

Content-Aware Neural Style Transfer

Rujie Yin

January 19, 2016

Motivation

In the emerging inter-disciplinary field of art and image processing, algorithms have been developed to assist the analysis of art work. In most applications, especially brush stroke analysis, high resolution digital images of paintings are required to capture subtle patterns and details in the high frequency range of the spectrum. Algorithms have been developed to learn styles of painters from their digitized paintings to help identify authenticity of controversial paintings. However, high quality testing datasets containing both original and forgery are limited to confidential image files provided by museums, which is not publicly available, and a small sets of original/copy paintings painted by the same artist, where copies were deferred to two weeks after the originals were finished. Up to date, no synthesized painting by computers from a real painting has been used as a negative test case, mainly due to the limitation of prevailing style transfer algorithms.

There are two main types of style transfer algorithms, either transferring the tone (color, contrast, saturation, etc.) of an image, preserving its patterns and details, or distorting the texture uniformly of an image to create “style”. In this paper, we are interested in a higher level of style transfer, particularly, transferring a source natural image (e.g. a photo) to a high resolution painting given a reference painting of similar object. The transferred natural image would have a similar presentation of the original object to that of the reference painting. In general, an object is painted in a different style of brush strokes than that of the background, hence the desired style transferring algorithm should be able to recognize the object in the source natural image and transfer brush stroke styles in the reference painting in a *content-aware* way such that the styles of the foreground and the background, and moreover different parts of the foreground in the transferred image, are consistent to that in the reference painting.

Recently, an algorithm based on deep convolutional neural network has been developed to transfer artistic style from an art painting to a photo [2]. Successful as it is in transferring styles from impressionist paintings of artists such as Vincent van Gogh to photos of various scenes, the algorithm is prone to distorting the structure of the content in the source image and introducing artifacts/new

content in the background of the transferred image. We investigate conditions and methods to improve this neural style transfer algorithm to be content-aware along with the goal of synthesizing a high resolution painting in more realism style.

Background

Deep convolutional neural network (CNN) has been extensively applied to various image processing following its initial success in image classification of ImageNet [1]. A CNN typically consists of a sequence of sets of convolutional layers followed by non-linear rectifier and local pooling. Like classical image processing, each convolutional layer performs linear spatial filtering on the input image or the output feature maps from the previous layer; the local pooling down-samples spatial information similar to dyadic down-sampling in classic wavelet decomposition of images. On the other hand, the non-linearity introduced by the rectifier enables CNN to approximate complicated non-linear transforms, mapping images to a feature space where different classes of images become linearly separable. The depth of CNN is proportional to the capacity of the function space that can be approximated, hence a deeper CNN is more “expressive” than a shallower one, encoding more feature information of image classes.

Neural Style Transfer

Due to the richness of features that a deep CNN can possess, pre-trained CNN models on massive datasets like ImageNet with high accuracy in image classification have been used as feature mappings from image domain to abstract feature spaces for other image processing. In particular, the neural style transfer algorithm in [2] is based on VGG-Network(VGG-Net) [5], a near-human performance deep CNN model in image recognition.

Given an input image \mathbf{x} , let $\mathbf{f}_{l,k}(\mathbf{x})$ be the output N_l dimensional feature map¹ of VGG-Net at the l^{th} layer and the k^{th} filter, where N_l is the size of the corresponding spatial grid, then the neural style transfer is formulated as the following optimization problem,

$$\begin{aligned} \mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}'} & \sum_{l \in \mathcal{L}_c} \sum_k \|\mathbf{f}_{l,k}(\mathbf{x}_c) - \mathbf{f}_{l,k}(\mathbf{x}')\|_2^2 \\ & + \lambda \sum_{l \in \mathcal{L}_s} w_l \sum_{i,j} |\mathbf{f}_{l,i}^\top \mathbf{f}_{l,j}(\mathbf{x}_s) - \mathbf{f}_{l,i}^\top \mathbf{f}_{l,j}(\mathbf{x}')|^2 \end{aligned} \quad (1)$$

where \mathbf{x}_c is the source (content) image, \mathbf{x}_s is the reference painting (style image) and \mathbf{x}^* is the style transferred source image. The first term in (1) penalizes the difference between the feature vectors of source image and transferred image at a set of layers \mathcal{L}_c in VGG-Net, such that the content of these two images are

