Figure 1: Different style transfer methods. Left to right: content [17], style [3], Gatys *et al.* [7], Li and Wand [14], ours.

# Improving the Neural Algorithm of Artistic Style

**Roman Novak** and **Yaroslav Nikulin**
Department of Mathematics
École normale supérieure de Cachan
94230 Cachan, France
{rnovak, yaroslav.nikulin}@ens-cachan.fr

## Abstract

In this work we investigate different avenues of improving the Neural Algorithm of Artistic Style [7].

While showing great results when transferring homogeneous and repetitive patterns, the original style representation often fails to capture more complex properties, like having separate styles of foreground and background. This leads to visual artifacts and undesirable textures appearing in unexpected regions when performing style transfer.

We tackle this issue with a variety of approaches, mostly by modifying the style representation in order for it to capture more information and impose a tighter constraint on the style transfer result.

In our experiments, we subjectively evaluate our best method as producing from barely noticeable to significant improvements in the quality of style transfer.

## 1 Introduction

We first briefly describe the original algorithm of style transfer presented in [7].

In the following we denote $F^l$ the 3D output (activation / feature volume) of a convolutional layer $l$ of size $K_l \times X_l \times Y_l$ with $K_l$ being the number of convolutional filters (features, kernels) in layer $l$, $X_l$ and $Y_l$ being the spatial dimensions.

Using a pre-trained object recognition deep convolutional neural network the content of an image is defined as an activation volume $F^{l_c}$ of a fixed (usually relatively deep) layer $l_c$ of the network given the image as an input. In [7] the VGG [21] network was used and the output of conv4_2 layer was considered as content.

Style is defined on top of activation volumes as a set of Gram correlation matrices $\{\mathcal{G}_l\}$ with the entries $\mathcal{G}_{ij}^l = \langle F_i^l, F_j^l \rangle$, where $l$ is the network layer, $i$ and $j$ are two filters of the layer $l$ and $F_k^l$ is

the activation map (indexed by spatial coordinates) of filter $k$ in layer $l$. In other words, $\mathcal{G}_{ij}^l$ value says how often features $i$ and $j$ in the layer $l$ appear together in the input image.

Once defined, the content and style spaces can be used to find a synthesized image $I$ matching both given content $C$ and style $S$ images, resulting in high-quality repainting of the content image with preserved semantic information but covered with textures extracted from the style image. The resulting image $I$ is synthesized by back-propagating the loss

$$\alpha \texttt{StyleLoss}(S,I) + \texttt{ContentLoss}(C,I) = \alpha \sum_l w_l \left\| \mathcal{G}^l(S) - \mathcal{G}^l(I) \right\|^2 + \left\| F^{l_c}(C) - F^{l_c}(I) \right\|^2$$

into pixels of $I$ ($\alpha$ is the style/content trade-off hyper-parameter and $w_l$ is a style weight assigned to a specific layer, set to 1 for all layers by default).

The suggested synthesizing algorithm generally produces great results in transferring repetitive artistic styles. Unfortunately, otherwise generated images often don't meet human expectations of how the style should be transferred. This is the problem we attempt to tackle in this work.

Further, when visually inspecting style transfer results and commenting on the quality, we use two complimentary criteria that we expect a good style transfer algorithm to meet:

1. Similar areas of the content image should be repainted in a similar way.
2. Different areas should be repainted differently.

As it turns out, it is often difficult to satisfy both simultaneously.

The rest of the work is structured as follows: section 2 gives an overview of other research built on top of the original style transfer algorithm. Section 3 presents our suggested improvements to the style transfer algorithm. Section 5 describes and presents the results of experiments conducted with a visual comparison of different approaches to style transfer. We give concluding remarks in section 6.

## 2 Related Work

Several works related to [7] have been published recently, both introducing new applications of style transfer and improving the algorithm itself.

In [24] the impact of resolution and content-style alignment are studied, however, a significant part of suggested pipeline was done manually in photo-editing software and seems to be applicable only to a pair of strongly similar content-style images. Another novelty presented in this paper is super-resolution style transfer when a higher resolution style image is iteratively transferred to a lower resolution content image, upsampled at every iteration.

To improve the style transfer and photorealism of the generated images, a combination of dCNNs and Markov Random Fields is used in [14] to leverage matching feature patches in content and style images. Instead of global bag-of-features-like approach as in [7], [14] works with best corresponding patches of feature maps: MRF prior gives better local plausibility of the image. Visually, results produced by [14] resemble repainting the content image with appropriate pieces from the style image, while the original algorithm [7] can be informally described as generating textures of different scale and complexity over the content image. It should be noted that, while giving more photorealistic results (with a photorealistic style image as input), a good semantic correspondence between two images is needed.

An interesting application of [14] can be found in [4], where a rough sketch can be transformed into a painting using high-level annotations of an example image. The sketch can have a very different layout of image parts like sky, water, forest, etc. with respect to the original image, and still be repainted into a very credible result.

Used in unguided setting (without content image), the algorithm can produce high-quality textures based on a given example [8]. Moreover, using the style space from [7] one can define an appropriate loss function and train a network that learns a given style [23]. It alleviates the computational cost

of generating a new texture (delegating it to the training stage). Basically, in [23] the network learns to process input noise image in a way to minimize the style loss function introduced by [7] for a chosen style image.

In the following, we compare our approach with the original one [7] and with the CNN-MRF method from [14]. The latter pursues similar goals to ours, improving the quality of style transfer via patch matching, while we suggest to enforce a tighter statistical constraint in an augmented style space.

## 3 Helpful Modifications

Below we describe all modifications to the original style transfer algorithm in [7] that turned out to be helpful.

