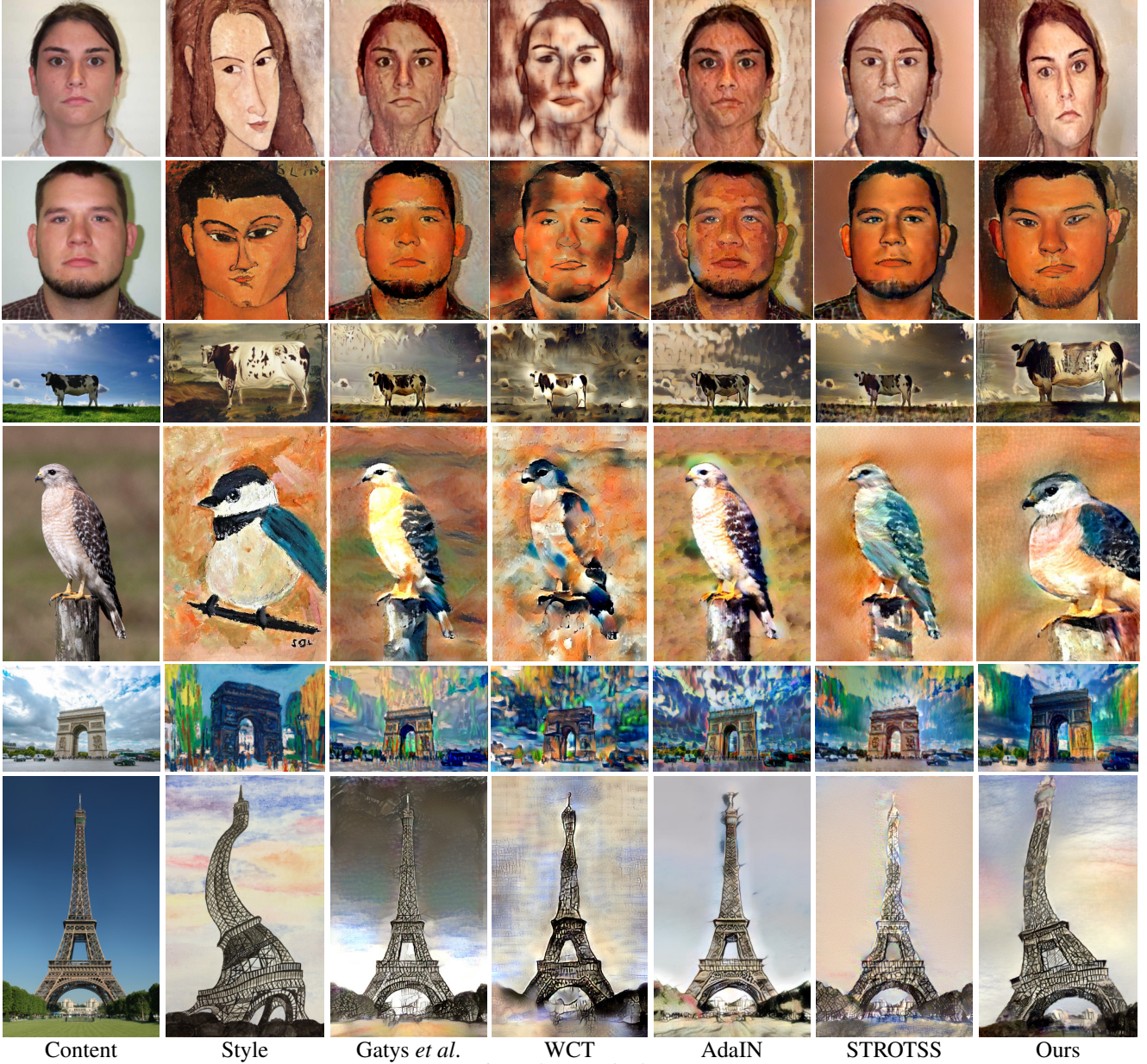| Content | Style | Gatys *et al.* | WCT | AdaIN | STROTSS | Ours |

Figure 5: Comparison to related methods. Gatys *et al.* [14] and STROTSS [32] are image optimization based, AdaIN [26] and WCT [38] are feature transformation based.

where $\alpha$ and $\beta$ are the weighting factors.

