# A Flexible Convolutional Solver with Application to Photorealistic Style Transfer

Gilles Puy and Patrick Pérez
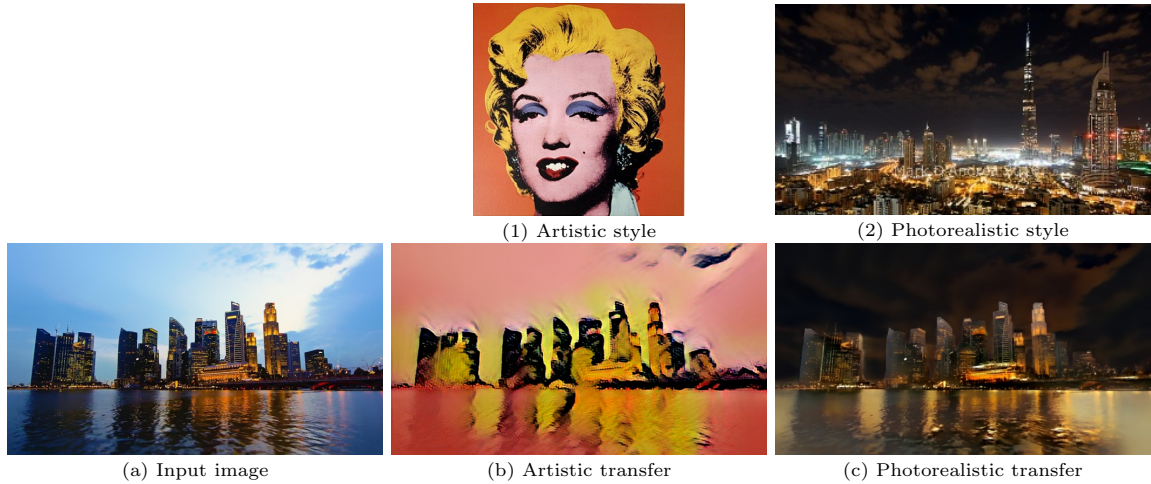
(1) Artistic style

(2) Photorealistic style

(a) Input image

(b) Artistic transfer

(c) Photorealistic transfer

Figure 1. *Given an input image (a), the proposed convolutional neural network yields fast artistic stylisation results,* e.g., *(b) for style (1), and can be modified at runtime – hence without retraining – to yield photorealistic stylisation results,* e.g., *(c) for style (2).*

**Abstract.** We propose a new flexible deep convolutional neural network (convnet) to perform fast visual style transfer. In contrast to existing convnets that address the same task, our architecture derives directly from the structure of the gradient descent originally used to solve the style transfer problem [1]. Like existing convnets, ours approximately solves the original problem much faster than the gradient descent. However, our network is uniquely flexible by design: it can be manipulated at runtime to enforce new constraints on the final solution. In particular, we show how to modify it to obtain a photorealistic result with no retraining. We study the modifications made by [2] to the original cost function of [1] to achieve photorealistic style transfer. These modifications affect directly the gradient descent and can be reported on-the-fly in our network. These modifications are possible as the proposed architecture stems from unrolling the gradient descent.

## 1. Introduction

Transferring the style of one image onto another one, while keeping the content of the second, is a difficult task but for which visually impressive results have recently been obtained thanks to deep neural networks.

In [1], the authors adapt their texture synthesis method [3] to perform style transfer. The stylised image is obtained by minimisation of a loss function built using feature maps obtained from a pre-trained deep neural network. While yielding good results, this method is nevertheless slow as the stylised image is a solution of a long iterative optimisation problem. This issue has however been overcome by researchers who have proposed to train deep networks that estimate rapidly a solution of the style transfer problem. The first architectures were trained to transform any natural image, given as input to the network, to a fixed style [4, 5]. These networks thus need to be retrained for every new style. This drawback was partly addressed by works that showed that the style given to an image was

Gilles Puy and Patrick Pérez are with Technicolor, 975 Avenue des Champs Blancs, 35576 Cesson-Sévigné, France.

controlled by few filters [6] or parameters in the network [7]. It thus becomes possible to do a training for several styles jointly. At runtime, an image can then be transformed to any of the preselected styles by choosing the corresponding set of filters or parameters. Finally, recent works show that it is possible to train a deep neural network that takes as inputs any pair of style and content images and produces a stylized image, even for style images not viewed at training time [8, 9, 10]. Furthermore, it was also showed that it is possible to control the level of stylisation or to mix different styles by doing a linear interpolation between the parameters of the deep network that control the corresponding styles.

While the systems described so far provide impressive results when the style is a painting, the results are less impressive when the style is a natural image. Another line of research concerns the adaptation of the previous techniques to perform photorealistic style transfer. In [2], the authors adapt the loss of [1] with an additional term that favours photorealistic results. This additional term ensures a better preservation of the geometry of the content image. As in [1], this new method is iterative and slow. Up until relatively recently, no method providing a fast solution of this photorealistic style transfer problem existed. A technical report proposing a fast method has however been made public [11] while writing this report. We discuss this alternative method in Section 6.

**Contributions.** First, we propose a new network architecture to perform fast style transfer. The architecture is directly derived from the structure of the gradient descent algorithm used in [1] to solve the style transfer problem. We confirm using experiments that this network can provide fast artistic style transfer results. Second, we show that the network's architecture is flexible enough to be modified at runtime – thus *with no retraining* – to incorporate modification of the initial training loss function. Third, we propose a slight modification of the loss function proposed by [2] to perform photorealistic style transfer and modify our trained network accordingly. Fourth, we show that we achieve results similar to those obtain by [2] with a better preservation of fine details. Finally, we show that our network can be used to perform fast videorealistic style transfer.

