

# Multimodal Transfer: A Hierarchical Deep Convolutional Neural Network for Fast Artistic Style Transfer

Xin Wang<sup>1,2</sup>, Geoffrey Oxholm<sup>2</sup>, Da Zhang<sup>1</sup>, Yuan-Fang Wang<sup>1</sup>

<sup>1</sup>University of California, Santa Barbara, CA

<sup>2</sup>Adobe Research, San Francisco, CA

{xwang, dazhang, yfwang}@cs.ucsb.edu, oxholm@adobe.com

## Abstract

*Transferring artistic styles onto everyday photographs has become an extremely popular task in both academia and industry. Recently, offline training has replaced on-line iterative optimization, enabling nearly real-time stylization. When those stylization networks are applied directly to high-resolution images, however, the style of localized regions often appears less similar to the desired artistic style. This is because the transfer process fails to capture small, intricate textures and maintain correct texture scales of the artworks. Here we propose a multimodal convolutional neural network that takes into consideration faithful representations of both color and luminance channels, and performs stylization hierarchically with multiple losses of increasing scales. Compared to state-of-the-art networks, our network can also perform style transfer in nearly real-time by conducting much more sophisticated training offline. By properly handling style and texture cues at multiple scales using several modalities, we can transfer not just large-scale, obvious style cues but also subtle, exquisite ones. That is, our scheme can generate results that are visually pleasing and more similar to multiple desired artistic styles with color and texture cues at multiple scales.*

## 1. Introduction

Style transfer, or to repaint an existing photograph with the style of another, is considered a challenging but interesting problem in arts. Recently, this task has become an active topic both in academia and industry due to the influential work by Gatys *et al.* [8], where a pre-trained deep learning network for visual recognition is used to capture both style and content representations, and achieves visually stunning results. Unfortunately, the transfer run time is prohibitively long because of the online iterative optimization procedure. To resolve this issue, a feed-forward

network can be trained offline with the same loss criterion to generate stylized results that are visually close (but still somewhat inferior). In this way, only one single inference pass of the feed-forward network is needed at the application time. This results in a computational algorithm that is hundreds of times faster [13, 25].

Though past work creates visually pleasing results for many different types of artworks, two important drawbacks stand out: (1) The current feed-forward networks [13, 25] are trained on a specific resolution of the style image, so deviating from that resolution (bigger or smaller) results in a scale mismatch. For example, applying a model trained with a style guide of size 256 on higher-resolution images would generate results whose texture scale is smaller than that of the artistic style, and (2) Current networks often fail to capture small, intricate textures, like brushwork, of many kinds of artworks on high-resolution images. While it has been shown that these feed-forward networks function quite well on artworks with abstract, large-scale textures and easily discernible strokes, *e.g.*, *The Starry Night* by Vincent van Gogh, artistic styles are much more encompassing than what has been demonstrated. That is, different artistic styles may be characterized by exquisite, subtle brushes and strokes, and hence, our observations are that results of these style-transfer networks are often not satisfactory for a large variety of artistic styles.

In this paper, we propose a novel hierarchical deep convolutional neural network architecture for fast style transfer. Our contribution is fourfold: (1) We introduce a hierarchical network and design an associated training scheme that is able to learn both coarse, large-scale texture distortion and fine, exquisite brushwork of an artistic style by utilizing multiple scales of a style image; (2) Our hierarchical training scheme and end-to-end CNN network architecture allow us to combine multiple models into one network to handle increasingly larger image sizes; (3) Instead of taking