**Multi-Scale Strategy:** After geometric style transfer (Section 5), some parts of the image will be enlarged in some cases (see Figure 10), and the resolution will decrease accordingly. Furthermore, as pointed out in previous work [43], the effective receptive field of network neurons is fixed and relatively small, which limits the scale of synthesized features. In order to avoid the stylized effect being affected by the decreased image resolution we adopt a multi-scale strategy to do style transfer, inspired by prior art [21, 20, 51, 15] that shows good texture synthesis with a bank of multi-scale filters. Specifically, as shown in Figure 4, we first downsample $I^s$ and $I^c$ by feeding them into a Gaussian pyramid [1] (here we use 3 levels), and use $I^s_{l3}$ and $I^c_{l3}$ to perform stylization and get $I^o_{l3}$. Then we upsample $I^o_{l3}$ as a initialization, utilize $I^s_{l2}$ and $I^c_{l2}$ to generate $I^o_{l2}$, and so on. In this way, during the stylization, the optimization will simultaneously match features at every pyramid layer, which guarantees the generation of high-resolution outputs (we further discuss this in Section 7.3).

Figure 6: Comparison on one content with multiple style images, the geometric styles will change with different style images.

# 7. Experimental Results

Our experimental results come in three parts: (i) qualitative results designed to show the reader the difference that geometry transfer makes to style transfer; (ii) quantitative experiments that place a subjective measure over the degree of difference; (iii) ablation studies to show the difference between affine and TPS geometric style transfer and of multi-scale texture transfer. Experiments and results are explained below, before which we provide implementation details.

Our implementation uses PyTorch [46]. We use the pre-trained VGG-19 exactly as described in [50]. The geometric warping network is trained on the Microsoft COCO dataset [41]. We resize each of the training images to $240 \times 240$ and train the network with a batch size of 8. We use Adam [31] with a learning rate of $1 \times 10^{-3}$; training takes roughly 8 hours on a single GTX 1080Ti GPU; Sample points $x_{jk}$ in Equation 2 are from a $20 \times 20$ uniform grid. For multi-scale texture transfer (Section 6), we weight each layer equally in Equation 3 ($\omega_l$=1/5), the ratio $\alpha/\beta$ in Equation 5 is $5 \times 10^{-3}$. Run-times are comparable with other image optimization methods, generating a $512 \times 512$ image takes around 50 seconds.

## 7.1. Qualitative Comparisons

We qualitatively compare our method with some closely related NST approaches. We could not compare with NST methods that deal with geometric transfer [55, 59, 35, 4, 49, 54] because (a) they each deal with one object class only (faces or text) and (b) they tend not to conform to the 'content/style' input paradigm. Instead, we compare with well known methods that like us input a single content photograph and a single style image, and which are intended to be general purpose. We compare to Gatys *et al.* [14] and STROTSS [32] which are

image optimization methods; and to AdaIN [26] and WCT [38] which are feature transformation methods. Results are shown in Figures 5 and 6.

From the comparison results, we should first notice that for other methods, the output sizes are the same as that of content images, while the size of our results match that of style images. Second, from the view of artistic effects, all the results keep texture features and color distributions well. The most striking difference between our output and that of all other algorithms is that only ours changes shapes within the content image. Our output portraits lengthen the face when necessary, skew facial features when that is part of the style, makes livestock larger, fattens birds, and bends towers and houses.

## 7.2. Quantitative Comparisons

Here we present quantitative results relating to the quality of output, the impact of geometric style, and computational efficiency.