Our core contributions are:

1. A better per-layer content/style weighting scheme (section 3.1);
2. Using more layers to capture more style properties (section 3.2);
3. Using shifted activations when computing Gram matrices to eliminate sparsity and make individual entries more informative (section 3.3) and also speed-up style transfer convergence;
4. Targeting correlations of features belonging to different layers to capture more feature interactions (sections 3.4, 3.5 and 3.6).

Modifications leading to questionable results are presented in section 4.

The impact of our changes (both helpful and not) is presented in section 5.

### 3.1 Layer Weight Adjustment

In [7] content layer `conv4_2` and 5 style layers $\{$`conv1_1, conv2_1, conv3_1, conv4_1, conv5_1`$\}$ (of VGG-19) with unit weight $w_l$ each are used. We try to soften the "content – style" separation by using the same set of layers for both content and style, but with a geometric weighting scheme: for each layer $l$ we set the style $w_l^s$ and content $w_l^c$ weights to

$$w_l^s = 2^{D-d(l)}, \quad w_l^c = 2^{d(l)},$$

where $D$ is the total number of layers used and $d(l)$ is the depth of layer $l$ with respect to all the other layers used. Thus we try to preserve more information about style and content (as neither is limited to being contained in a single layer) and also indicate that the most important style properties (e.g. colors, simple patterns) are captured by the bottom layers while content is mostly represented by activations in the upper layers.

For the setup described in section 5, we have found this approach to sometimes produce a slight improvement in quality (in particular, the object tends to be better separated from the background), although oftentimes the results seem barely different.

### 3.2 Using More Layers

In [7] only five layers are used for the calculation of Gram matrices. We try to enrich the style representation by calculating Gram matrices for all 16 convolutional layers of VGG-19. We observe that this approach leads to a consistent improvement in quality across the majority of style images considered.

### 3.3 Activation Shift

Convolutional layers of the normalized version [12] of VGG-19 output non-negative activations with mean 1. For a typical image outputs are also sparse: in all layers, each filter has usually few non-zero activations across spatial dimensions. This results in Gram matrices being sparse as well.

3

We argue that sparsity is detrimental to style transfer quality. Precisely, it allows for unexpected patterns to appear in regions one would typically expect to be filled with, for example, a uniform background color.

Our informal reasoning is as follows: zero entries of a Gram matrix are easy to misinterpret. Indeed, $\mathcal{G}_{ij}^l = 0$ could mean that features $i$ and $j$ of layer $l$ are either both absent, have one of them absent or both present but never appear together. And since, empirically, Gram matrices contain a lot of zero entries, they leave too much freedom for the optimization procedure to "interpret" them in a wrong way.

We therefore decide to eliminate sparsity by calculating Gram matrices using shifted activations: instead of setting

$$\mathcal{G}^l = F^l F^{l^T},$$

(here $F^l$ is considered linearized into a 2D matrix of size $K_l \times X_l Y_l$) we put

$$\mathcal{G}^l = \left(F^l + s\right)\left(F^l + s\right)^T,$$

where $s$ is the shift value added to matrices element-wise. In this case, the gradient contributions in [7] should be changed respectively:

$$\frac{\partial \mathcal{G}^l}{\partial F^l} = 2(F^l + s).$$

Through several experiments we have concluded that putting $s = -1$ (i.e. centering activations at $0$) yields the best result. We have observed that shifting activations by this amount consistently improves the quality of style transfer across most of style images and style transfer methods. It also helps the optimization procedure to escape the starting point, where it can sometimes get stuck for hundreds of iterations otherwise.

Finally, activation shift appears to promote uniform image repainting. When descending from content, the original style transfer often applies style with varying intensity, i.e. strongly repaints certain regions while barely modifying others (see, for example, second column of table 4). It may require fine-tuning the number of iterations and style weight to properly color all the regions without distorting content excessively. We encounter such problems less frequently when using an activation shift.

### 3.4 Correlations of Features from Different Layers

In [7] style information is captured by a set of Gram matrices $\left\{\mathcal{G}^l\right\}$ of feature correlations within the same layer $l$. In our implementation we extend this definition to a set of Gram matrices $\left\{\mathcal{G}^{lk}\right\}$ of feature correlations belonging to possibly different layers $l$ and $k$. Since feature maps $F^l$ and $F^k$ may have different spatial dimensions, we always upsample the smaller map to match the size of the bigger one and thus define

$$\mathcal{G}^{lk} = F^l \left[\text{up}(F^k)\right]^T,$$

if $X_k \leqslant X_l$. $\mathcal{G}^{lk}$ has the size of $K_l \times K_k$ and, as earlier, an entry $\mathcal{G}_{ij}^{lk}$ represents the rate of co-occurrence of features $i$ and $j$ but now belonging to possibly different layers $l$ and $k$ respectively.

Having this definition the derivatives with respect to activation volumes change if $l \neq k$ :

$$\frac{\partial \mathcal{G}^{lk}}{\partial F^l} = \text{up}(F^k), \quad \frac{\partial \mathcal{G}^{lk}}{\partial F^k} = \left[\text{up}^{-1}(F^l)\right]^T.$$

With 16 convolutional layers in VGG-19 we have $2^{16^2}$ ways to "tie" them into a definition of style (i.e. a set of matrices $\mathcal{G}^{lk}$). One of the approaches we tried first was to tie all layers with a single

4

content layer, motivated by our perception of style as of the way in which low-level features correlate with high-level features. Precisely, instead of using $\left\{ \mathcal{G}^l \,|\, l = 1 \ldots 16 \right\}$ (with layers `conv1_1` to `conv5_4` numbered from 1 to 16) as our style representation, we use $\left\{ \mathcal{G}^{cl} \,|\, l = 1 \ldots 16 \right\}$ with $c$ corresponding to a single high-level layer (`conv4_2` for example). However, as seen in section 5, this particular style representation didn't yield any improvement. In general, we observe that tying distant in terms of depth layers produces poor results.