¹this is the same as the term *filter response* used in [2]

close. To measure the closeness of styles, the second term involves the covariance of feature maps \mathbf{f}_l , corresponding to different filters at a fixed layer l in the set \mathcal{L}_s weighted by w_l . This formulation is also used in another closely related neural algorithm for texture synthesis [3], where only a texture (style) image is used to generate a synthetic image presenting similar texture. In the style transfer algorithm, the sets of constraint layers are chosen as $\mathcal{L}_c : \text{conv4_2}$ and $\mathcal{L}_s^{\text{style}} : \text{conv1_1}, \text{conv2_1}, \text{conv3_1}, \text{conv4_1}, \text{conv5_1}$. In the texture synthesis algorithm, $\mathcal{L}_s^{\text{txtr}} : \text{conv1_1}, \text{pool1}, \text{conv2_1}, \text{pool2}, \text{conv3_1}, \text{pool3}, \text{conv4_1}, \text{pool4}$. We defer the comparison of these two settings of constraint layers to the discussion of methods in the next section.

Methods

We propose a content-aware neural style transfer algorithm based on the framework of (1). We first present the base algorithm for the simple case where the source content image and the reference style image are of the same resolution. The algorithm is then generalized to the case where the content image is in lower resolution than the style image. The output of our algorithm is a synthesized painting that could have been painted by the same painter of the reference painting from the content image. We always assume that the content image (photo) and the style image (painting) contain similar objects to minimize the ambiguity of how the generated painting should look like.

Style Layers \mathcal{L}_s

For each layer l in the network, let $s_l = N_{l-1}/N_l$ be the scaling factor of that layer. For a convolutional layer, the scaling factor is 1; for a pooling layer, the scaling factor is its kernel size. Therefore, the scale of features characterized at level l is $\sigma_l = s_1 \cdots s_l$. Specifically, a convolutional layer has the same feature scale as the last pooling layer ahead of it, and the final pooling layer has the maximum feature scale in a network. It is observed in both [3] and [2] that a texture image can be synthesized by fitting the covariance of feature maps \mathbf{f}_l , at layer $l \in \mathcal{L}_s$, namely minimizing the second term in (1) with random initialization. Moreover, the generated texture image is spatially uniform and it represents similar texture of the reference style image at the scales of style layers \mathcal{L}_s . In style transfer [2] and texture synthesis [3], different style layers \mathcal{L}_s are used, but the scales of these layers are the same, namely $1, \sigma_{\text{pool1}}, \sigma_{\text{pool2}}, \sigma_{\text{pool3}}$ and σ_{pool4} . Using layers at different scales collaboratively enforces a hierarchical structural constraint on the texture in the synthesized image. Given a reference style image whose texture is at scale σ_s , that is its energy concentrated within a circle of radius σ_s^{-1} centered at the origin of the frequency spectrum, \mathcal{L}_s should contain layers at scale larger than σ_s to retain full texture pattern.

Although restricting \mathcal{L}_s to layers at smaller scales than σ_s results in less satisfactory result, the statistics of the reference style pattern at finer scales are still preserved in the generated image. In fact, a style pattern at scale σ_s consists

of sub-patterns at finer scales, such as color dots, single or group of brush strokes, etc. Therefore, preserving statistics of $\mathcal{L}_s = \{\text{conv1_1}\}$ results in a synthesized style image consists of color dots, where the portion of areas in different colors is roughly the same as that in the reference image, see results in [2] and [3]. More generally, the sub-patterns at scale $\sigma_l, l \in \mathcal{L}_s$ are preserved in the synthesized image, which is a *re-arrangement* of the composition of reference style image at unconstrained scales. Together with the fact that VGG-Net is not rotation invariant, this insight of re-arrangement leads to an important condition on the input content image and the style image for content-aware style transfer: the object in content image and that in style image should be in the same orientation and in the same relative scale to image size, which is also numerically validated in section Experiments.

Content Layers \mathcal{L}_c and Initialization

In contrast to collaborative multi-scale style layer constraints, a single deep layer is typically used in content constraint. As discussed in [2], the output feature maps from a shallow layer in VGG-Net is highly redundant and sufficient to reconstruct the input image. The similarity between the reconstruction and the input image decreases as the layer whose output feature maps are constrained becomes deeper. Choosing a shallow content layer (e.g. `conv2_1`) thus does not provide much freedom in transferring the content image into a new style, whereas choosing a very deep content layer (e.g. `conv5_1`) does not penalize big distortion of object in the synthesized image. In [2], `conv4_2` is used as content layer constraint together with random initialization of the optimization (1).