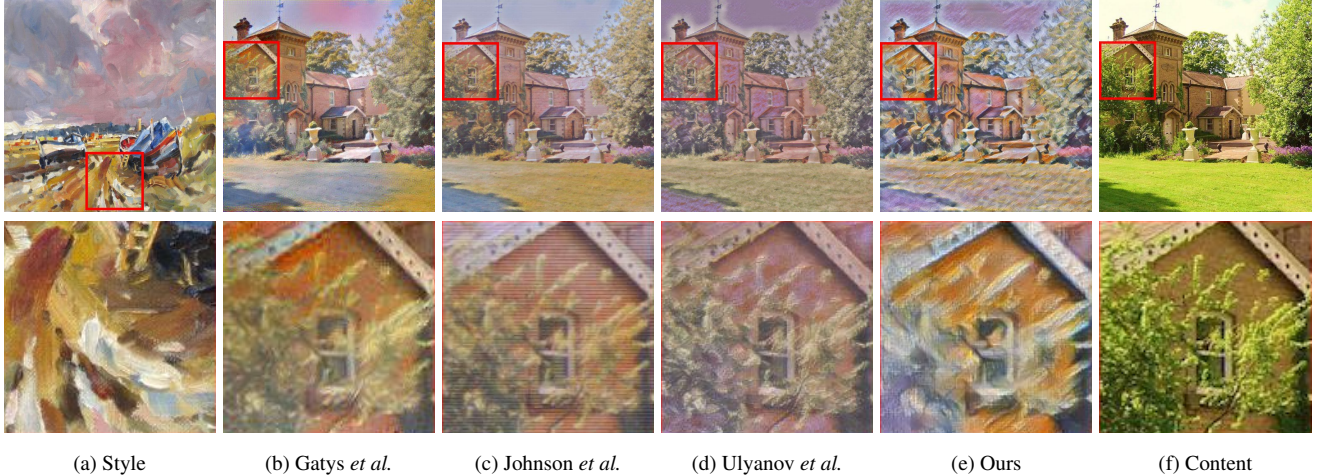


Figure 1: Top row: (a) The style guide is *At the Close of Day* by Tomas King, and (f) is the content image. (b) is the result of Gatys *et al.*'s optimization-based method. (Result size is 512 due to memory limitation of the method) (c), (d) and (e) are results generated by different feed-forward networks (all are of size 1024). Bottom row: the zoom-in display of the regions enclosed in the red boxes from the top row. As can be seen, all results are repainted with the color of the style image. However, a closer examination shows that the brush strokes are not captured well in (c) and (d). The zoom-in region in (b) is a little blurry. Comparing with the others, our multimodal transfer (e) is capable of simulating more closely the brushwork of the original artwork on high-resolution images.

only RGB color channels into consideration, our network utilizes representations of both color and luminance channels for style transfer; and (4) Through experimentation, we show that our hierarchical style transfer network can better capture both coarse and intricate texture patterns.

Our hierarchical style transfer network is trained with multiple stylization losses at different scales using a mixture of modalities, so we distinguish it as *multimodal transfer* from the feed-forward style transfer networks with only one stylization loss [13, 25], which we call *singular transfer*. In Fig. 1 we give an example that compares results from our multimodal transfer network with those from the current state-of-the-art singular transfer networks. Fig. 1 shows the advantages of multimodal transfer on learning different levels of textures, including style, color, large texture distortion and fine brushwork. Note specifically that our method can simulate more closely the brushwork of the artwork. In Sec. 4 we will show that multimodal transfer can also be used to train a combination model to stylize a single image with multiple, distinct artistic styles.

## 2. Related Work

**Understanding representations of deep neural networks.** Recently, seminal work was done on understanding deep neural networks. The DeconvNet method of Zeiler and Fergus [30] learns how certain network outputs are obtained by identifying which image patches are responsible for certain neural activation. Yosinski *et al.* [29] aims to understand what computation is performed by deep networks through visualizing the internal neurons. Mahendran and

Vedaldi [19] invert the image representations of certain layers to learn what information is preserved by the networks. The latter two approaches generate visualization images with an optimization procedure whose objective is for perceptual understanding of network functions. Similar optimization procedure is also adopted in other cases [23, 20].

Based on the better understanding of the powerful representations of deep convolutional networks [15], many traditional vision tasks have been addressed with much more improved outcomes. Optimization-based style transfer is one such example. Different from previous texture synthesis algorithms that are usually non-parametric methods [5, 28, 4, 11, 1, 16, 17], Gatys *et al.* first proposed an optimization method of synthesizing texture images where the objective loss is computed based on the representations of a pre-trained convolutional neural network [6]. This texture loss is then combined with content loss derived from Mahendran and Vedaldi [19] to perform the style transfer task [8].

**Feed-forward networks for image generation.** The optimization-based methods for image generation are computationally expensive due to the iterative optimization procedure. On the contrary, many deep learning methods use the perceptual objective computed from a neural network as the loss function to construct feed-forward neural networks to synthesize images [3, 9, 2, 22].

Fast style transfer has achieved great results and is receiving a lot of attention. Johnson *et al.* [13] proposed a feed-forward network for both fast style transfer and super-