**Output Quality:** There is no objective measure by which to assess the quality of outputs, therefore we followed others by conducting a questionnaire investigation to survey the preferences of different approaches. Every questionnaire included 10 pairs of content-style pairs, and participants are asked to vote for their favorite results. We collected questionnaires from 50 respondents and computed the percentage of every method with regard to the preferences. Results are shown in Figure 7. Our results are preferred more than any other, at 47% we are more than twice as likely as the next most popular, of 18%. However, popularity makes not statement about success in reaching the target style, our next experiment was designed to address this.

**Output Similarity:** The preferences of participants measure popularity but say nothing about the closeness of output to the
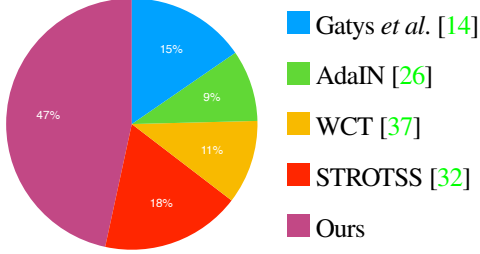
Figure 7: Participant preferences – our output is preferred more than twice as much as the nearest alternative.
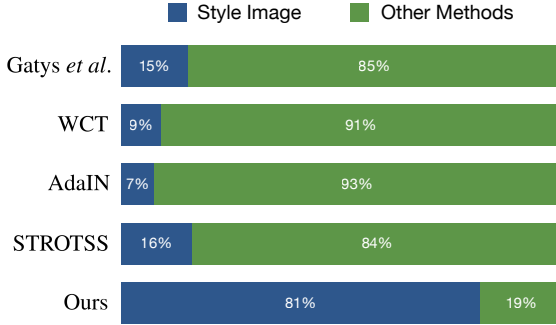


Figure 8: Results from our similarity experiment: output from four NST algorithms without GST are more likely to be assessed as more similar to each other than to the style image. Our GST output is judged as being closer to the style image.

target style. Asking questions about the similarity of output to the target style helps give us a handle on success in that regard. We showed 50 participants 3 images; one of the 3 was the style image, the other two images came from one of five algorithms, Gatys *et al.* [14], STROTSS [32] AdaIN [26], WCT [38], and ours. The three images, call them A,B,C were shown as three side-by-side pairs (A,B), (A,C), (B,C), each pair on a separate row. The choice of algorithm, the location of the images pairs and the ordering of a pair were all subject to randomization. We give each participant this simple instruction: *check the pair of images you think are most similar.* No other information was given to the participant nor did we ask any participant to explain their preference.

We recorded all preferences, and the number of times an image from an algorithm was picked. We collated this data into (a) the fractional number of times an output image from an algorithm was regarded as more similar to the style image, and (b) the fractional number of times the output images were regarded as more similar to each other. Figure 8 shows results in percentage form. Participants judged our output to be closer to the style target about 81% of the time, compared to at most 16% for any other. Furthermore, the other algorithms are more likely to be judged as more similar to each other than the style image. The standard deviation is about 4%.

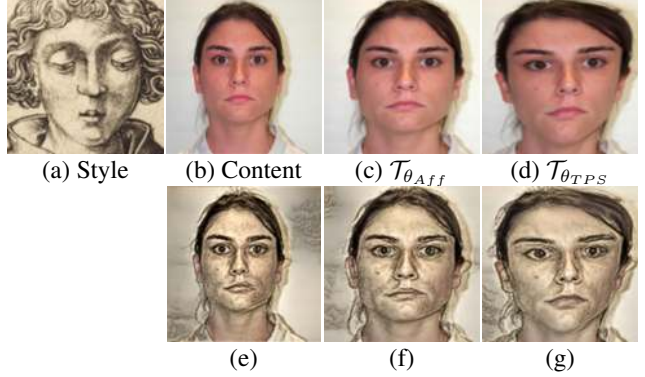The fact that the participants judge other algorithms' outputs



Figure 9: Effects of geometric warping. (a) style image, (b) content image, (c) warping only using an affine transformation $\mathcal{T}_{\theta_{Aff}}$, (d) warping with TPS transformations $\mathcal{T}_{\theta_{TPS}}$. (e) (f) and (g) are the texture transfer results of (b) (c) and (d) respectively. Affine transformation provides a rough warping but keeps some invariants (*e.g.* parallel lines, ratio of areas), TPS refines the transformation to a better match.

to be closer to each other than to the style image is important: the algorithms produce different textural output, so whatever criteria the participants used to assess similarity, it must have a stronger influence than texture. The fact our output was very much more likely to be chosen as most similar to the style target, strongly suggests that geometric style transfer explains the result.