This choice of content layer works well for indefinite style transfer where the style and the content are almost independent. As the weight λ before the style penalty term in (1) changes, the result synthesized image varies from presenting highly similar object in the content image with minimal fusion of the reference style to showing only the reference style without the knowledge of content. However, simply adjusting λ does not provide a synthesized image that retains adequate information in both content and style, which is required in content-aware style transfer.

To resolve this problem in content-aware style transfer, we initialize the optimization algorithm with an image close to the content image instead of a random image. It is well known that optimization on deep neural network like (1) using stochastic gradient descent achieves local minimum hence starting from an image close to the content image puts a strong prior on the final optimization result. On the other hand, this change of initialization degrades the style constraint as the potential barrier between the initial point and the set of points (images) \mathcal{I}_s in the reference style is higher than that between \mathcal{I}_s and a random point. In numerical experiments, we observe that solving (1) with fixed λ , and the same $\mathcal{L}_s, \mathcal{L}_c$ as that used in [2], the synthesized image is much more similar to the content image if initialized with the content image than with random initialization. Therefore, we enhance the style constraint and weaken the content penalty in (1) by setting $\mathcal{L}_s : \text{conv1_1, pool1, conv2_1, pool2, conv3_1, pool3, conv4_1, pool4, conv5_1}$

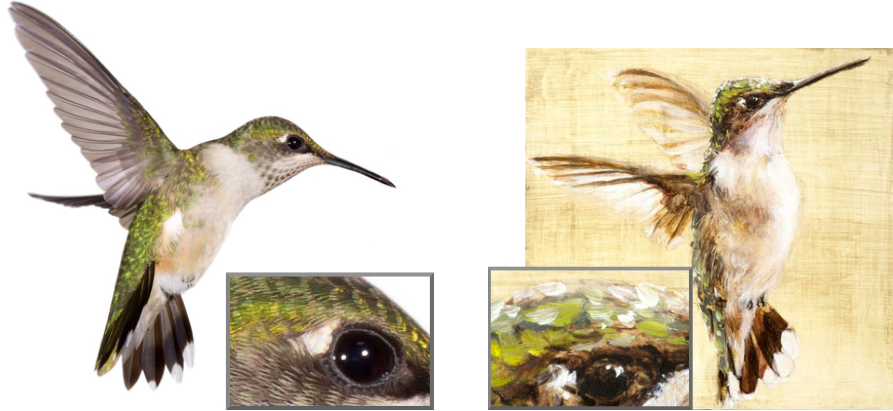


Figure 1: Left: content image, a photo of humming bird with white background (1896×1864 pixel), right: reference painting, a humming bird painted on a square wood board (2024×1947 pixel)

and $\mathcal{L}_c : \text{conv5_2}$, where \mathcal{L}_s contains layers at all scales and \mathcal{L}_c contains one layer at coarsest scale. As $\lambda \rightarrow \infty$, (1) degenerated to the texture synthesis formulation. As long as we don't initialize randomly and keep a mild penalization on content in the objective function, the result image fully combines style and content from the inputs.

Experiments

In this section, we demonstrate style transfer algorithms with different configurations of style layers \mathcal{L}_s , content layers \mathcal{L}_c , optimization weight parameter λ , initialization and input alignment. We first show results on same-resolution style transfer, where the style pattern of the input style image is in almost the same resolution as that of the content image. Based on the best configuration for content-aware style transfer in this easier case, we then discuss its extension to the super-resolution case, where the style pattern of the input style image is in much higher resolution than that of the content image.

Data and Preprocessing

In the following numerical experiments, a pair of content image and style image shown in Fig.1 have been used. The style image is a high resolution scan of a real painting of a humming bird by Charlotte Caspers,² and it's available to download at http://services.math.duke.edu/~rachel/charlotte_new2/

² this painting is selected from a collection of paintings painted by Charlotte Caspers in 2012, served as a public research dataset for digital art authenticity analysis. The dataset is available at <http://services.math.duke.edu/~rachel/resources/charlotte2012.html>

Nr2_original.tiff. The content image was manually found through Google search, where the foreground humming bird is in similar type and pose to that in the painting. The content photo is originally 1896×1864 pixel and the painting is rescaled to 2024×1947 from 4049×3894 , such that it has the same resolution as the content image. We use the implementation of the style transfer algorithm in [2] on Github [4], which allows flexible configuration of the original algorithm including the one we propose.