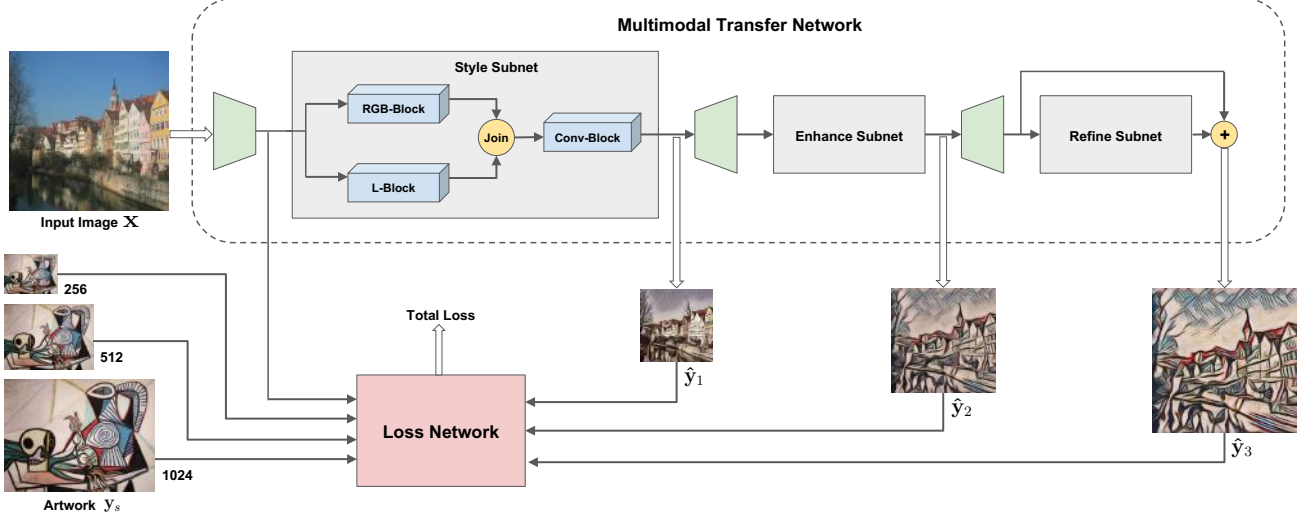


Figure 2: Overall Architecture. Please see Sec. 3.1 for explanation.

resolution using the perceptual losses defined in Gatys *et al.* [8]. A similar architecture *texture net* is introduced to synthesize textured and stylized images [25]. More recently, Ulyanov *et al.* [26] shows that replacing spatial batch normalization [12] in the feed-forward network with instance normalization can significantly improve the quality of generated images for fast style transfer. Here we present further improvement of such style transfer algorithms to handle progressively larger images using hierarchical networks with mixed modalities. Furthermore, it allows the use of multiple, distinct styles for repainting a single input image.

### 3. Multimodal Transfer Network

#### 3.1. Overall Architecture and Learning Schemes

Our proposed network, which is shown in Fig. 2, is comprised of two main components: a feed-forward multimodal network and a loss network. The feed-forward multimodal network (*MT Network*) is a hierarchical deep residual convolutional neural network. It consists of three subnetworks: *style subnet*, *enhance subnet* and *refine subnet*. These subnets are parameterized by  $\Theta_1$ ,  $\Theta_2$ , and  $\Theta_3$  respectively (these parameters will be made explicit later). At a high level, the MT Network takes an image  $\mathbf{x}$  as the input and is trained to generate multiple output images  $\hat{\mathbf{y}}_k$  of increasing sizes,

$$\hat{\mathbf{y}}_k = f(\cup_{i=1}^k \Theta_i, \mathbf{x}) \quad . \quad (1)$$

These output images are then taken separately as inputs to the loss network to calculate a stylization loss for each. The total loss is a weighted combination of all stylization losses. We will show later in Sec. 3.2 the loss network and the definition of the total loss.

At test time, in order to produce the same stylization effect and correct texture scale of the artworks when applied to larger images, the MT network stylizes the image hierarchically: The input image is first resized into 256 with a bi-linear downsampling layer and stylized by the *style subnet*, capturing the large color and texture traits of the artwork. Next the stylized result, which is the first output  $\hat{\mathbf{y}}_1$ , is up-sampled into 512 and transferred to the output  $\hat{\mathbf{y}}_2$  by the *enhance subnet*, which enhances the stylization strength. Then it is resized back to 1024. Finally, the *refine subnet* removes the local pixelization artifacts and further refines the result. The high-resolution and most visually appealing result  $\hat{\mathbf{y}}_3$  is obtained after these three-stage processing. Note that while we illustrate the process using a two-level hierarchy, the same concept can be extended recursively to enable stylization of progressively larger images.

#### 3.2. Loss Functions

In this section, we first introduce the single stylization loss function and then present a hierarchical stylization loss function that is adopted to train our multimodal transfer network.

##### 3.2.1 Single Stylization Loss Function

Similar to the loss definition in previous work for fast style transfer [13, 25], the stylization loss is also derived from Gatys *et al.* [8], where a loss network (a pre-trained VGG-19 network optimized for object recognition [24]) is used to extract the image representations.