## 7.3. Discussions and Ablation Studies

In this part we discuss two factors that highly influence the results: Geometric warping (Section 5.2) and Multi-scale strategy (Section 6), as well as the limitations of our method.

**Geometric Warping.** As stated in Section 5.2, the final estimated transformation either an affine transformation $\mathcal{T}_{\theta_{Aff}}$ or a thin-plate spline (TPS) transformation $\mathcal{T}_{\theta_{TPS}}$. Figure 9 illustrates the effectiveness of two transformations. We can see that affine transformation moves and scales the source image to roughly match the target, but the invariant properties of the affine transformation, (*e.g.* parallel lines, ratio of areas) prevent it from reaching a closer geometric mapping. The TPS refines the transformation to generate a better warping. As noted previously, there is nothing to prevent our approach reaching higher-order transforms, but we have not found it necessary.

**Multi-Scale Strategy.** This strategy was applied to improve the quality of the stylized results (Section 6). This is important to us, because sometimes the geometric alignment manipulation will enlarge parts of the image, which reduced image sharpness accordingly. This is seen in Figure 10: the output is not sufficiently sharp and some details are missing unless the multi-scale approach is used.

**Limitations and Interesting Cases:** Limiting assumptions are: (i) the content image and geometric style images share
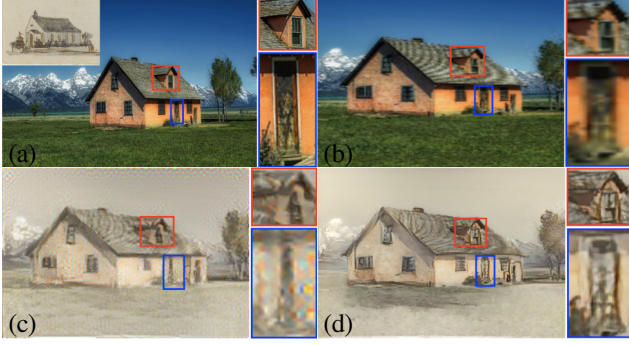
Figure 10: Style transfer with multi-scale strategy. (a) Content and style (upper left) images. (b) Output after the geometric warping. (c) Output without multi-scale synthesis. (d) Output with multi-scale strategy. The colored boxes show the magnified details. (b) shows that the resolution will decrease after the geometric warping, and details in (c) will lose accordingly. Multi-scale strategy (Section 6) will ensure the generation of high-resolution results without losing too much detail.
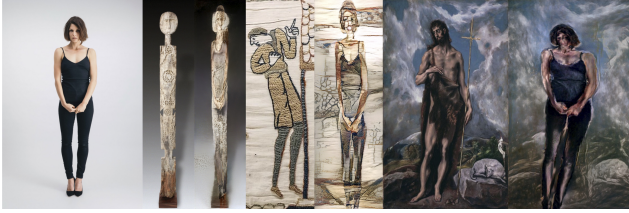


Figure 11: Some interesting cases, left to right, content image then 3 styles: African, Bayeux tapestry detail, el Grecco. In each case, the style image is on the left, output on the right.



Figure 12: Style transfer using three input images. Content image, top left; geometry style images top-middle and top-right; texture style images left-middle and left-bottom. Output images in the corresponding $2 \times 2$ array.

same semantic content and each show one major object; (ii) the geometric style can be matched using a continuous warping functions across the whole image. Assumption (i) is required for feature matching. Assumption (ii) confines the geometric styles we can reach beyond changing to higher-order maps. Some styles such as Cubism require piece-wise spatial mappings; other artists often change use of orthogonal projection, or use many vanishing points; pose can be un-natural, as in Egyptian art.