## 2. Related Work

**Visual style transfer.** Example-based image editing, whereby an input image is modified so as to acquire some of the characteristics of a reference picture, has a long history, *e.g.* [12, 13, 14] for textures and colors. Recently, the specific problem of example-based stylization of a photograph, with the reference being typically a painting, has been revisited with great success by Gatys *et al.* [1], using pre-trained deep convolutional features as the key ingredient. Through iterative optimization, a new image is created *ex-nihilo* such that it retains the main structure (the "content") of the original photo while adopting the look and feel (colors, textures, paint strokes, the "style" in short) of the reference. This two-fold goal is encapsulated in a cost function over deep features of the three images. One major limitation of the approach is the cost of the optimization, with each iteration requiring forward-backward computations through the deep convnet that extracts the feature maps. Several authors, starting with [4][5], have showed that another feed-forward deep convnet can be designed and trained without supervision, for one or several given target styles, to perform the task much faster with similar quality: given a new content, this image is directly transformed into the stylized content without resorting to costly optimization. Our approach is also of this type, *i.e.*, unsupervised learning of a "neural solver" for the original optimization problem, but with a very different take on the design of the architecture. Namely, we rely on the powerful idea of unrolling an iterative process to turn it into a trainable neural network. One key advantage of our approach is that important modifications of the original cost function can be mimicked in our architecture with no need for retraining. We demonstrate this in the case of photorealistic transfer [2]. Fast methods in [4, 5, 6, 7] might struggle with this extension of Gatys *et al.*'s cost function, in particular due to the matting Laplacian that plays a crucial role, and would require retraining in any case. In contrast, a simple runtime restructuring of our architecture can accommodate these major changes. Note that alternative optimisation-based techniques to [2], such as [15], exist for photorealistic color/style transfer. We use our flexible framework to rapidly solve [2] but the principles detailled here could be used to solve other
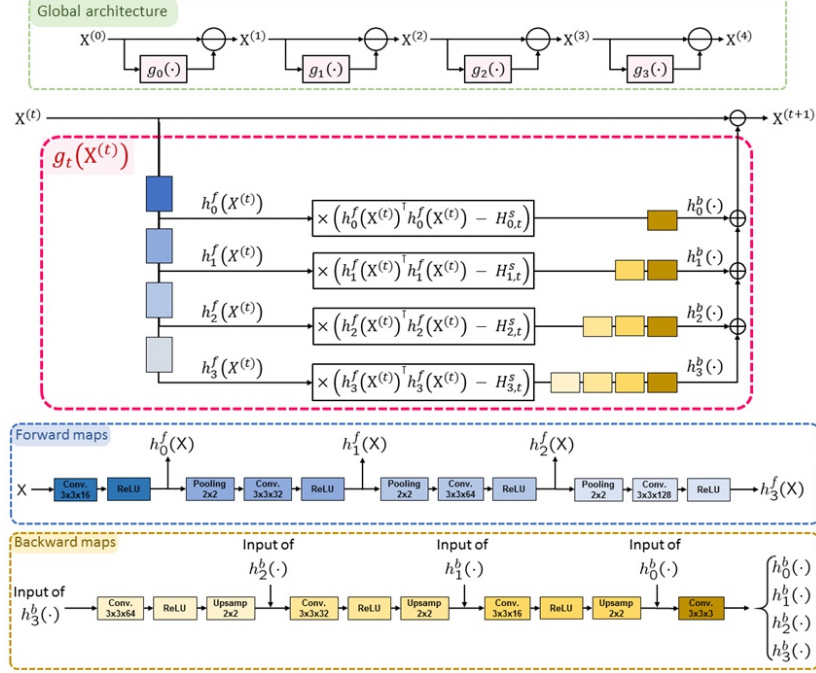
Figure 2. *First (top): Global structure of the trained convnet for fast style transfer. It is a stack of four networks with identical structures, each one transforming an input image into a new one through one step of learned descent. Second: The structure of the descent maps $\{g_t\}$. The coloured blocks correspond to subnetworks with same colors in the forward and backward maps $\{h_\ell^f\}$ and $\{h_\ell^b\}$. Third and fourth: Structure of the forward and backward maps for the computation of $\{g_t\}$.*

optimisation problems. Finally, we discuss in Section 6 the fast photorealistic style transfer method [11] which was made public while writing this report.

**Neural nets from algorithms unrolling.** A wide range of classic iterative solvers, *e.g.*, to solve optimization problems or partial differential equations, amount to the repeated application of linear and component-wise non-linear maps. Unfolding such an algorithm over a fixed number of iterations allows one to see it as several layers of a neural network with shared, predefined weights. This can be taken as the starting point for an actual, trainable neural net. This powerful idea was introduced by Gregor and LeCun [16] in the context of sparse coding: the popular iterative shrinkage thresholding algorithm (ISTA) [17] gives rise of a learnable network (LISTA) for fast approximate sparse coding. Each layer has independent trainable weights and the whole network is trained under the supervision of ISTA solutions. In the related context of linear inverse problems, including compressed sensing [18, 19, 20, 21] and image restoration [22, 23, 24], several works have followed this unrolling approach. In some cases, weight tying across layers is maintained as in the unrolled algorithm [22]. In some other cases, the architecture departs notably from the unrolled algorithm, *e.g.*, through using multiple convolutional filters per layer instead of only one as suggested by the unrolled iteration [24]. Let us stress out that, by essence of such inverse problems, full supervision is possible: input-output samples of the transform to be inverted are indeed readily obtained, *e.g.*, by sensing or artificially degrading natural images. By contrast, our approach is fully unsupervised, as other fast style transfer alternatives. In the absence of known outputs for corresponding input training samples, the training loss is defined as the original cost function.

**Unsupervised neural solvers.** Training without supervision a (fast) feed-forward network to solve approximately a complex optimization problem is a challenging problem, with important applications. Using the cost function of interest as training loss seems straightforward, but the capacity of the designed architecture and the amount of actual guidance provided by the loss must be sufficient

for such an unsupervised training to succeed. Besides works in fast style transfer, other recent works have demonstrated successful instances of this paradigm in other contexts. In [25], efficient fluid simulation is addressed. Within each time step of a classic solver of the incompressible Euler equation, the costly solving of the Poisson equation on the pressure field is replaced by a trained 3D convnet. In this case the loss amounts to the divergence square norm. In a very different context, [26] learn a deep neural solver for a complex inverse rendering problem: estimating the shape, expression and reflectance of a face in a single image, by minimization of the pixel-wise discrepancy between rendered and real faces. In both works, trained solvers provide good quality approximation of the solutions of the initial minimisation problem, but with massive acceleration compared to iterative solvers, and with no need for initialization in the latter work. In the present paper, we also propose to learn an unsupervised neural solver for a complex optimization problem, but following the unrolling principle explained above. As a consequence, the acceleration, already obtained by others on the same problem, is not our main contribution. More importantly, the unrolling approach gives the possibility to adapt at runtime our network in order to accommodate important changes to the original cost function.