## 3.5   Correlation Chain

In the framework of using inter-layer correlations for style as described in section 3.4, we now consider the following "chained" style representation: $\left\{ \mathcal{G}^{l,l-1} \,|\, l = 2 \ldots 16 \right\}$. We thus still constrain correlations of high- and low-level features, but in a local way, where only the correlation with immediate neighbors are considered.

This approach has lead to a consistent and often significant improvement in style transfer quality in most cases considered.

## 3.6   Blurred Correlations

Recall that in section 3.4 the smaller activation map is upsampled to match the bigger one and then pixel-wise correlations of their features are calculated. However, even having the same spatial dimensions, feature maps still correspond to features of different scales. For example, consider layers `conv1_1` and `conv1_2`: they have the same size, but features of `conv1_2` are composed of $3 \times 3$ patches of `conv1_1`-features (as VGG-19 uses $3 \times 3$ kernels). Therefore, it may be reasonable to not capture the correlation of `conv1_2`-activations the activations in pixels directly "underneath" them (which is the case in the default implementation of $\mathcal{G}^{lk}$), but rather measure its correlation with an average activation of the other feature in the respective $3 \times 3$-patch.

Precisely, we implement this as follows:

$$\mathcal{G}^{lk} = F^l \left[ \text{blur}^{l-k} \circ \text{up}(F^k) \right]^T,$$

i.e. we apply box blurring $l - k$ times to the upsampled version of the smaller activation. Then the new derivatives become

$$\frac{\partial \mathcal{G}^{lk}}{\partial F^l} = \text{blur}^{l-k} \circ \text{up}(F^k), \quad \frac{\partial \mathcal{G}^{lk}}{\partial F^k} = \left[ \text{up}^{-1} \circ \text{blur}^{l-k} \left( F^l \right) \right]^T.$$

This modification has produced mixed, but overall positive results. However, it does complicate the objective function which results in slow and unreliable convergence.

# 4   Other Approaches

Below we describe some other attempts to improve style transfer that didn't work out in the end.

## 4.1   Gradient Masking

One of our first attempts to improve style transfer quality was to prevent unexpected textures from appearing in the background of a content image.

We try a very simple idea: apply a binary mask to the style gradient from every layer. The mask was supposed to suppress undesirable textures in areas with weak content response, corresponding to monotone areas. Precisely, we replace the derivative:

$$\frac{\partial \mathcal{G}^l}{\partial F^l} = 2F^l$$

with

$$\frac{\partial \mathcal{G}^l}{\partial F^l} = 2F^l \odot S^l,$$

with mask $S$ being generated by simple thresholding of the feature maps of the content image $P^l$:

$$S_i^l(x, y) = \begin{cases} 1 & \text{if } P_i^l(x, y) > \texttt{threshhold}(l) \\ 0 & \text{otherwise} \end{cases}.$$

The $\texttt{threshhold}(l)$ value was set in a way to keep the desired portion of the output. From a few experiments we have found that keeping $100\%, 40\%, 20\%, 10\%, 10\%$ of biggest entries for the layers $\texttt{conv1\_1}, \texttt{conv2\_1}, \texttt{conv3\_1}, \texttt{conv4\_1}, \texttt{conv5\_1}$ respectively indeed suppressed texture generation on the background while preserving reasonable quality of style transfer to the object in the foreground. Notice that we let the mask pass more shallow layer activations and fewer deep layer activations, expecting the color and simple style features to replace the background in the content image. Examples of style transfer results using masked gradients are presented in figure 2.

However, applying a mask modifies the gradient and thus damages the minimization procedure. In practice we observed that after several hundreds of iterations the error began to stall or even increase. Also, the portion of elements one should set to zero was found experimentally for a given content image, and we do not expect this approach to generalize well.

## 4.2 Amplifying Activations

In an attempt to make the approach in section 4.1 generalize better, we decide to not have a hard threshold on activations, but rather amplify them so that regions with strong feature responses in the content image are prioritized. Precisely, we set

$$\mathcal{G}^l = F^{lp} \left( F^{lp} \right)^T,$$

so that the derivative becomes

$$\frac{\partial \mathcal{G}^l}{\partial F^l} = 2p F^{l^{2p-1}},$$

where the exponentiation is performed element-wise. Then, setting $p$ to 2 would result in $\frac{\partial \mathcal{G}^l}{\partial F^l} = 2p F^{l^3}$, amplifying and prioritizing change in strong activation regions.

With this objective we indeed observed how during style transfer feature-rich regions (e.g. eyes) were modified before and more strongly than the rest of the image. However, we could not fully evaluate this idea (and, notably, see if this modified objective makes sense) as it had serious problems with convergence. We think the reason is the objective function becoming highly non-linear, since convergence problems were less severe as we reduced $p$ from 2 down to 1.

## 4.3 Adjacent Activations Correlations

In [7], as well as in previous sections, pixel-wise feature correlations are considered (apart from a slight modification in section 3.6, where high-level features are correlated with average values of low-level feature patches). We now try to also constrain how features correlate with other features not only in the exact same spatial location, but in the 8 adjacent locations as well. Precisely, a $K_l \times K_k$ Gram matrix $\mathcal{G}^{lk}$ expands into a $3 \times 3 \times K_l \times K_k$ matrix with new dimensions corresponding to shifts along $X$ and $Y$ axis $(-1, 0, 1)$. If we introduce the following matrix of pixel shifts

$$S = \begin{bmatrix} (-1, -1) & (-1, 0) & (-1, 1) \\ (0, -1) & (0, 0) & (0, 1) \\ (1, -1) & (1, 0) & (1, 1) \end{bmatrix},$$

6

Figure 2: Comparison of the original algorithm in [7] (left) and the algorithm with masked style gradient (right). Content image (here used in $400 \times 300$ resolution): [17]. Style images (top to bottom): [20, 6, 1, 11, 9, 3, 2].