Limiting case (i) is less limiting than it sounds. First of all, the method itself is agnostic with respect to image content, but best results are to be had with similar semantic content. To see why this is the case recall the fact that human artists alter geometry for emphasis and in non-arbitrary ways – Stubbs exaggerated bull-like characteristics to depict bigger, stronger animals. Similarly, faces are altered to bring out some desired latent character, such as femininity or masculinity. This means that the manner of warp is class-conditional: human artists do not usually try to distort horses towards houses. It is, therefore, not at all unreasonable for the geometric style picture to contain an exemplar related to the object class in the content image.

Interesting examples arise even when these conditions are ad-

hered to, as Figure 11 shows. A global spatial transform means that local pose changes *etc*. are not well modeled, noticeable in the Bayeux tapestry detail. El Grecco, known for stylized elongation of bodies, also has a pose change but our output suffers less, probably because the change of pose has little impact on the overall profile. Similar remarks apply to the African sculpture.

Important detail is not always transferred well. Facial features are lost in all cases, the Bayeux tapestry detail is not a convincing tapestry; and the African sculpture appears weathered. Where the content image is blank, our algorithm copies texture more-or-less directly from the texture image. Some other algorithms also do this, see Figure 5 for examples.

Finally, we do not have to use a single style image but can instead use two style images: one for geometric style the other for texture style, Figure 12 shows examples. The figure has one content image, in the top-left. Geometric style images are placed top-middle and top-right, with texture style images on the left column. These reference images form a $2 \times 2$ array of corresponding output images.

Using three images extends the current paradigm in a novel and useful way. This adds versatility to the system because different pictures can be used to specify different components of style. For example, as in Figure 12, pure texture can be used to specify the texture style, and statues can be used to specify geometry style. This would not be possible using a single image to specify style. The principle might be extended in the future so that different elements of output style (*e.g.* texture, geometry,

composition) are characterized by different example images.

## 8. Conclusions

Our paper provides a novel method for image stylization: geometric style transfer. We provide a network to compute geometric style in a general setting, and use multi-scale texture transfer to maintain image quality throughout the transfer. Experimental results illustrate the qualitative expressiveness of our stylized results and greater quantitative similarity to target styles than other algorithms' outputs. These results are consistent with the Art History literature, where projection style has been used to characterise human art [52].

Our algorithm does have limits that provides plenty of future work. The content of the style image must be similar enough to the content of the content image for high-level features to be matched. This is more general than the requirement of strong models of faces or text [59, 35, 4, 49, 55, 54], but it is nonetheless a restriction. Our algorithm is global, whereas many styles will be local. This is known to be the case for texture [33] and facial features *et al.* [55]. Some styles, such as Cubism, are beyond out scope – but this is true of every other NST algorithm we know of. If NST is to progress to styles of that kind, then algorithms that include geometric style transfer are inevitable.

## References

[1] Edward H Adelson, Charles H Anderson, James R Bergen, Peter J Burt, and Joan M Ogden. Pyramid methods in image processing. *RCA engineer*, 29(6):33–41, 1984.

[2] Rudolf Arnheim. *Art and visual perception: A psychology of the creative eye*. Univ of California Press, 1954.

[3] Fred L. Bookstein. Principal warps: Thin-plate splines and the decomposition of deformations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(6):567–585, 1989.

[4] Kaidi Cao, Jing Liao, and Lu Yuan. Carigans: Unpaired photo-to-caricature translation. *ACM Transactions on Graphics (Proc. of Siggraph Asia 2018)*, 2018.

[5] Alex J. Champandard. Semantic style transfer and turning two-bit doodles into fine artworks. arXiv:1603.01768 [cs.CV], 2016.

[6] Dongdong Chen, Jing Liao, Lu Yuan, Nenghai Yu, and Gang Hua. Coherent online video style transfer. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1105–1114, 2017.