## 3. Network Architecture

### 3.1. Style Transfer by Gradient Descent

Style transfer consists in transforming an image $X_c \in \mathbb{R}^{n \times 3}$ – of $n$ pixels and 3 color channels[1] – to give it the "style" of another image $X_s \in \mathbb{R}^{n' \times 3}$ while preserving the content in $X_c$. The style transfer method of Gatys *et al.* [1] consists in solving a minimisation problem involving deep features obtained from a pre-trained convnet, typically the VGG-19 network [27], to obtained a stylised image $X \in \mathbb{R}^{n \times 3}$. We denote by $F_\ell, C_\ell \in \mathbb{R}^{n_\ell \times c_\ell}$ and $S_\ell \in \mathbb{R}^{n'_\ell \times c_\ell}$ the $\ell^{\text{th}}$ feature maps – in the VGG-19 network – obtained from the images $X, X_c$ and $X_s$, respectively. The minimised loss takes the form

$$(1) \qquad\qquad \mathcal{L}(X) = \lambda_c \, \mathcal{L}_c(X, X_c) + \lambda_s \, \mathcal{L}_s(X, X_s),$$

where the content loss $\mathcal{L}_c(X) = |\mathcal{I}_c|^{-1} \sum_{\ell \in \mathcal{I}_c} \|F_\ell - C_\ell\|_{\text{F}}^2 / (n_\ell \, c_\ell)$ ensures transfer of the content of $X_c$ in the final image ($\|\cdot\|_{\text{F}}$ stands for matrix Frobenius norm), while the style loss

$$(2) \qquad\qquad \mathcal{L}_s(X) = \frac{1}{|\mathcal{I}_s|} \sum_{\ell \in \mathcal{I}_s} \frac{1}{c_\ell^2} \left\| \frac{1}{n_\ell} F_\ell^{\mathsf{T}} F_\ell - \frac{1}{n'_\ell} S_\ell^{\mathsf{T}} S_\ell \right\|_{\text{F}}^2$$

ensures transfer of the style of $X_s$ to the final image, and $\lambda_c, \lambda_s > 0$. As in [1], the layers used to encode the style and the content are $\mathcal{I}_s = \{\text{conv}_{1\_1}, \text{conv}_{2\_1}, \text{conv}_{3\_1}, \text{conv}_{4\_1}, \text{conv}_{5\_1}\}$ and $\mathcal{I}_c = \{\text{conv}_{4\_2}\}$, respectively.

A stylised image $X^*$ is usually obtained by starting from a random image $X^{(0)}$ and by progressively updating it to minimise the loss $\mathcal{L}$ by gradient descent:

$$(3) \qquad\qquad X^{(t+1)} = X^{(t)} - \mu \left[ \lambda_c \, \nabla \mathcal{L}_c(X^{(t)}) + \lambda_s \, \nabla \mathcal{L}_s(X^{(t)}) \right],$$

where $\mu > 0$ is the stepsize. While achieving impressive results, this method is nevertheless slow. To overcome this issue, researchers have designed deep convnets trained to solve (1) approximately at a much lower computational cost [4, 5, 10], that is a neural map $f(\cdot)$ such that $\mathcal{L}\big(f(X^{(0)})\big) \approx \min_X \mathcal{L}(X)$ provided that $X_c$ is a natural image and that the network input is suitable, *e.g.*, $X^{(0)} = X_c$ in [4, 10] or a combination of $X_c$ and additive white noise in [5].

### 3.2. Learned Gradient Descent for Style Transfer

We propose a new deep neural network to achieve fast style transfer. The architecture of this network is directly inspired by the structure of the gradient descent algorithm, following the idea of unrolling [16]. This structure gives a better interpretability to the role of each filter as they correspond to

---

[1]Throughout, the two spatial dimensions of images and features maps are flattened into a single vector dimension for notational convenience.

specific parts of the original gradient descent algorithm. This property permits us to augment, at runtime, the network with new desirable features, for instance to perform photorealistic style transfer.

**Global architecture.** The proposed convnet follows the update rule of the gradient descent algorithm (3) where the actual gradient is replaced by a learned descent direction:[2]

$$(4) \qquad\qquad \mathsf{X}^{(t+1)} = \mathsf{X}^{(t)} - g_t\left(\mathsf{X}^{(t)}\right), \quad t = 0, 1, 2, 3.$$

The global architecture of the network is illustrated in Fig. 2. Note that the proposed network can be viewed as a residual network [28]. We restrict the computation to 4 iterations to maintain the computational cost low and constant. By comparing (3) to (4), we see that $g_t\left(\mathsf{X}^{(t)}\right)$ plays the role of the gradient. We explain next how $g_t\left(\mathsf{X}^{(t)}\right)$ is computed.

**Computing the descent direction.** The gradient in (3) is made of a content term and a style term. To simplify the network and reduce computational costs, we design $g_t$ by mimicking the computational architecture of $\nabla \mathcal{L}_s$ only. However, as in [4], we initialise $\mathsf{X}^{(0)}$ with the RGB image[3] $\mathsf{X}_c$ – after scaling the pixel values in $[0, 1]$ – instead of the classic random initialization of the gradient descent, and we use the complete cost $\mathcal{L}$ as training loss. This allows us to overcome the lack of $\nabla \mathcal{L}_c$ in the network and keep the content of $\mathsf{X}_c$ in the final image $\mathsf{X}^{(4)}$.

Let us now concentrate on $\nabla \mathcal{L}_s$. We recall that

$$(5) \qquad\qquad \frac{\partial \mathcal{L}_s}{\partial \mathsf{F}_\ell} \propto \mathsf{F}_\ell \cdot \left[ \frac{1}{n_\ell} \mathsf{F}_\ell^{\mathsf{T}} \mathsf{F}_\ell - \frac{1}{n_\ell'} \mathsf{S}_\ell^{\mathsf{T}} \mathsf{S}_\ell \right].$$

The role of the difference of Gram matrices in (5) is to correct the feature maps $\mathsf{F}_\ell$ so that, after backpropagation, the image $\mathsf{X}^{(t)}$ gets closer to the style of $\mathsf{X}_s$. The filter resulting from the difference of Gram matrices is called *style correction filter* hereafter. The gradient of $\mathcal{L}_s$ is obtained by

  (i) computing the feature maps $\mathsf{F}_\ell$ (forward maps);
  (ii) computing the partial derivatives (5);
  (iii) back-propagating these partial derivatives to the input (backward maps).