Two perceptual losses are defined to measure to what extent the generated image  $\hat{\mathbf{y}}_k$  combines the content of the content target  $\mathbf{y}_c$  with the texture and style cues of the style target  $\mathbf{y}_s$  (see Fig. 3).

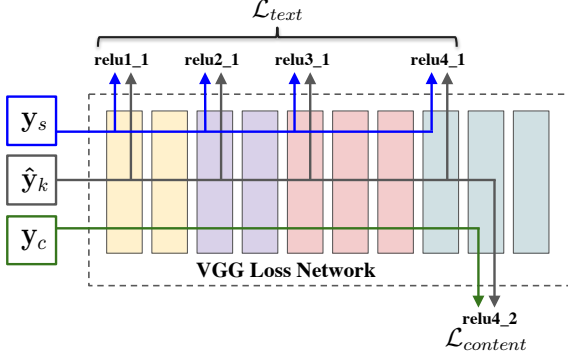


Figure 3: Loss network. Please see Sec. 3.2 for explanation.

**Content Loss** The content loss function is used to measure the dissimilarity between  $\hat{\mathbf{y}}_k$  and  $\mathbf{y}_c$ . Let  $F_i^l(\mathbf{x})$  denote the  $i$ -th feature map in the  $l$ -th layer of the loss network applied to image  $\mathbf{x}$ . The content loss is the squared-error loss between the two feature representations at layer  $l$

$$\mathcal{L}_{content}(\hat{\mathbf{y}}_k, \mathbf{y}_c, l) = \sum_{i=1}^{N_l} \|F_i^l(\hat{\mathbf{y}}_k) - F_i^l(\mathbf{y}_c)\|_2^2 \quad (2)$$

That is, the content loss directly compares the feature maps computed from the corresponding layers and thus is suitable for characterizing spatial content similarity.

**Texture or Style Loss** Gatys *et al.* propose that the correlations between feature maps in each layer of the loss network can be seen as texture representations of an image [6, 8]. Those correlations are given by the Gram matrix, whose elements are pairwise scalar products between those feature maps:

$$G_{ij}^l(\mathbf{x}) = \langle F_i^l(\mathbf{x}), F_j^l(\mathbf{x}) \rangle \quad (3)$$

A set of Gram matrices  $G^l, l \in L$  is used as the texture representations, which discard the spatial information but retain the statistic profiles of color and intensity distribution of an input image. So the texture loss function is defined as

$$\mathcal{L}_{text}(\hat{\mathbf{y}}_k, \mathbf{y}_s) = \sum_{l \in L} \|G^l(\hat{\mathbf{y}}_k) - G^l(\mathbf{y}_s)\|_2^2 \quad (4)$$

Finally, the stylization loss for each output  $\hat{\mathbf{y}}_k$  from the MT network is defined as a weighted sum of the content loss and the texture loss

$$\mathcal{L}_S(\hat{\mathbf{y}}_k, \mathbf{y}_c, \mathbf{y}_s) = \alpha \mathcal{L}_{content}(\hat{\mathbf{y}}_k, \mathbf{y}_c) + \beta \mathcal{L}_{text}(\hat{\mathbf{y}}_k, \mathbf{y}_s), \quad (5)$$

where  $\alpha$  and  $\beta$  are the weights of the content loss and texture loss, respectively.

### 3.2.2 Hierarchical Stylization Loss Function

The multimodal transfer network can generate  $K$  output results of  $K$  increasing sizes ( $K = 3$  in the network shown in Fig. 2). Then a stylization loss is computed for each output result  $\hat{\mathbf{y}}_k$

$$\mathcal{L}_S^k(\hat{\mathbf{y}}_k, \mathbf{y}_c, \mathbf{y}_s) = \alpha \mathcal{L}_{content}(\hat{\mathbf{y}}_k, \mathbf{y}_c) + \beta \mathcal{L}_{text}(\hat{\mathbf{y}}_k, \mathbf{y}_s) \quad (6)$$

where  $\mathbf{y}_c^k$  and  $\mathbf{y}_s^k$  are the corresponding content target and style target, which are the input to the subnet that outputs  $\hat{\mathbf{y}}_k$ , and are the scaled versions of the artwork  $\mathbf{y}_s$ . By training the subnets with different style scales, we control the types of artistic features that are learned for different subnets. Again, we want to emphasize that the concept can be easily extended for more layers.