Figure 2: synthesized image by basic style transfer using the content image and style image in Fig.1 (both input images are down-sampled to 512×503 pixel)

Due to the size of the input content image and style image and the limitation of our computation resource, it is impossible to directly take the content image and style image as input and synthesize a style transferred image. To understand the potential problems of applying the original style transfer algorithm to the full content image and style image, we run the algorithm on input images of decreased size (512×503). The result is shown in Fig.2, where the background is not well synthesized and the color of the body of the humming bird is not close to that in the painting. These problems are not related to the resolution of the input images, but they are closely related to the initialization as we will show below, as we will show in experiments of same-resolution style transfer.

In order to generate synthesized image of the same resolution of the input image, the content image is segmented into major parts such as head, body, tail and wings, see Fig. 3. Each part overlaps with its adjacent parts such that a full-size style transferred image can be reconstructed by merging the synthesized image on parts. The foreground object in the painting is segmented into parts and accordingly the content image. The segmented parts in both images are in the same relative spatial configuration except that for part (f), the secondary wing in the content image, is matched with a sub-region of of part (c), the major wing in the style image for better consistency. We thus focusing on content-aware style transfer on each part instead of the full image in our experiments.

Same-resolution Style Transfer

The goal of same-resolution style transfer is to transfer the style pattern from the style image in its original resolution to the content image. Moreover, when the style is content-dependent, e.g. the brush stroke and wood grain pattern in background is different from the brush stroke used to depict foreground object, the algorithm should transfer the background pattern to the background of the

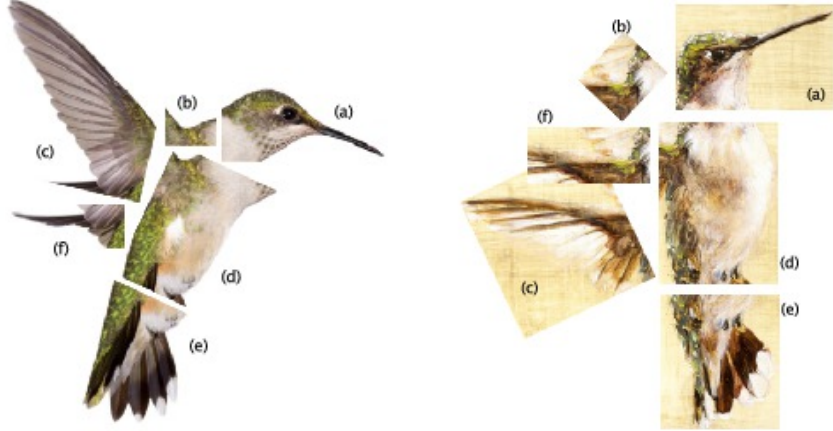


Figure 3: cropped corresponding parts from content image and style image in Fig.1: (a) head, (b) nape, (c) full wing, (d) body, (e) tail, (f) second wing.

content image and the foreground style to the foreground.

We first show the necessity of input alignment and non-random initialization. Consider transferring style on part (a), the head of the bird. a-1 and a-2 in Fig.4 are cropped part (a) from the full content image and style image respectively. The difference of detailed features between a-1 and a-2 are shown in the middle columns of Fig.4. The painting doesn't capture very fine details in the photo, such as the pink lump on the tip of the beak and the vein of feathers on the head. Besides the local dissimilarity in texture, the two input images are not globally aligned in the same direction. For input with alignment, we use a-1' in Fig.4 as the content image together with a-2 in our experiment.

We consider the following four configurations solving (1): (i) basic style transfer, i.e. the original algorithm in [2], (ii) texture enhanced style transfer, where the style layers $\mathcal{L}_s^{\text{txtr}}$ in [3] is used instead of $\mathcal{L}_s^{\text{style}}$, (iii) texture enhanced with input alignment, which is the same as (ii) but use a-1' instead of a-1 as content image, (iv) content-aware style transfer without input alignment, (v) content-aware style transfer, our proposed configuration. Fig.5 shows the results and a table of these four configurations.