We construct $g_t$ by mimicking these 3 steps with a simpler convnet than VGG-19.

[*Transforming Steps (i) and (iii)*] The structure of the simplified forward and backward maps are presented in Fig. 2. The new feature maps playing the role of $\{\mathsf{F}_\ell\}$ are $\{h_\ell^f(\mathsf{X})\}$. As in the original style loss, the spatial dimension between two consecutive layers is divided by 2 in both directions and that the number of feature channels is multiplied by 2. Let us highlight that we use only 4 layers instead of 5 as originally in $\mathcal{L}_s$. The reason is only computational. We noticed that this was enough to yield satisfactory results. The backward maps $\{h_\ell^b(\cdot)\}$ are obtained from the forward maps by symmetry. Let us highlight that the last layer of the backward maps does not involve any ReLU to allow negative descent direction. The pooling layers use average pooling and the upsampling layers use bilinear upsampling. All convolutions are computed using reflection padding. All the filters in these maps are learned and are *identical* for all styles.

[*Transforming Step (ii)*] The computation of the partial derivatives (5) is directly transformed to

$$(6) \qquad\qquad h_\ell^f(\mathsf{X}^{(t)}) \cdot \left[ \frac{1}{n_\ell} h_\ell^f(\mathsf{X}^{(t)})^{\mathsf{T}} h_\ell^f(\mathsf{X}^{(t)}) - \mathsf{H}_{\ell,t}^s, \right]$$

where the learned matrix $\mathsf{H}_{\ell,t}^s$ plays the role of the Gram matrix $\mathsf{S}_\ell^{\mathsf{T}} \mathsf{S}_\ell$ that encodes the style. Note however that the matrix $\mathsf{H}_{\ell,t}^s$ is different at each iteration $t$, unlike in (5). The matrices $\{\mathsf{H}_{\ell,t}^s\}$ are specific to each style. For simplicity, we do not force $\mathsf{H}_{\ell,t}^s$ to be symmetric and positive semi-definite during training.

The complete structure of $g_t$ is presented in Fig. 2. Let us highlight that the maps $\{h_\ell^f\}$ and $\{h_\ell^b\}$ are independent of $t$; they are shared across the four iteration networks. Note that the matrix multiplication with (6) can be implemented as a convolution with a filter of spatial size $1 \times 1$, so that the whole network is indeed convolutional. We also notice that this filter is instance dependent: its weights depend on $h_\ell^f(\mathsf{X}^{(t)})$. The forward and backward maps $\{h_\ell^f\}$ and $\{h_\ell^b\}$ contain $194\,755$ weights

---

[2]There is no mathematical guarantee that $g_t\left(\mathsf{X}^{(t)}\right)$ is indeed a descent direction.

[3]At training time, a zero-mean uniform noise of amplitude drawn uniformly in $[0, 0.1]$ is added to $\mathsf{X}^{(0)}$.

to learn. At each iteration $t$, the style matrices contain $21\,760$ trainable weights. The total number of parameters is thus $194\,755 + 21\,760 \cdot 4 = 281\,795$.

**Training.** Following [4], we augment $\mathcal{L}$ defined in (1) with a Total Variation regularisation term to promote piecewise smoothness of the transformed images:

$$(7) \qquad \lambda_c \, \mathcal{L}_c(\mathsf{X}^{(4)}) + \lambda_s \, \mathcal{L}_s(\mathsf{X}^{(4)}) + \lambda_{\mathrm{TV}} \, \mathrm{TV}(\mathsf{X}^{(4)}),$$

where the values in the RGB image $\mathsf{X}^{(4)}$ are clipped between 0 and 1. All filters in $\{h_\ell^f\}$, $\{h_\ell^b\}$ are initialised using Xavier initialisation [29] and the matrices $\mathsf{H}_{\ell,t}^s$ are initialised at zero. The filters are trained using Adam [30] with a constant stepsize of $10^{-5}$ and a batchsize of 1. The network is trained for 17 epochs over the 2014 MS-COCO training dataset [31].[4] The training images are cropped to the largest possible square, resized to $320 \times 320$, and used as content images for training. We use 69 style images in total, including 18 images of paintings. The remaining 51 images are photographs used by the authors of [2] and for which the masks to segment the images into different semantic regions are available.[5] Each semantic region of a photograph defines a style. We thus consider 153 styles in total, coming from either paintings or photographs. The style term $\mathcal{L}_s$ is therefore modified with semantic masks that apply on the style images during training. Let $\mathsf{M}^s \in \mathbb{R}^{n' \times n'}$ be the binary diagonal matrix that encodes a semantic mask on the style image $\mathsf{X}_s$, the style term becomes

$$(8) \qquad \mathcal{L}_s(\mathsf{X}) = \frac{1}{|\mathcal{I}_s|} \sum_{\ell \in \mathcal{I}_s} \frac{1}{c_\ell^2} \left\| \frac{1}{n_\ell} \mathsf{F}_\ell^{\mathsf{T}} \mathsf{F}_\ell - \frac{1}{\mathrm{Tr}(\mathsf{M}_\ell^s)} (\mathsf{M}_\ell^s \mathsf{S}_\ell)^{\mathsf{T}} (\mathsf{M}_\ell^s \mathsf{S}_\ell) \right\|_{\mathrm{F}}^2$$

where $\mathsf{M}_\ell^s$ is the semantic mask propagated at the $\ell^{\mathrm{th}}$ layer of the VGG-19. Let us remark that in absence of semantic mask, then $\mathsf{M}_\ell^s$ is the identity, $\mathrm{Tr}(\mathsf{M}_\ell^s) = n'_\ell$, and we recover the original style loss. Note that the content images are not segmented during training. We use $\lambda_c = 0.025$, $\lambda_{\mathrm{TV}} = 0.5$, and

$$(9) \qquad \lambda_s = \left[ \frac{1}{|\mathcal{I}_s|} \sum_{\ell \in \mathcal{I}_s} \frac{1}{c_\ell^2} \left\| \frac{1}{\mathrm{Tr}(\mathsf{M}_\ell^s)} (\mathsf{M}_\ell^s \mathsf{S}_\ell)^{\mathsf{T}} (\mathsf{M}_\ell^s \mathsf{S}_\ell) \right\|_{\mathrm{F}}^2 \right]^{-1}.$$

We notice that this choice of $\lambda_s$ allows similar degree of stylization across learned styles. The evolution of the training loss computed on 20 images of the 2014 MS-COCO validation set is presented in Fig. 3. We remark that the loss starts to stabilise after around 14 epochs.