Since such stylization losses are computed based on the outputs of different layers of the whole network, a total loss (e.g., a weighted combination of all stylization losses) cannot be used here to directly propagate and update the weights backward. Thus, a parallel criterion is adopted so that different stylization losses are used to back-propagate the weights for different ranges of layers. We define the hierarchical stylization loss function  $\mathcal{L}_H$ , which is a weighted sum of such stylization losses, as

$$\mathcal{L}_H = \sum_{k=1}^K \lambda_k \mathcal{L}_S^k(\hat{\mathbf{y}}_k, \mathbf{y}_c^k, \mathbf{y}_s^k) \quad (7)$$

where  $\lambda_k$  is the weight of stylization loss  $\mathcal{L}_S^k$ .

Therefore, during the end-to-end learning on natural images  $\mathbf{x} \sim \mathcal{X}$ , each subnet denoted by  $\Theta_k$  is trained to minimize the parallel weighted stylization losses that are computed from the latter outputs  $\hat{\mathbf{y}}_i$  ( $i \geq k$ ) (latter means it comes later in the feed-forward direction) as in

$$\Theta_k = \arg \min_{\Theta_k} E_{\mathbf{x} \sim \mathcal{X}} \left[ \sum_{i \geq k}^K \lambda_i \mathcal{L}_S^i(f(\cup_{j=1}^k \Theta_j, \mathbf{x}), \mathbf{y}_c^i, \mathbf{y}_s^i) \right] \quad (8)$$

In practice, suppose the general back-propagation function is denoted by  $f^{-1}$ , then for every iteration, the weight updates (gradients) of the subnet  $\Theta_k$  can be written as

$$\Delta \Theta_k = \begin{cases} f^{-1}(\lambda_k \mathcal{L}_S^k) & k = K \\ f^{-1}(\lambda_k \mathcal{L}_S^k, \Delta \Theta_{k+1}) & 1 \leq k < K \end{cases} \quad (9)$$

so the weights of the current subnet  $\Theta_k$  are influenced by both the stylization loss at the current level  $\mathcal{L}_S^k$  and the gradients of the latter subnets.

From Eq. (8), we can see that even though all those subnets are designed for different purposes, they are not totally independent. Former subnets also contribute to minimize losses of the latter. Thus, shallower CNN structure can be used for latter subnets, which saves both computing memory and running time.



### 3.3. Network Architecture

One of the key drawbacks of singular transfer networks (e.g. [13, 25]) is that the scale at which the singular transfer network is trained limits the range of style details that are captured. Since it is trained with one particular scale of the style image, during training we need to choose if it learns the coarse texture or the fine brushwork. That is, it learns one at the expense of the other.

To remedy this problem, we design the hierarchical architecture where different subnets are trained with different scales of the style image to learn different levels of artistic texture cues. This design enables a test image to be transferred using different levels of the style in increasing resolutions. Furthermore, because all these subnets are combined into one network and trained hierarchically, the latter subnets are also able to enhance and refine the results from previous ones, making ours a collaborative scheme for improved efficiency and robustness.

We have experimented with several architectures that have varying levels of hierarchy and different internal structures. Here we introduce the general architecture of the network shown in Fig. 2, which has the best stylization quality from our experience.

As stated before, the multimodal transfer network consists of three learnable subnetworks, style subnet, enhance subnet and refine subnet, each following a fixed bilinear up-sampling/downsampling layer. Note that the upsampling layer between enhance subnet and refine subnet is only inserted at test time, so during training the input to refine subnet is still of size 512, which hugely reduces the required memory and speeds up the training process. The salient features of these networks are explained below.

#### 3.3.1 Style Subnet

**Luminance-Color Joint Learning** To better address the issue of preserving small intricate textures, our network utilizes representations of both color and luminance channels, because visual perception is far more sensitive to changes in luminance than in color [27, 11, 7]. We separate the luminance channel from the RGB color image and use two independent branches (*RGB-Block* and *L-Block*) to learn their representations distinctively. The feature maps calculated from both branches are then joined together along the depth dimension and further processed by the ensuing *Conv-Block*.

*RGB-Block* comprises three strided convolutional layers ( $9 \times 9, 3 \times 3, 3 \times 3$  respectively, the latter two are used for downsampling) and three residual blocks [10], while *L-Block* has a similar structure except that the depth of convolution is different. *Conv-Block* is composed of three residual blocks, two *resize-convolution* layers for upsampling and the last  $3 \times 3$  convolutional layer to obtain the output RGB



Figure 4: (a) is the content image and (b) is the style image. (c)(d)(e) show the outputs of the three subnets,  $\hat{y}_1$ ,  $\hat{y}_2$  and  $\hat{y}_3$ , whose sizes are 256, 512 and 1024 respectively. The third row depicts the absolute difference between: (f) the content image and the output image  $\hat{y}_1$ , (g) the output images  $\hat{y}_1$  and  $\hat{y}_2$ , and (h) the output images  $\hat{y}_2$  and  $\hat{y}_3$ .

image  $\hat{y}$ . All non-residual convolutional layers are followed by instance normalization [26] and ReLU nonlinearity. Part of our style subnet is designed based on the work [13, 22].