[7] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. StyleBank: an explicit representation for neural image style transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1897–1906, 2017.

[8] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stereoscopic neural style transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6654–6663, 2018.

[9] Yi-Lei Chen and Chiou-Ting Hsu. Towards deep style transfer: A content-aware perspective. In *British Machine Vision Conference (BMVC)*, 2016.

[10] John P Collomosse and Peter M Hall. Cubist style rendering from photographs. *IEEE Transactions on Visualization and Computer Graphics*, 9(4):443–453, 2003.

[11] Doug DeCarlo and Anthony Santella. Stylization and abstraction of photographs. *ACM Transactions on Graphics (ToG)*, 21(3):769–776, 2002.

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.

[13] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. arXiv:1508.06576[cs.CV], 2015.

[14] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2414–2423, 2016.

[15] Leon A Gatys, Alexander S Ecker, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. Controlling perceptual factors in neural style transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3985–3993, 2017.

[16] A. Gupta, Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Characterizing and improving stability in neural style transfer. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4067–4076, 2017.

[17] Paul Haeberli. Paint by numbers: Abstract image representations. In *Proceedings of the 17th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, pages 207–214. ACM, 1990.

[18] P.M. Hall. Painting by example. In *Eurographics UK*, pages 159–167, 1998.

[19] Peter M Hall, John P Collomosse, Yi-Zhe Song, Peiyi Shen, and Chuan Li. Rtcams: A new perspective on nonphotorealistic rendering from photographs. *IEEE Transactions on Visualization and Computer Graphics*, 13(5):966–979, 2007.

[20] Charles Han, Eric Risser, Ravi Ramamoorthi, and Eitan Grinspun. Multiscale texture synthesis. *ACM Transactions on Graphics (ToG)*, 27(3):51:1–51:8, 2008.

[21] David J Heeger and James R Bergen. Pyramid-based texture analysis/synthesis. In *Proceedings of the 22nd Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, pages 229–238. ACM, 1995.

[22] Aaron Hertzmann. A survey of stroke-based rendering. *IEEE Computer Graphics and Applications*, 23(4):70–81, 2003.

[23] Aaron Hertzmann, Charles E Jacobs, Nuria Oliver, Brian Curless, and David H Salesin. Image analogies. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, pages 327–340. ACM, 2001.

[24] Haozhi Huang, Hao Wang, Wenhan Luo, Lin Ma, Wenhao Jiang, Xiaolong Zhu, Zhifeng Li, and Wei Liu. Real-time neural style transfer for videos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 783–791, 2017.

[25] Hua Huang, Lei Zhang, and Hong-Chao Zhang. Arcimboldo-like collage using internet images. *ACM Transactions on Graphics (ToG)*, 30(6):155, 2011.

[26] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1501–1510, 2017.

[27] Yongcheng Jing, Yang Liu, Yezhou Yang, Zunlei Feng, Yizhou Yu, Dacheng Tao, and Mingli Song. Stroke controllable fast style transfer with adaptive receptive fields. In *European Conference on Computer Vision (ECCV)*, pages 238–254, 2018.

[28] Y. Jing, Yezhou Yang, Zunlei Feng, Jingwen Ye, Yizhou Yu, and Mingli Song. Neural style transfer: A review. *IEEE Transactions on Visualization and Computer Graphics*, 2019.

[29] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision (ECCV)*, pages 694–711, 2016.

[30] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[31] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. `arXiv:1412.6980[cs.LG]`, 2014.

[32] Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. Style transfer by relaxed optimal transport and self-similarity. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10051–10060, 2019.

[33] Dmytro Kotovenko, Artsiom Sanakoyeu, Sabine Lang, and Björn Ommer. Content and style disentanglement for artistic style transfer. In *IEEE International Conference on Computer Vision (ICCV)*, 2019.

[34] Dmytro Kotovenko, Artsiom Sanakoyeu, Pingchuan Ma, Sabine Lang, and Björn Ommer. A content transformation block for image style transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10032–10041, 2019.