As shown in Fig.5, the results with random initialization, the first to third rows, suffer from un-even lightening in the background despite of the difference in style constraint and input alignment. This is because that the covariance of feature maps, which encodes style information, is estimated globally across the full image in (1) not distinguishing the foreground and the background. As the style penalty term is not content-aware, the style of the foreground and the background in the synthesized image can be mixed. In particular, the background of the painting is uniform of small variance, and that of the



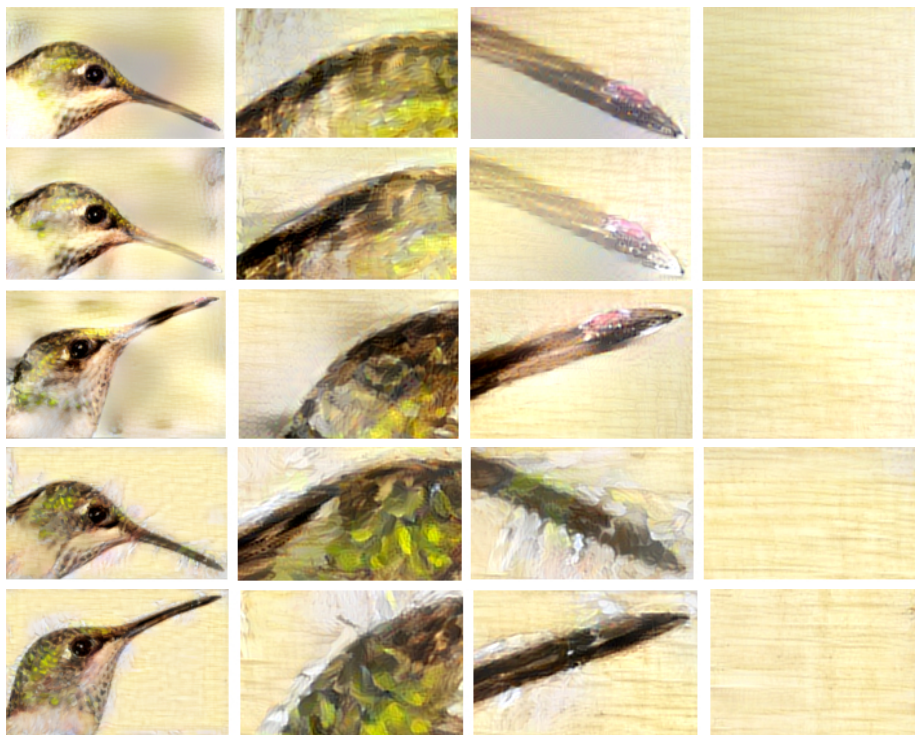
Figure 4: a-1: part (a) of content image, a-2: part (a) of style image, a-1': rotated a-1 aligned in the same direction of a-2; second column: a zoomed in patch at the upper left of the head; third column: a zoomed in patch of the beak; fourth column in the bottom row: a zoomed in background patch in a-2.

synthesized image has larger variance due to the constraint on variance of color channels in style layers \mathcal{L}_s mainly contributed by the foreground object in the painting. Therefore, the spatial information contained in the initial point is a strong regularization, and the results initialized with the content images have uniform lightening in the background as in the style image.

Furthermore, the results with input alignment (i.e. using a-1' and a-2 as input), the third and fifth rows, have better brush stroke simulation than those without alignment (i.e. using a-1 and a-2 as input). In particular, the direction of brush strokes of the beak aligns in the same direction of the beak when the inputs are aligned, but not when mis-aligned, see the third column of Fig.5. This is due to the fact that VGG-Net is not rotation-invariant and the same object with different orientations are encoded in different filter channels in the layers of the network, especially in the shallow layers, where feature information is not highly aggregated. When the content is consistent with the style image, it cannot be consistent with the style image, so are their feature maps in shallow layers, unless both input images are in the same orientation. On the other hand, the shallow layers cannot be dropped from \mathcal{L}_s for reconstruction of texture in fine scale.