A nearest neighbor interpolation upsampling layer and a convolutional layer called *resize-convolution layer* is used here instead of deconvolutions to avoid the checkerboard artifacts of generated images [21].

#### 3.3.2 Enhance Subnet and Refine Subnet

Although the style subnet is intended to stylize the input image with large texture distortion to match that of the style guide, we have found that it is difficult to optimally adjust texture and content weights to achieve style transfer while preserving the content for a large variety of styles. Thus, we allow the style subnet to perform texture mapping with an eye toward preserving the content, and train a separate enhance subnet with a large texture weight to further enhance the stylization. Fig. 4 illustrates the specific role of each subnet. Evidently, the style subnet changes both color and texture heavily, but the enhance subnet also contributes greatly to the texture mapping while adding more detail. The refine net further refines and adds more detail into the final result.

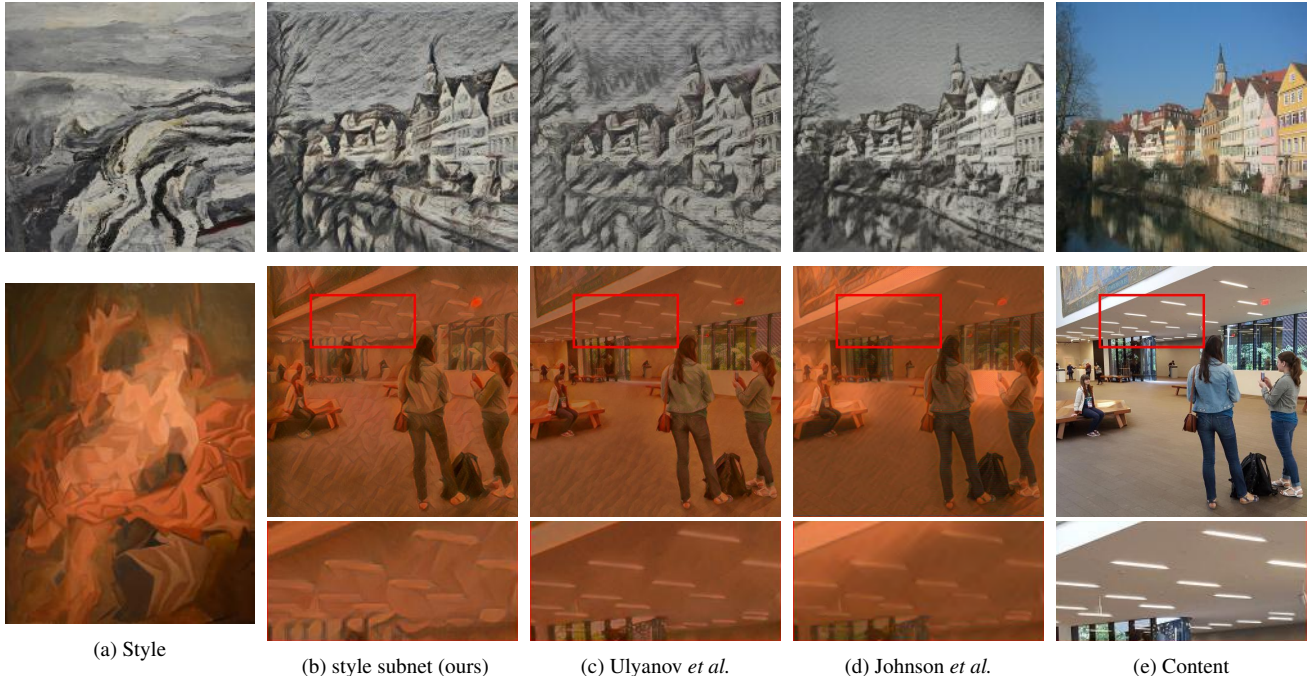


Figure 5: Comparison between our style subnet (b) and other singular transfer networks (c) and (d). They were tested on two styles here, *Mountain No. 2* by Jay DeFeo (top) and *Venus* by Manierre Dawson (bottom). All results were of size 512. Notice that for the second style, we also zoomed in on the region in the red box to better compare the fine texture.

Accordingly, for the sake of enhancing the stylization, we adopt a similar structure as the style subnet for the enhance subnet. The only difference is that the enhance subnet has one more convolutional layer for downsampling and one more resize-convolution layer for upsampling, which enlarges the receptive field sizes. This is needed because the input to the enhance subnet is twice larger than that to the style subnet.