[35] Wenbin Li, Wei Xiong, Haofu Liao, Jing Huo, Yang Gao, and Jiebo Luo. Carigan: Caricature generation through weakly paired adversarial learning. `arXiv:1811.00445[cs.CV]`, 2018.

[36] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast arbitrary style transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[37] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Diversified texture synthesis with feed-forward networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3920–3928, 2017.

[38] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems (NIPS)*, pages 386–396, 2017.

[39] Y. Li, Ming-Yu Liu, Xueting Li, Ming-Hsuan Yang, and Jan Kautz. A closed-form solution to photorealistic image stylization. In *European Conference on Computer Vision (ECCV)*, 2018.

[40] Jing Liao, Yuan Yao, Lu Yuan, Gang Hua, and Sing Bing Kang. Visual attribute transfer through deep image analogy. *ACM Transactions on Graphics (ToG)*, 36(4):120:1–120:15, 2017.

[41] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755, 2014.

[42] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep photo style transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6997–7005, 2017.

[43] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 4898–4906, 2016.

[44] Roey Mechrez, Eli Shechtman, and Lihi Zelnik-Manor. Photorealistic style transfer with screened poisson equation. In *British Machine Vision Conference (BMVC)*, 2017.

[45] Roey Mechrez, Itamar Talmi, and Lihi Zelnik-Manor. The contextual loss for image transformation with non-aligned data. In *European Conference on Computer Vision (ECCV)*, pages 768–783, 2018.

[46] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.

[47] Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. Artistic style transfer for videos and spherical images. *International Journal of Computer Vision*, 126:1199–1219, 2018.

[48] Ahmed Selim, Mohamed Elgharib, and Linda Doyle. Painting style transfer for head portraits using convolutional neural networks. *ACM Transactions on Graphics (ToG)*, 35(4):129, 2016.

[49] Yichun Shi, Debayan Deb, and Anil K Jain. Warpgan: Automatic caricature generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10762–10771, 2019.

[50] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *The International Conference on Learning Representations (ICLR)*, 2015.

[51] Xavier Snelgrove. High-resolution multi-scale neural texture synthesis. In *SIGGRAPH ASIA 2017 Technical Briefs*. ACM, 2017.

[52] John Willats. *Art and representation: New principles in the analysis of pictures*. Princeton University Press, 1997.

[53] Zheng Xu, Michael Wilber, Chen Fang, Aaron Hertzmann, and Hailin Jin. Learning from multi-domain artistic images for arbitrary style transfer. `arXiv:1805.09987[cs.CV]`, 2018.

[54] Shuai Yang, Zhangyang Wang, Zhaowen Wang, Ning Xu, Jiaying Liu, and Zongming Guo. Controllable artistic text style transfer via shape-matching gan. In *IEEE International Conference on Computer Vision (ICCV)*, 2019.

[55] Jordan Yaniv, Yael Newman, and Ariel Shamir. The face of art: Landmark detection and geometric style in portraits. *ACM Transactions on Graphics (ToG)*, 38(4):60, 2019.

[56] Yuan Yao, Jianqiang Ren, Xuansong Xie, Weidong Liu, Yong-Jin Liu, and Jun Wang. Attention-aware multi-stroke style transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[57] Jaejun Yoo, Youngjung Uh, Sanghyuk Chun, Byeongkyu Kang, and Jung-Woo Ha. Photorealistic style transfer via wavelet transforms. In *IEEE International Conference on Computer Vision (ICCV)*, 2019.

[58] Jingyi Yu and Leonard McMillan. A framework for multiperspective rendering. In *Rendering Techniques*, pages 61–68, 2004.

[59] Ziqiang Zheng, Wang Chao, Zhibin Yu, Nan Wang, Haiyong Zheng, and Bing Zheng. Unpaired photo-to-caricature translation on faces in the wild. *Neurocomputing*, 2019.