## 4. Runtime Restructuring

In this section, we use properties of our trained network to obtain photorealistic results via modifications made at runtime. We recall what changes need to me made to the original loss function for style transfer to obtain photorealistic results. These changes directly impact the gradient descent algorithm used to solve this problem. As our network has the same structure as this algorithm, it is straightforward to make the same modifications in the network to obtain the desired results with no retraining.

### 4.1. Laplacian Regularisation for Photorealism

Minimizing (1) yields very good results when transferring the style of a painting to a photograph [1]. Unfortunately, the results are less impressive when transferring the style of a photograph to another photograph as the result is often not photorealistic. The authors of [2] propose to solve this issue by favouring transformations from the content image to the final image that are locally affine. In practice, this is done by adding a penalty term to $\mathcal{L}$, which becomes

$$(10) \qquad \mathcal{L}(\mathsf{X}) + \lambda_{\mathrm{L}} \, \mathrm{Tr}(\mathsf{X}^{\mathsf{T}} \mathsf{L} \mathsf{X}),$$

where $\mathsf{L} \in \mathbb{R}^{n \times n}$ is the Matting Laplacian [32] of the content image $\mathsf{X}_c$. Even though our network was trained without this Laplacian term, we can enforce an equivalent constraint at runtime to obtain the desired results.

---

[4]http://cocodataset.org/

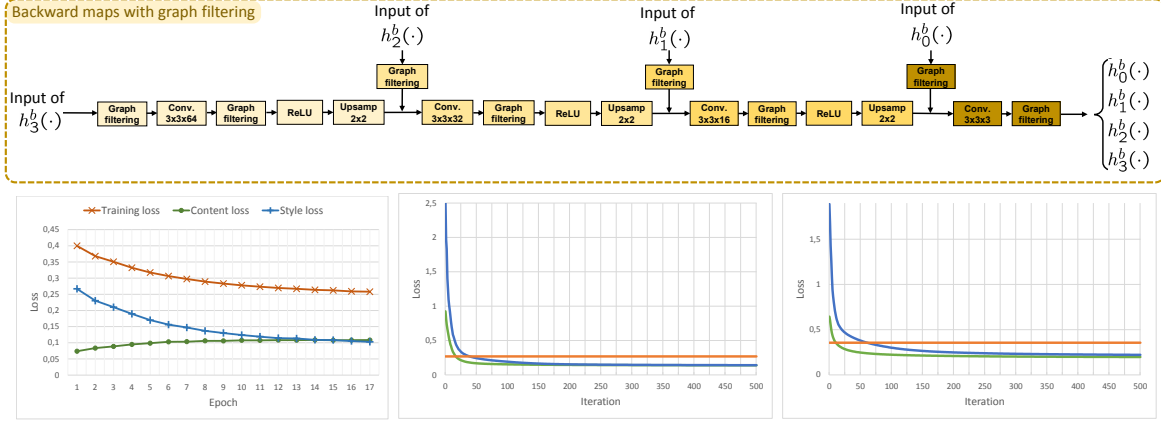[5]https://github.com/luanfujun/deep-photo-styletransfer

Figure 3. *Top: Structure of the backward maps modified at runtime with graph filtering. Bottom left: Evolution of the training loss (orange), and of its summands related to $\mathcal{L}_c$ (green) and $\mathcal{L}_s$ (blue), on a validation set of 20 images as a function of the number of epochs. Bottom, middle and right: Evolution of (7) vs. the number of iterations of L-BFGS when initialised with noise (blue) or $\mathsf{X}_c$ (green). Value of (7) attained with the solution of our trained network (orange). Curves are obtained using an image of* The scream *(middle) or of the* Portrait of Dora Maar *(right) as styles.*

In graph signal processing terms, finding an image $\mathsf{X}$ for which $\mathrm{Tr}(\mathsf{X}^{\mathsf{T}}\mathsf{L}\mathsf{X})$ is small means finding an image that is smooth on the graph $\mathcal{G}$ encoded by the Laplacian matrix $\mathsf{L}$. Let us recall the concept of smoothness on graphs.

**Smooth graph signals.** The graph Fourier transform of any graph signal $\boldsymbol{x} \in \mathbb{R}^n$ is defined as $\hat{\boldsymbol{x}} = \mathsf{U}^{\mathsf{T}}\boldsymbol{x}$, where $\mathsf{U} = (\boldsymbol{u}_1, \ldots, \boldsymbol{u}_n) \in \mathbb{R}^{n \times n}$ is the matrix of eigenvectors of $\mathsf{L}$: $\mathsf{L} = \mathsf{U}\Lambda\mathsf{U}^{\mathsf{T}}$, where $\Lambda = \mathrm{diag}(\lambda_1, \ldots, \lambda_n) \in \mathbb{R}^{n \times n}$ and $\lambda_1 \leqslant \ldots \leqslant \lambda_n$ [33]. A smooth graph signal is a signal whose energy is concentrated on the first few Fourier coefficients (Fourier coefficients at the lowest graph frequency). A $k$-bandlimited graph signal is a signal whose energy is entirely concentrated on the first $k$ Fourier coefficients, and is thus a special type of smooth signal when $k \ll n$. The set of $k$-bandlimited graph signals is denoted $\mathrm{span}(\mathsf{U}_k)$ where $\mathsf{U}_k = (\boldsymbol{u}_1, \ldots, \boldsymbol{u}_k) \in \mathbb{R}^{n \times k}$ is the restriction of $\mathsf{U}$ to its first $k$ columns.