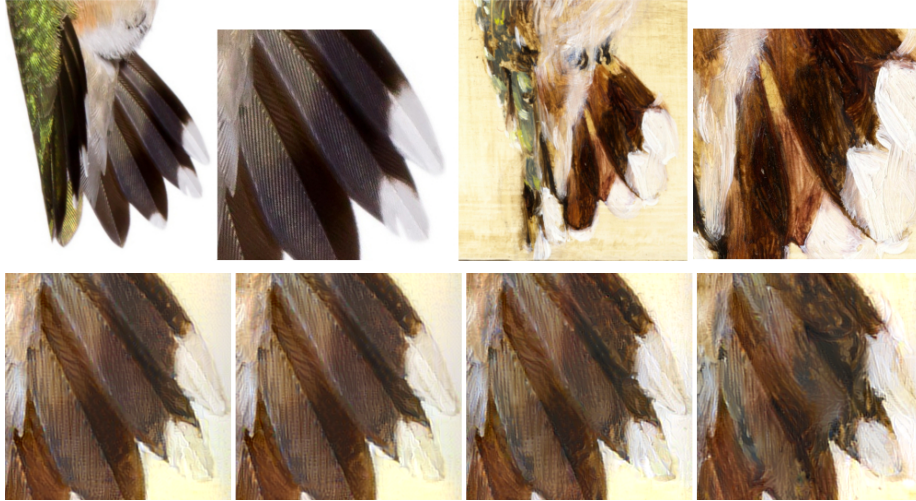
In addition, including pooling layers in \mathcal{L}_s helps to enhance textures from the style image and diminish fine details from the content image in the synthesized image. See the first and second rows in Fig.5 for comparison of basic configuration and the texture enhanced configuration.

Once the style-transfer algorithm is initialized with the content image, this strong prior on content limits the searching region of the algorithm in the image space. To ensure that a good style-transferred image is achievable, it's necessary to weaken the content constraint term in (1) and set \mathcal{L}_c to a deeper layer at a



configuration	style layers	content layers	λ	initialization	input alignment
I	$\mathcal{L}_s^{\text{style}}$	conv4_2	20	random	no
II	$\mathcal{L}_s^{\text{txtr}}$	conv4_2	20	random	no
III	$\mathcal{L}_s^{\text{txtr}}$	conv4_2	20	random	yes
IV	$\mathcal{L}_s^{\text{txtr}}$ + conv5_1	conv5_2	200	content image	no
V	$\mathcal{L}_s^{\text{txtr}}$ + conv5_1	conv5_2	200	content image	yes

Figure 5: Figures shown in columns from left to right: synthesized image from part (a) in the content photo and the reference painting, a zoomed in patch at the upper left of the head, a zoomed in patch of the beak, a zoomed in background patch in the upper right (the first, second and fourth rows)/bottom right (the third and fifth rows) corner. The results in rows from top to bottom correspond to configurations I to V in the table above.



configuration	style layers	content layers	λ	initialization	input alignment
VI	$\mathcal{L}_s^{\text{txtr}}$	conv4_2	20	content image	yes
VII	$\mathcal{L}_s^{\text{txtr}} + \text{conv5}_1$	conv4_2	20		
VIII	$\mathcal{L}_s^{\text{txtr}} + \text{conv5}_1$	conv4_2	200		
V	$\mathcal{L}_s^{\text{txtr}} + \text{conv5}_1$	conv5_2	200		

Figure 6: Top left: part (e) of content image and a zoomed in patch of tail, top right: part (e) of style image and a zoomed in patch of tail; bottom row: zoomed in patch of tail in synthesized image at the same location as that in the top row corresponding to configurations VI, VII, VIII and V.