then we define the new Gram matrix as

$$\tilde{\mathcal{G}}^{lk} = \text{map}\left(\mathcal{G}^{lk}, S\right) = \begin{bmatrix} \mathcal{G}^{lk}(-1,-1) & \mathcal{G}^{lk}(-1,0) & \mathcal{G}^{lk}(-1,1) \\ \mathcal{G}^{lk}(0,-1) & \mathcal{G}^{lk}(0,0) & \mathcal{G}^{lk}(0,1) \\ \mathcal{G}^{lk}(1,-1) & \mathcal{G}^{lk}(1,0) & \mathcal{G}^{lk}(1,1) \end{bmatrix},$$

where

$$\mathcal{G}^{lk}(dx, dy) = F^l \text{shift}(F^k, dx, dy)^T,$$

with the "shift" operation shifting (with 0 padding) the activation map $dx$ and $dy$ pixels along $X$ and $Y$ dimensions.

With this new construction the derivatives with respect to activation maps become

$$\frac{\partial \tilde{\mathcal{G}}^{lk}}{\partial F^l} = \text{map}\left(\text{shift}(F^{k^T}, \cdot, \cdot), S\right), \quad \frac{\partial \tilde{\mathcal{G}}^{lk}}{\partial F^k} = \text{map}\left(\text{shift}(F^l, \cdot, \cdot), -S\right).$$

This modification has a significant impact on the style transfer result, but in a negative way more often than not.

## 4.4 Content-aware Gram Matrices

We now describe an attempt to make the style more "content-aware". Precisely, during style transfer, we want different styles to be applied to different regions of an image, depending on the content of the region. In a simple case we want styles of background and foreground to be translated separately (contrary to being blended into a single statistic over the whole image). In a more complex case we want the style to vary as a function of content and result in an overall improvement of style transfer quality.

We first pick a content layer, for example $l_c = \texttt{conv4\_2}$. It has $K_{l_c} = 512$ activation maps $F_k^{l_c}$, i.e. 512 "types of content" that we distinguish. For each type we want to have separate styles, i.e. separate sets of Gram matrices $\mathcal{G}^l$ (for whichever $l$s we decide to represent style).

We store this information by introducing depth to Gram matrices: a third dimension of size $K_{l_c}$. We thus introduce "content-aware Gram matrices" $\mathcal{G}_{l_c}^l$ of size $K_{l_c} \times K_l \times K_l$ as

$$\left(\mathcal{G}_{l_c}^l\right)_{k,\cdot,\cdot} = \texttt{styleBy}\left(F^l, F_k^{l_c}\right),$$

where $\texttt{styleBy}$ is a function used to compute a style statistic of $F_l$ that best describe content of type $k$ based on the supplied activation map $F_k^{l_c}$. Having non-negative activations (see section 3.3) we suggest defining $\texttt{styleBy}$ as

$$\left[\texttt{styleBy}\left(F^l, F_k^{l_c}\right)\right]_{ij} = \sum_{x,y} F_i^l(x,y) \cdot F_j^l(x,y) \cdot F_k^{l_c}(x,y),$$

i.e. the $K_{l_c}$ slices of $\mathcal{G}_{l_c}^l$ are Gram matrices with correlation contributions weighted by response of feature $k$ in layer $l_c$ (for simplicity it is assumed here that $F^l$ and $F^{l_c}$ have the same spatial dimensions, which is generally not true. In this case the smaller map is upsampled as described in section 3.4).

Thus each feature collects style statistics present in regions where the feature has a high response.

Formally

$$\left(\mathcal{G}_{l_c}^l\right)_{k,i,j} = \left(F_i^l \odot F_k^{l_c}\right) F_j^{l^T},$$

and the derivatives are

8

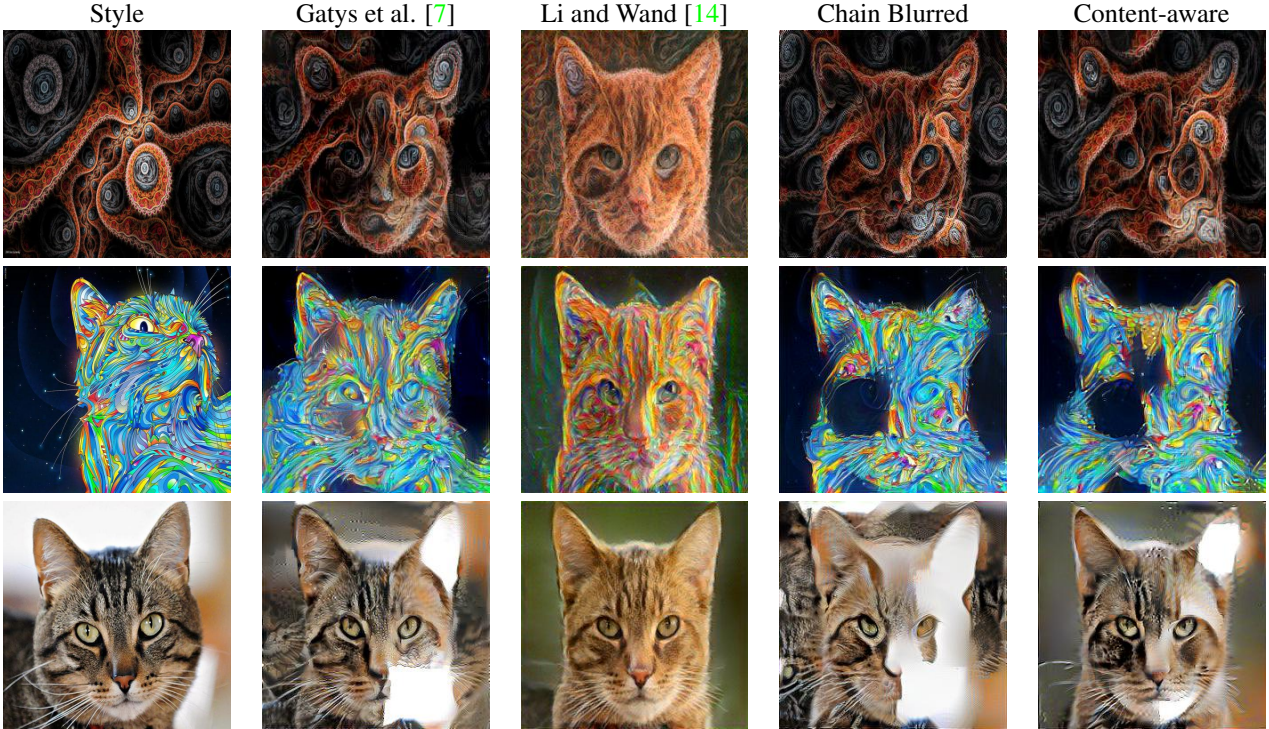| Style | Gatys et al. [7] | Li and Wand [14] | Chain Blurred | Content-aware |
|-------|------------------|------------------|---------------|---------------|