Another fundamental tool in graph signal processing is filtering. One can define spectral graph filtering as follows $\boldsymbol{x} * u = \mathsf{U}\,\mathrm{diag}(u(\lambda_1), \ldots, u(\lambda_n))\,\mathsf{U}^{\mathsf{T}}\boldsymbol{x}$, where $\boldsymbol{x} * u$ is the signal $\boldsymbol{x}$ filtered by $u : \mathbb{R}_+ \to \mathbb{R}$. Projecting a signal $\boldsymbol{x}$ onto the set of $k$-bandlimited graph signal is done by filtering with

$$(11) \qquad i_{\lambda^*}(\lambda) = \left\{ \begin{array}{ll} 1 & \text{if } \lambda \leqslant \lambda^*, \\ 0 & \text{oth,} \end{array} \right. \qquad \text{with } \lambda^* \in [\lambda_k, \lambda_{k+1}).$$

**Constrained Version of** (10). Minimising $\mathrm{Tr}(\mathsf{X}^{\mathsf{T}}\mathsf{L}\mathsf{X})$ as in (10) means finding an image that is smooth on the graph $\mathcal{G}$ encoded by the Laplacian matrix $\mathsf{L}$. Instead of adding the penalisation term $\mathrm{Tr}(\mathsf{X}^{\mathsf{T}}\mathsf{L}\mathsf{X})$ to $\mathcal{L}$, one can also search for an image $\mathsf{X}$ smooth on $\mathcal{G}$ by solving the constrained problem

$$(12) \qquad \min_{\mathsf{X}} \mathcal{L}(\mathsf{X}) \quad \text{subject to} \quad \mathsf{X} \in \mathrm{span}(\mathsf{U}_k).$$

This problem can be solved by projected gradient descent whose iterations satisfy:

$$(13) \qquad \mathsf{X}^{(t+1)} = \left[\mathsf{X}^{(t)} - \mu\,\nabla\mathcal{L}\left(\mathsf{X}^{(t)}\right)\right] * i_{\lambda^*},$$

where the filtering with $i_{\lambda^*}$ projects the iterates onto the space of $k$-bandlimited signals on $\mathcal{G}$. The only difference with (3) is thus the extra projection after each gradient update. It is straightforward to add this constraint in our deep CNN by changing (4) to

$$(14) \qquad \mathsf{X}^{(t+1)} = \mathsf{X}^{(t)} - g_t\left(\mathsf{X}^{(t)}\right) * i_{\lambda^*}.$$

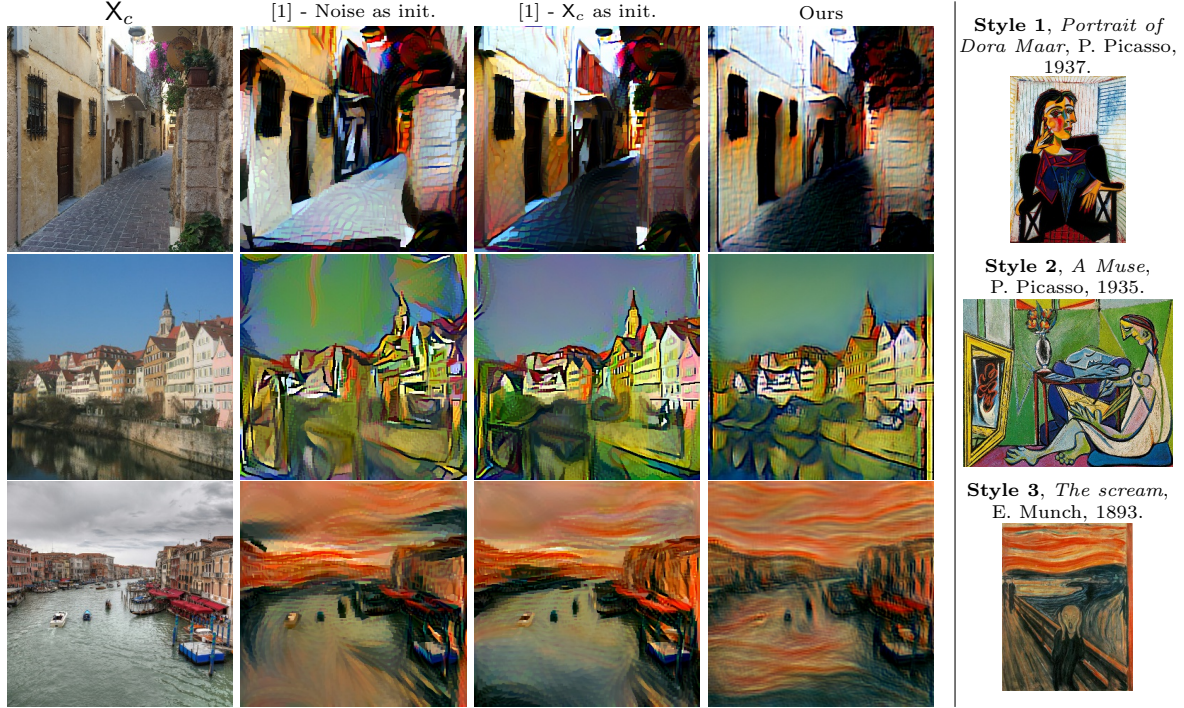|  $\mathsf{X}_c$  | [1] - Noise as init. | [1] - $\mathsf{X}_c$ as init. | Ours | |

Figure 4. *Artistic style transfer results obtained with* [1] *using two different initialisations (second and third columns) and our network (fourth column), for 3 different content images (first column) and corresponding styles (last column).*

As we initialise $\mathsf{X}^{(0)}$ with $\mathsf{X}_c$, and because we can safely assume that $\mathsf{X}_c \approx \mathsf{X}_c * i_{\lambda^*}$ (otherwise the smoothness constraint would not be a good photorealistic prior), it is enough to filter $g_t(\mathsf{X}^{(t)})$ at every iteration to ensure that all $\mathsf{X}^{(t)}$ are $k$-bandlimited.[6] Adding this graph filtering step at runtime is sufficient to transform a network trained to solve (10) to a new network solving (12).