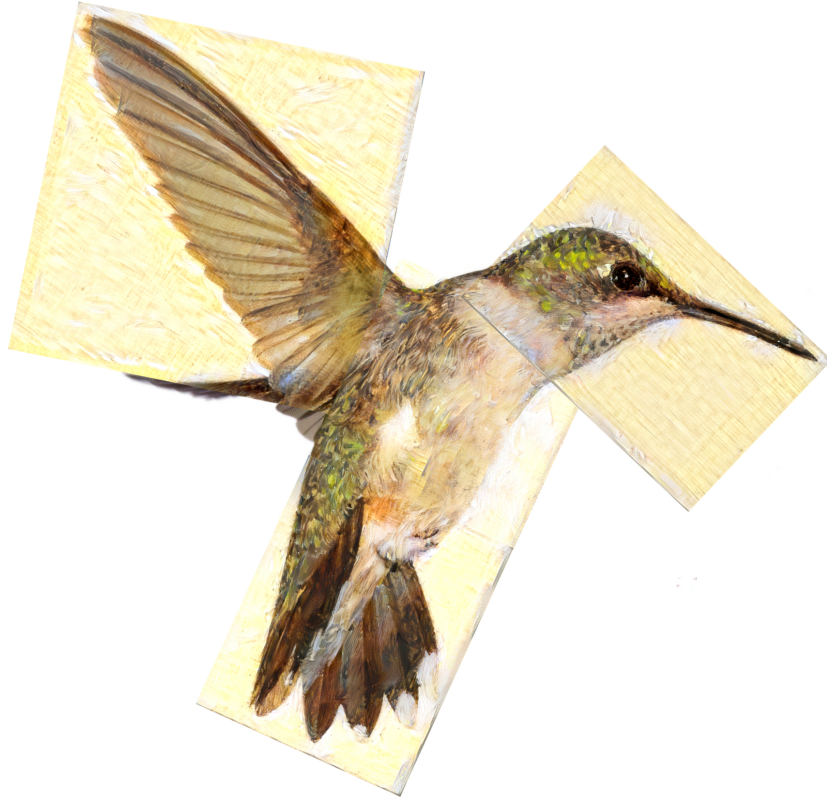


Figure 7: stitched synthesized parts (a), (c), (d) and (e) of content image

coarser scale than `conv4_2`. The second row of Fig.6 compares results of different configurations with aligned input and content-based initialization. Configuration VI to VIII provide almost the same results (the first three columns) where the feather veins are more prominent than that generated by configuration V, our proposed content-aware style-transfer configuration (the last column).

We apply the content-aware style transfer to pairs of content and style images for part (a) to (f). Fig.7 shows an overlay of the synthesized result on part (a), (c), (d) and (e) up on the content photo. Fig.8 shows style transfer result on the whole object, which is obtained by overlaying the synthesized results of each part over the content photo and merging the overlapping regions between parts³.

³The merging process is done in GIMP using the blending tool, which creates a mask on each part image with a smooth decay on the overlapping boundaries with other part images



Figure 8: merged synthesized parts (a), (b), (c), (d), (e) and (f) of content image

There are some minor artifacts in the result shown in Fig.8. The background color is not consistent within parts and there are white brush strokes in the background of part (a), (c), (d) and (f). Compare Fig.8 with the photo and the painting in Fig.1, the humming bird in the photo is slimmer than the painted bird who wears more white feathers, and the extra white feathers in the painting are redistributed to the background in the synthesized image. Although the white feathers is not in the photo, it is introduced by the style constraint to match the covariance of feature maps, resulting in a bias toward the object in the style image. The framework (1) does not achieve a strict separation of content and style information in a neural network. Futhermore, because of the spatial-independent property of the style constraint, there is no restriction where the white brush strokes could appear in the background. Next, we extend this algorithm for same-resolution content-aware style transfer to the super-resolution case.

Super-resolution Style Transfer

In the super-resolution case, we use the same photo and painting as previously discussed but downscale the photo to one-fourth of its original size. We apply the following iterative algorithm to generate a synthesized image of the same resolution of the style image:

Input: content image I_c^0 at scale α_0 , style image I_s^K at scale α_K

1. downsize I_s^K to I_s^k at scale α_k , $k < K$, s.t. $\alpha_{k-1} < \alpha_k < \alpha_{k+1}$

for $k = 0, \dots, K$

2.1 generate synthesized image I_g^k transferring the style of I_s^k to I_c^k

2.2 upscale I_g^k to α_{k+1} and set as I_c^{k+1}

Output: synthesized image I_g^K

Fig.9 shows the generation of I_g^k along in the super-resolution algorithm. The final output is more similar to the painting than the synthesized image in the same-resolution case. In fact, the extra white feathers in the painting emerges gradually in I_g^k as the scale increases. The iterative algorithm can regularize where to add the extra feature because the global structure of the object is captured and preserved through scales. In the first iteration, the VGG-Net encodes global structure of the downsized input images I_c^0 , I_s^0 at deep layers after several spatial pooling layers, and this structural information is kept in I_g^0 and passed on to the next iteration. In the k^{th} iteration, the result of the last iteration is upscaled by a factor of α_k/α_{k-1} and used as the new content image to generate I_g^k in higher resolution. On the other hand, some fine content features such as gray spots on the neck of the bird diminishes along the iteration. The level of degradation in the k^{th} iteration is proportional to α_k/α_{k-1} , or equivalently how underdetermined is the optimization.