Table 1: Content-aware Gram matrices mixed performance compared to other approaches on $256 \times 256$ images. Content: figure 3. Style images (top to bottom): [11, 3].

$$\frac{\partial \left(\mathcal{G}^l_{l_c}\right)_{k,i,j}}{\partial F^l_i} = F^l_j \odot F^{l_c}_k, \quad \frac{\partial \left(\mathcal{G}^l_{l_c}\right)_{k,i,j}}{\partial F^l_j} = F^l_i \odot F^{l_c}_k, \quad \frac{\partial \left(\mathcal{G}^l_{l_c}\right)_{k,\cdot,\cdot}}{\partial F^{l_c}_k} = F^l_i \odot F^l_j.$$

On a practical note, we have also observed that shifting style layer activations (see section 3.3) was helpful in this context as well.

We were only able to test this approach on $256 \times 256$ images with `conv4_2` as the content layer and `conv1_1`, `conv2_1`, `conv3_1`, `conv4_1` (almost as in [7]) as style layers. While in the photorealistic context it has lead to an improvement (see table 1), it did not generalize well to other styles and generally performed comparably to other statistic-based approaches (while being much slower).

### 4.5 Gram Cubes

In the spirit of section 4.4 we can also try to define style as a set of "Gram cubes" $\tilde{\mathcal{G}}^l$ of size $K_l \times K_l \times K_l$ storing triple correlations of features in each layer:

$$\left(\tilde{\mathcal{G}}^l\right)_{k,i,j} = \left(F^l_i \odot F^l_j\right) F^{l^T}_k.$$

In practice, however, it did not lead to any promising results and was not deemed worth the computational complexity.

## 5 Experiments

Suggestions described in section 3 were implemented in Torch [5] on top of Kai Sheng Tai's implementation [22] of the original style transfer [7]. The network used was VGG-19 [21], precisely, the normalized version with weights scaled to have a unit mean neuron activations over images and positions [12].

| Method name | Classic | Classic Shifted | Classic Dense | All-To-Content | Chain | Chain Uniform | Chain Unshifted | Chain Blurred | Chain Extended | CNN-MRF |
|---|---|---|---|---|---|---|---|---|---|---|
| Content layers | conv4_2 | | all convolutional | | | | | | | Default |
| Style layers | conv1_1, conv2_1, conv3_1, conv4_1, conv5_1 | | all convolutional | conv4_2 - conv5_4, conv4_2 - conv5_3, ..., conv4_2 - conv1_2, conv4_2 - conv1_1 | conv5_4 - conv5_3, conv5_3 - conv5_2, ..., conv2_1 - conv1_2, conv1_2 - conv1_1 | | | | | |
| Weighting scheme | uniform | | geometric | | | uniform | geometric | | | |
| Activation shift | 0 | -1 | | | | | 0 | -1 | | |
| Blurred correlation | no | | | | | | | yes | no | |
| Adjacent activations correlation | no | | | | | | | | yes | |

Table 2: Style transfer methods compared.

All style and content images were rescaled to the size of $512 \times 512$ pixels. For each presented setup 270 iterations of L-BFGS optimization algorithm were performed starting from the content image.

By default the style weight was set to $2 \times 10^9$. However, in some cases the algorithm had trouble escaping the starting image for hundreds of iterations, in which case we either modified the style weight and re-run the computation or saved a result of a later iteration. The CNN-MRF [14] images were generated using default settings of the implementation shared by authors in [13].

The content image is displayed in figure 3.



Figure 3: Content image [17] considered for experiments.