**Fast graph filtering.** Filtering with $i_{\lambda^*}$ requires the computation of the matrix $\mathsf{U}$. To avoid this intractable computation, we use a polynomial approximation $\tilde{i}_{\lambda^*}$ of $i_{\lambda^*}$ to perform filtering. Polynomial filters are computationally efficient as they only require matrix vector multiplication with the sparse Laplacian matrix $\mathsf{L}$. We let the reader refer to, *e.g.*, [34] for more explanations on this subject. In practice, we use Jackson-Chebychev polynomials [35] of order 5 computed to approximate $i_{\lambda^*}$ with $\lambda^* = 0.2\,\lambda_n$.[7] Note that polynomial graph filters are also used to build efficient graph CNNs [36, 37].

We noticed that to achieve more pleasing results visually, it is better to filter $g_t(\mathsf{X}^{(t)})$ at different resolutions. It permits to preserve better large scale and low scale structures. We take advantage of the multiscale structure of the trained network to compute this multiscale filtering efficiently. Each feature map after computing (6) and in the four levels of the backward pass is graph filtered individually (channelwise). This graph filtering is placed after each convolutional layer and before the ReLU, if present, as presented in Fig. 3. The four used Laplacian matrices are computed on the initial content image $\mathsf{X}_c$ downsampled at the corresponding resolutions of the feature maps in $\{h_\ell^b\}$. Note that, for fair comparison, we post-process $\mathsf{X}^{(4)}$ using a guided filter as done in [2, 11]. The effect of this post-processing is studied in the Appendix. It permits to remove artefacts that sometimes appear near edges, to sharpen some edges, but also to correct compression artefacts when present.
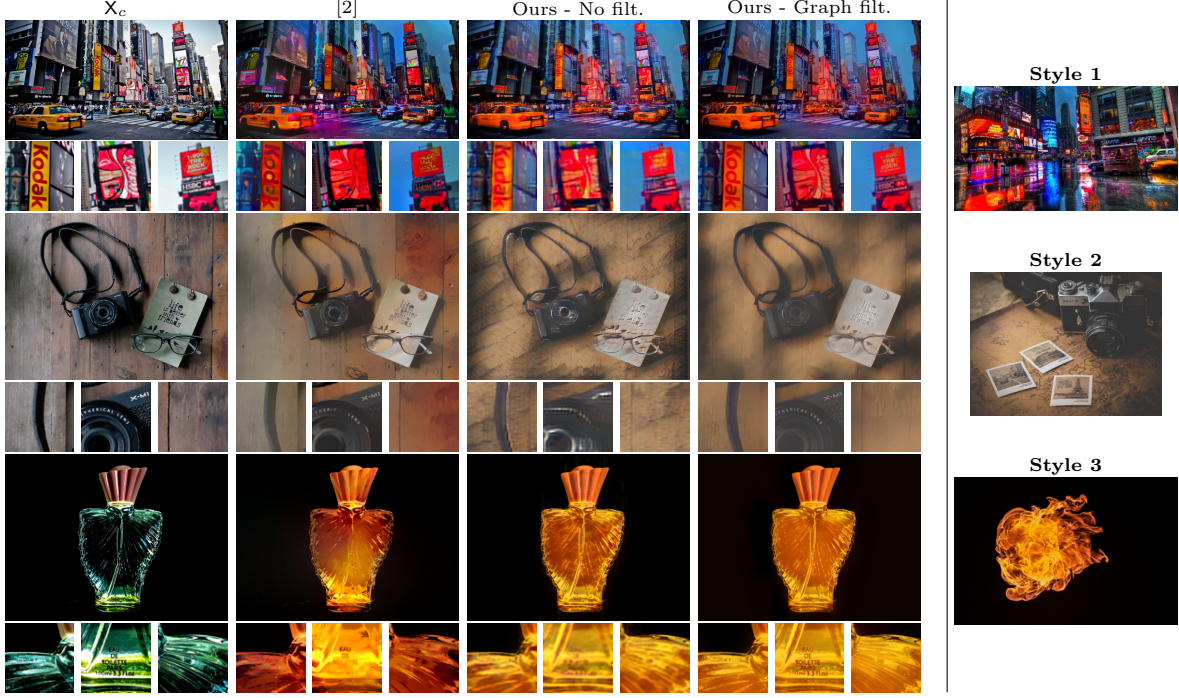
Figure 5. *Photorealistic style transfer results obtained by [2] (second column) and our network without (third column) and with graph filtering (fourth column), for 3 different content images (first column) and corresponding style (last column). Zoom in different regions are presented below each result for better visualisation of fine details.*

## 4.2. Masking at Runtime

Another fundamental part of the method in [2] is to match parts of the content and style images that have the same semantic, *e.g.*, sky to sky, buildings to buildings, roads to roads, etc. This matching can be achieved by using masks in the style loss $\mathcal{L}_s$. Let $\mathsf{M}^c \in \mathbb{R}^{n \times n}$ and $\mathsf{M}^s \in \mathbb{R}^{n' \times n'}$ denote the masks applying on $\mathsf{X}_c$ and $\mathsf{X}_s$ to match similar semantic regions. The masked style loss satisfies

$$\mathcal{L}_s(\mathsf{X}) = \frac{1}{|\mathcal{I}_s|} \sum_{\ell \in \mathcal{I}_s} \frac{1}{c_\ell^2} \left\| \frac{1}{\mathrm{Tr}(\mathsf{M}_\ell^c)} (\mathsf{M}_\ell^c \mathsf{F}_\ell)^\mathsf{T} (\mathsf{M}_\ell^c \mathsf{F}_\ell) - \frac{1}{\mathrm{Tr}(\mathsf{M}_\ell^s)} (\mathsf{M}_\ell^s \mathsf{S}_\ell)^\mathsf{T} (\mathsf{M}_\ell^s \mathsf{S}_\ell) \right\|_\mathrm{F}^2 ,$$

where $\mathsf{M}_\ell^c$ and $\mathsf{M}_\ell^s$ are the masks propagated at the $\ell^{\mathrm{th}}$ layer of the VGG-19. The partial derivative (5) becomes

$$(15) \qquad \frac{\partial \mathcal{L}_s}{\partial \mathsf{F}_\ell} \propto (\mathsf{M}_\ell^c \mathsf{F}_\ell) \cdot \left[ \frac{1}{\mathrm{Tr}(\mathsf{M}_\ell^c)} (\mathsf{M}_\ell^c \mathsf{F}_\ell)^\mathsf{T} (\mathsf{M}_\ell^c \mathsf{F}_\ell) - \frac{1}{\mathrm{Tr}(\mathsf{M}_\ell^s)} (\mathsf{M}_\ell^s \mathsf{S}_\ell)^\mathsf{T} (\mathsf{M}_\ell^c \mathsf{S}_\ell) \right] .$$