Finally, the refine net consists of three convolutional layers, three residual blocks, two resize-convolution layers and one last convolutional layer to obtain the final output, which is much shallower than the style and enhance subnet. This is because from Eq. (1) and (8) we know former subnets can also contribute to the learning tasks of the latter. Shortening the refine subnet is advantageous. It significantly reduces memory and computational complexity, which is critical for images of size 1024. Furthermore, we add an identity connection from its beginning to the end, forcing it to learn just the difference between its input and output.

## 4. Experiments

**Training Details** The MT Network was trained on a subset of the Microsoft COCO dataset [18], which contained 32,059 images (whose width and height were  $\geq 480$ ). We cropped those images and resized them to  $512 \times 512$ . Adam optimization [14] was used to train models for 10,000 iterations with batch size 1. The learning rate was initially set

as  $1 \times 10^{-3}$  and then reduced by a factor 0.8 every 2,000 iterations. Content losses were computed at layer relu4\_2 of VGG-19 and texture losses at layers relu1\_1, relu2\_1, relu3\_1 and relu4\_1 for all subnets. The content weights were all set to 1, while the texture weights depended on different styles, because a universal ratio of texture to content did not fit all artistic styles. As for the weights of stylization losses, we set  $\lambda_1 : \lambda_2 : \lambda_3 = 1 : 0.5 : 0.25$ . The rationale was that, during training, the parameters of former subnets were updated to incorporate the current and latter stylization losses. The latter losses should have smaller weights in order not to totally dominate the optimization process of the former subnets. Experiments revealed, however, that results were fairly robust to changes in  $\lambda_k$ . It took around one hour to fully train a hierarchical model on an NVIDIA GTX 1080. The model size on disk was about 35MB.

### Comparison among Different Singular Transfer Networks

As mentioned before, singular transfer was a feed-forward style-transfer network with a single stylization loss and multimodal transfer was the hierarchical network with multiple stylization losses and a mixture of modalities (color and luminance). Here we separated the style subnet from our MT network as a singular transfer network and compared it with the other state-of-the-art networks by Johnson *et al.* [13] and Ulyanov *et al.* [26]. All three networks were trained on images of size 256 with the same





(a) Style

(b) Singular Transfer (style size 256)

(c) Singular Transfer (style size 1024)

(d) Multimodal Transfer

Figure 6: We compared our multimodal transfer with two singular transfer networks trained with different scales of the style image, 256 and 1024 (the singular transfer networks had the same architecture and had the same number of parameters as our multimodal transfer network). All generated results were  $1024 \times 1024$  pixels. As can be seen, there was a big texture scale mismatch issue in the results of singular transfer 256 (column (b)) that the texture scale was apparently smaller than that of the original artworks. While singular transfer 1024 (column (c)) failed to learn the texture distortion and brushwork although being rendering with the correct color. The results of multimodal transfer (column (d)) resolved those issues and managed to learn both coarse texture and fine, intricate brushwork.





Figure 7: Multimodal transfer with two styles. The model was trained with style (a) *Still Life with Skull, Leeks and Pitcher* by Pablo Picasso, and (b) *At the Close of Day* by Tomas King. (c) was the test image, and (f) was the final stylized result that had large texture distortion from (a) and small, detailed brushwork from (b). For comparison, we gave the results transferred by the models trained on a single style in (d) and (e).

content to texture weights and used instance normalization. Generally speaking, our style subnet generated qualitatively comparable results. Particularly, it performed better than others on capturing texture details in some cases. In Fig. 5 we gave two comparison examples. In the first example, the result of Ulyanov *et al.* was visibly darker than the style image in color, and the texture scale in Johnson *et al.*'s result did not match that of the style image very well. The result of our style subnet seemed better in both aspects. In the second example, comparing with the other two networks, our style subnet performed better on simulating small, detailed texture. Therefore, we chose the style subnet as a representative of singular transfer to be compared with multimodal transfer next.

**Singular Transfer VS Multimodal Transfer on High-resolution Images** We tested our method on numerous artistic styles. In Fig. 6, we compared our multimodal transfer network on high-resolution images ( $1024 \times 1024$ ) with a singular transfer network with the same number of learning weights. (More exactly, we duplicated our style subnet to a deeper network that had the same number of parameters to provide a fair comparison). Examining the results shown in Fig. 6, comparing with singular transfer, multimodal transfer results were visually more similar to the original artistic styles both in coarse texture structure and fine brushwork, while singular transfer with the style size 256 caused a texture scale mismatch that the texture scale was much smaller than that of the original artwork. Furthermore, singular transfer with style size 1024 failed to learn the distortion

and fine brushwork.