Conclusion

In this paper, we show that the formulation of the optimization (1) in [2] doesn't account for spatial dependency of style patterns in the reference style image, and the basic style transfer algorithm is not content-aware, which is resolved by setting a strong initial prior on content while weaken the content constraint in the optimization algorithm. In addition, we show that it's necessary to align the input images to the same orientation, otherwise the induced content constraint and style constraint are not consistent as the VGG-Net is not rotation invariant. Furthermore, we extend this content-aware configuration into an iterative algorithm that transfers high resolution style to a low-resolution content image. Our proposed configuration for content-aware style transfer significantly improves the basic algorithm, and the extra features from the reference style image present in the synthesized image introduced by the style constraint means that the content and the style information is not completely separated by neural network.

References

- [1] Jia Deng, Wei Dong, Richard Socher, Li jia Li, Kai Li, and Li Fei-fei. Imagenet: A large-scale hierarchical image database. In *In CVPR*, 2009.
- [2] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. Aug 2015.
- [3] L. A. Gatys, A. S. Ecker, and M. Bethge. Texture synthesis using convolutional neural networks. In *Advances in Neural Information Processing Systems 28*, May 2015.
- [4] jcjohnson. neural-style. <https://github.com/jcjohnson/neural-style>, 2015.
- [5] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

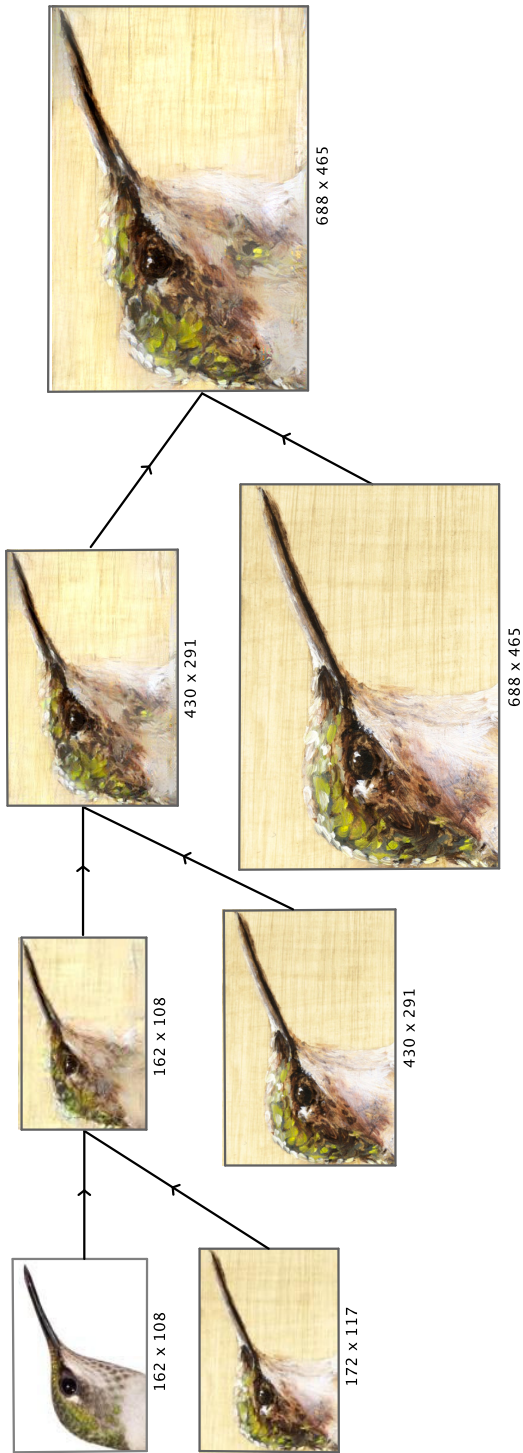


Figure 9: super-resolution style transfer on reduced-size content photo with a series of painting increasing in size