Style transfer methods compared in this work are presented in table 2. In tables 3 and 4 we compare `Classic` [7], `CNN-MRF` [14] and `Chain Blurred` to demonstrate the performance of our approach compared to existing solutions. In table 5 we compare all configurations to demonstrate the impact of individual modifications described in section 3.
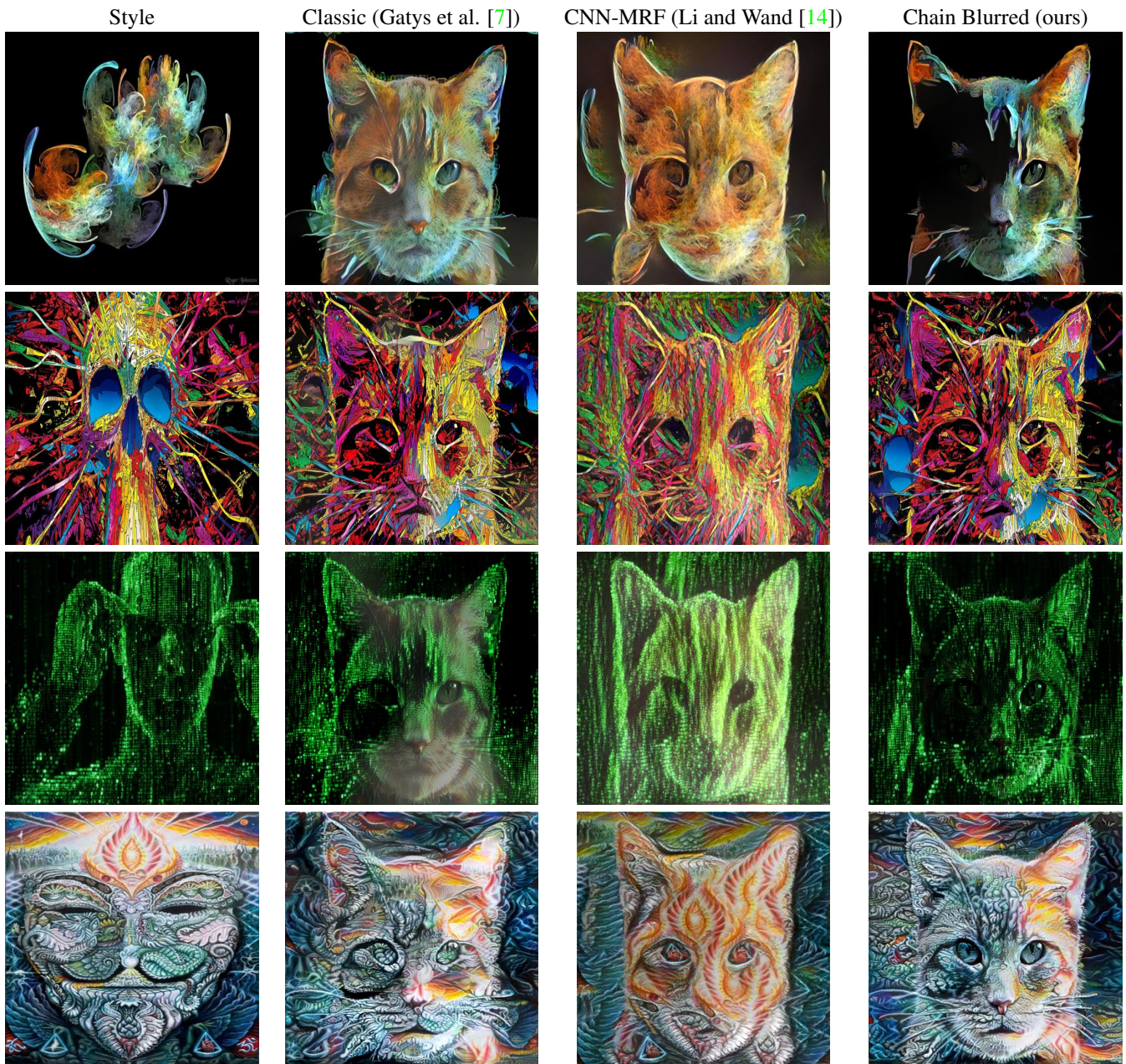
| Style | Classic (Gatys et al. [7]) | CNN-MRF (Li and Wand [14]) | Chain Blurred (ours) |



Table 3: Comparing our results to existing style transfer solutions (part 1). Style images (top to bottom): [16, 2, 15, 20].