**Multimodal Transfer with Multiple Styles** Our multimodal transfer allowed an interesting application that was not possible before: It could be trained with multiple styles such that the final stylized result fused the content of one test image, the coarse texture distortion of one style image, and the fine brushwork of another style image. Here we gave an example in Fig. 7 where the model was trained with two different styles.

**Processing Speed and Memory Use** We compared quantitatively the speed and memory usage of our multimodal transfer network (MT Net) with other singular transfer networks (Here we used Johnson Net, the network by Johnson *et al.*). We also constructed a deep singular transfer network for comparison (called DS Net), which had the same structure as the MT Net. We took the average of the test time for 1,000 generations (excluding model loading time).

Network	Test Time	Memory Usage
MT Net	0.54s	3100 MB
Johnson Net	0.42s	2400 MB
DS Net	0.63s	6700 MB

Table 1: Comparison results of speed and memory use when applied on  $1024 \times 1024$  images.

As shown in Table 1, although MT Net was more than twice deeper than Johnson Net, its speed and memory usage were close to those of Johnson Net (0.54s vs 0.42s, 3100 MB vs 2400 MB) when generating high-resolution images, which benefited from the hierarchical transfer procedure where most computation was done on low resolutions. The singular transfer network DS Net performed the worst even with same number of parameters. Therefore, multimodal transfer is suitable for real-world applications that usually required high image resolution, because it was able to generate results more similar to the desired artistic styles with a small cost.

## 5. Conclusion

Here we present a hierarchical training scheme (*multimodal transfer*) for fast style transfer to learn artistic style cues at multiple scales, including color, coarse texture structure and fine, exquisite brushwork. The scheme solves the texture scale mismatch issue and generates much more visually appealing stylized results on high-resolution images.

In the future, we plan to investigate other losses that can better capture the artistic style at different scales. We also want to explore alternate loss networks that costs less memory to extend our scheme onto much larger images.



## References

- [1] N. Ashikhmin. Fast texture transfer. *IEEE Computer Graphics and Applications*, 23(4):38–43, 2003.
- [2] E. L. Denton, S. Chintala, R. Fergus, et al. Deep generative image models using a laplacian pyramid of adversarial networks. In *Advances in neural information processing systems*, pages 1486–1494, 2015.
- [3] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1538–1546, 2015.
- [4] A. A. Efros and W. T. Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 341–346. ACM, 2001.
- [5] A. A. Efros and T. K. Leung. Texture synthesis by non-parametric sampling. In *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*, volume 2, pages 1033–1038. IEEE, 1999.
- [6] L. Gatys, A. S. Ecker, and M. Bethge. Texture synthesis using convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 262–270, 2015.
- [7] L. A. Gatys, M. Bethge, A. Hertzmann, and E. Shechtman. Preserving color in neural artistic style transfer. *arXiv preprint arXiv:1606.05897*, 2016.
- [8] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2414–2423, 2016.
- [9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [11] A. Hertzmann, C. E. Jacobs, N. Oliver, B. Curless, and D. H. Salesin. Image analogies. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 327–340. ACM, 2001.
- [12] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [13] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer, 2016.
- [14] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [16] V. Kwatra, A. Schödl, I. Essa, G. Turk, and A. Bobick. Graphcut textures: image and video synthesis using graph cuts. In *ACM Transactions on Graphics (ToG)*, volume 22, pages 277–286. ACM, 2003.
- [17] H. Lee, S. Seo, S. Ryoo, and K. Yoon. Directional texture transfer. In *Proceedings of the 8th International Symposium on Non-Photorealistic Animation and Rendering*, pages 43–48. ACM, 2010.
- [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014.
- [19] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *2015 IEEE conference on computer vision and pattern recognition (CVPR)*, pages 5188–5196. IEEE, 2015.
- [20] A. Nguyen, J. Yosinski, and J. Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 427–436. IEEE, 2015.
- [21] A. Odena, V. Dumoulin, and C. Olah. Deconvolution and checkerboard artifacts. <http://distill.pub/2016/deconv-checkerboard/>, 2016.
- [22] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [23] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [24] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [25] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *Int. Conf. on Machine Learning (ICML)*, 2016.
- [26] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [27] B. A. Wandell. *Foundations of vision*. Sinauer Associates, 1995.
- [28] L.-Y. Wei and M. Levoy. Fast texture synthesis using tree-structured vector quantization. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 479–488. ACM Press/Addison-Wesley Publishing Co., 2000.
- [29] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015.
- [30] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833. Springer, 2014.