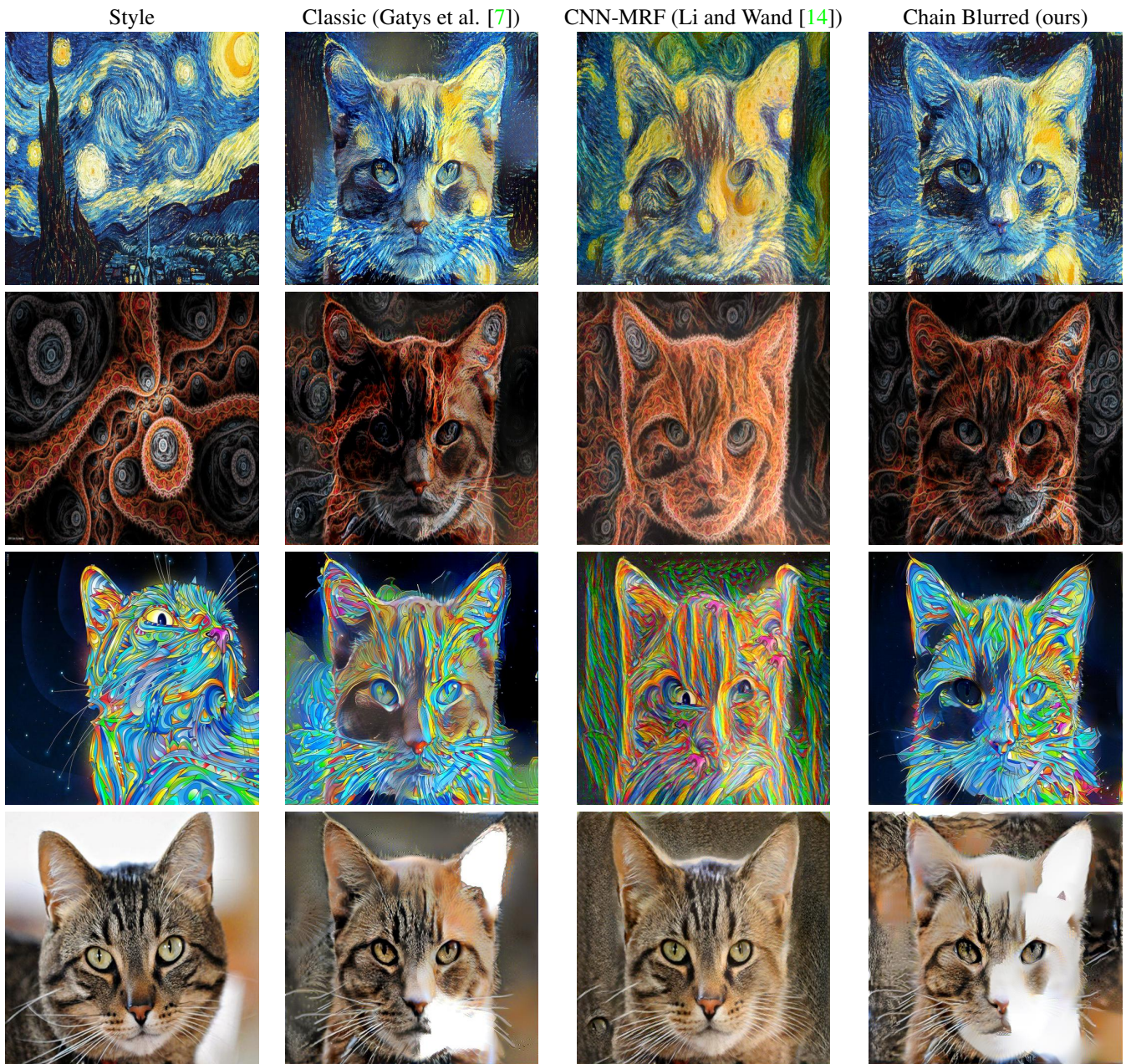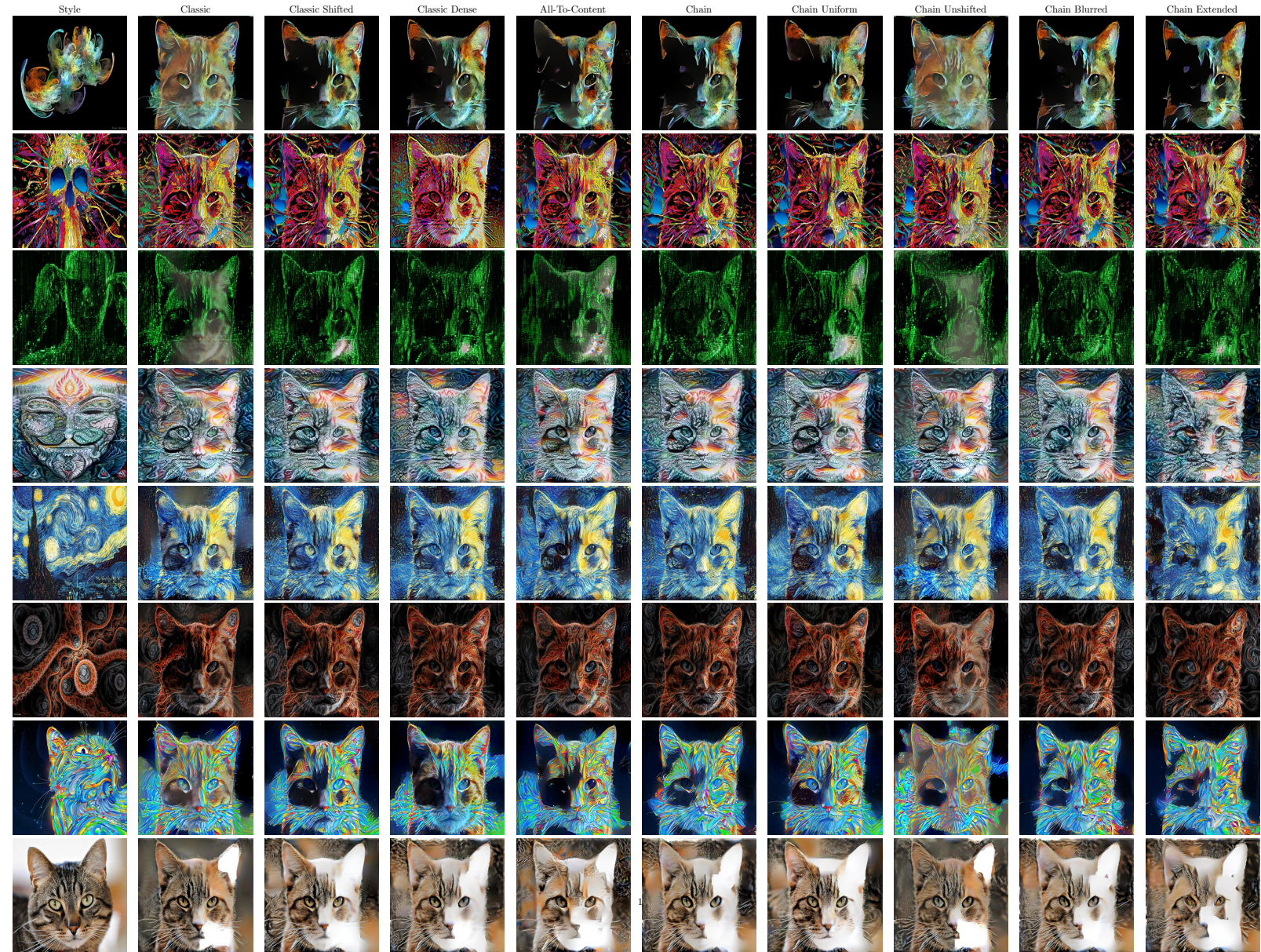| Style | Classic (Gatys et al. [7]) | CNN-MRF (Li and Wand [14]) | Chain Blurred (ours) |



Table 4: Comparing our results to existing style transfer solutions (part 2). Style images (top to bottom): [9, 11, 3, 18].

| Style | Classic | Classic Shifted | Classic Dense | All-To-Content | Chain | Chain Uniform | Chain Unshifted | Chain Blurred | Chain Extended |

Table 5: Demonstration of individual modifications contributions.

# 6 Conclusion

In this work several ways to improve the style transfer algorithm suggested in [7] were investigated. The direction of our experiments was mostly inspired by the Julesz Conjecture [10, 19], according to which two textures sharing a certain set of statistics should be visually indistinguishable. Similarly, our conjecture was that there exists a set of statistics describing the style of an image with desirable visual properties, and we used feature correlation matrices from [7] as a starting point.

Our most useful contributions turned out to be, in order of decreasing impact

- activation shift (see section 3.3) that eliminates the ambiguity of zero entries in Gram matrices and improves results while accelerating convergence across different style images and style transfer methods;
- augmenting the style representation by using 16 layers (contrary to 5 in [7]) and by considering a chain of inter-layer feature correlations (see sections 3.4, 3.5 and 3.6);
- geometric weighting scheme to soften the style/content separation and to prioritize simple style features when repainting (see section 3.1).

These changes have consistently yielded improvements on the majority of style images in our experiments, ranging from marginal up to very significant, although at the cost of increased memory and computational cost compared to [7].

Other suggestions discussed in section 4 gave mixed results and may require further research in order for them to bare fruit.

We believe our modifications to be most useful when using a style of "average complexity".

Simple, repetitive textures are very well transferred by the original algorithm in [7]. However, it does fail often on more complex styles, e.g. images with a background and a foreground.

On the other end of complexity range, in the context of photorealistic style transfer, [14] remains unchallenged. However, it does require a very good content match between content and style. Otherwise results may be inferior to our approach, with violations of criterion 2 caused by certain content parts disappearing due to lack of distinct corresponding patches in the style image.

However, as one can see in section 5, our approach is still far from our goal of satisfying the two criteria in section 1: uniform background sometimes gets fragmented into differently stylized regions (violating criterion 1), and the foreground object may sometimes blend with background (violating criterion 2). It is merely a step in this direction, that we hope will be useful for future research.

# References

[1] URL: http://hqwallbase.pw/105645-strange-birds/.

[2] Matei Apostolescu. *Farewell Mr. Eldritch*. URL: http://www.013a.com/html/skull_color.htm.

[3] Matei Apostolescu. *Midnight Cat*. URL: http://www.013a.com/html/cat.htm.

[4] Alex J. Champandard. *Semantic Style Transfer and Turning Two-Bit Doodles into Fine Artworks*. 2016. URL: http://arxiv.org/abs/1603.01768.

[5] R. Collobert, K. Kavukcuoglu, and C. Farabet. "Torch7: A Matlab-like Environment for Machine Learning". In: *BigLearn, NIPS Workshop*. 2011. URL: http://torch.ch.

[6] Thomas C. Fedro. *Cubist 9*. 1969. URL: http://www.ebsqart.com/Art-Galleries/Contemporary-Cubism/43/Cubist-9/204218/.

[7] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. *A Neural Algorithm of Artistic Style*. 2015. URL: http://arxiv.org/abs/1508.06576.

[8] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. *Texture Synthesis Using Convolutional Neural Networks*. 2015. URL: http://arxiv.org/abs/1505.07376.

[9] Vincent Van Gogh. *Starry Night*. 1889. URL: https://en.wikipedia.org/wiki/The_Starry_Night.

[10] Béla Julesz. *Visual pattern discrimination*. 1962. URL: http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=1057698.

[11] Dan Kuzmenka. *Hatosia*. 2008. URL: https://www.fractalus.com/dan/Galleries/Fractals/Fractal%20Gallery%208/slides/Hatosia.html.

[12] Bethge Lab. *Artistic Stylization using Convolutional Neural Networks*. 2015. URL: http://bethgelab.org/deepneuralart/.

[13] Chuan Li and Michael Wand. *CNNMRF*. 2016. URL: https://github.com/chuanli11/CNNMRF.

[14] Chuan Li and Michael Wand. *Combining Markov Random Fields and Convolutional Neural Networks for Image Synthesis*. 2016. URL: http://arxiv.org/abs/1601.04589.

[15] Pando. *Square VP Jared Fliesler joins Matrix Partners as a general partner*. 2013. URL: https://pando.com/2013/03/07/square-vp-jared-fliesler-joins-matrix-partners-as-a-general-partner/.

[16] Elusive Parallelograms. *And Everything Changes*. 2009. URL: http://www.amazon.com/Everything-Changes-Elusive-Parallelograms/dp/B0025WUJHK.

[17] PetFinder. *The Special Grooming Needs of a Senior Cat*. 2015. URL: https://www.petfinder.com/cats/cat-grooming/grooming-needs-senior-cat/.

[18] PetFinder. *What Is Fostering a Cat?* 2015. URL: https://www.petfinder.com/animal-shelters-and-rescues/fostering-cats/what-is-cat-fostering/.

[19] Javier Portilla and Eero P. Simoncelli. *A Parametric Texture Model Based on Joint Statistics of Complex Wavelet Coefficients*. 1999. URL: http://www.cns.nyu.edu/pub/lcv/portilla99-reprint.pdf.

[20] Randal Roberts. *Fawkes*. 2012. URL: http://www.allofthisisforyou.com/gallery-2/fawkes2012/.

[21] K. Simonyan and A. Zisserman. *Very Deep Convolutional Networks for Large-Scale Image Recognition*. 2014. URL: http://arxiv.org/pdf/1409.1556.

[22] Kai Sheng Tai. *style-transfer*. 2015. URL: https://github.com/kaishengtai/neuralart.

[23] Dmitry Ulyanov et al. *Texture Networks: Feed-forward Synthesis of Textures and Stylized Images*. 2016. URL: http://arxiv.org/abs/1603.03417.

[24] Rujie Yin. *Content Aware Neural Style Transfer*. 2016. URL: http://arxiv.org/abs/1601.04568.