The *style correction filter* is thus computed by taking into account the semantic regions of interest only. By copying these modifications in (6), we obtain

$$(16) \qquad h_\ell^f(\mathsf{X}^{(t)}) \cdot \left[ \frac{1}{\mathrm{Tr}(\mathsf{M}_\ell^c)} \left( \mathsf{M}_\ell^c h_\ell^f(\mathsf{X}^{(t)}) \right)^\mathsf{T} \left( \mathsf{M}_\ell^c h_\ell^f(\mathsf{X}^{(t)}) \right) - \mathsf{H}_{\ell,t}^s \right] .$$

Note that $\mathsf{H}_{\ell,t}^s$ encodes the style of a single semantic region by construction of the training procedure, hence the absence of modification of this matrix at runtime. Let us highlight that we did not multiply $h_\ell^f(\mathsf{X}^{(t)})$ with $\mathsf{M}_\ell^c$ on the left-hand side of (16). We noticed that this masking creates artefacts at the boundary of the mask. Instead, we mask $\mathsf{X}^{(4)}$ to keep the effect of the style transfer on the semantic

---

[6]Remark that the space of $k$-bandlimited signals is linear.

[7]Fixing $k$ is also possible but requires additional computations to estimate $\lambda^*$.
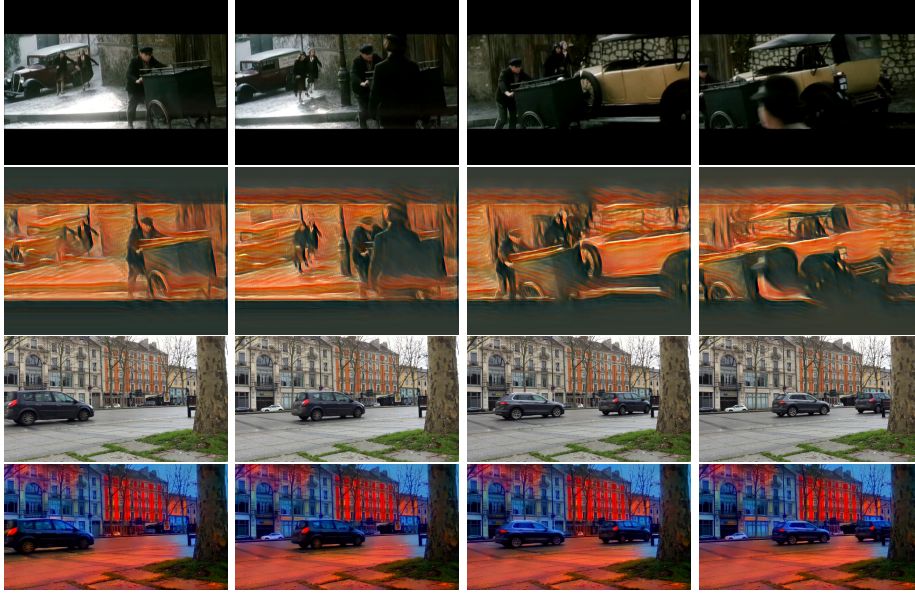
Figure 6. *Style transfer results on two video sequences (first and third rows) for artistic style* 1 *of Fig. 4 (second row) and photorealistic style* 1 *of Fig. 5 (last row).*

region of interest only. Masking $h_\ell^f(\mathsf{X}^{(t)})$ to compute the *style correction filter* is however necessary as, otherwise, one takes into account the whole image in the amount of correction instead of just the semantic region to correct. Let us highlight that to recombine the different stylised semantic regions, we do not apply binary masks on $\mathsf{X}^{(4)}$ but alpha mattes computed using the method of *Levin et al.* [32] and user-defined strokes.

## 5. Style Transfer Experiments

**Fast artistic style transfer.** We present in Fig. 4 artistic style transfer results. We first show results obtained using the method of Gatys *et al.* [1] by minimising (7) using the L-BFGS algorithm. We use random noise as initialisation of this algorithm, as originally done in [1]. We also initialise the L-BFGS algorithm with $\mathsf{X}_c$ as our network is trained to minimise (7) starting form this initial image. Finally, we also present in the same figure the corresponding results obtained using our trained network. We notice that we obtain images qualitatively similar to those obtained by [1]. As a quantitative result, we show in Fig. 3 the evolution of (7) as a function of the number of iterations of the L-BFGS algorithm and the value of (7) reached by our network. The curves show values averaged over 5 test images using for 2 different styles. Our method achieves a loss comparable to around 40 iterations of L-BFGS when initialized with noise. Similar results are obtained in [4]. This shows that our network achieves results close to state-of-the-art methods for fast artistic style transfer.

**Fast photorealistic style transfer.** We present in Fig. 5 style transfer results for styles that are photographs. The results are presented both with and without multiscale graph filtering at runtime. For comparison, we also present the results[8] obtained by Luan *et al.* [2]. First, we notice that graph filtering permits to correct geometric distortions and improve the photorealism. Second, we claim that we achieve results of similar overall quality as Luan *et al.* [2]. We notice however that we preserve better fine details. This is more visible in close-ups below each image, especially in those containing text.

**Video style transfer.** We present in Fig. 6 artistic and photorealistic style transfer results on four frames of two video sequences. Even though each frame is processed individually, we remark that the results are relatively consistent temporally. The processing time was 0.10 s and 8.4 s per frame per

---

[8]Downloaded at `https://github.com/luanfujun/deep-photo-styletransfer`

| $\mathsf{X}_c$ | [2] | [11] | Ours | $\mathsf{X}_s$ |
|---|---|---|---|---|



Figure 7. *Style transfer results obtained by the optimisation method of* [2]*, the fast method of* [11] *and us.*

style for the artistic – on frames of $640 \times 480$ pixels – and photorealistic – on frames of $640 \times 360$ pixels – style transfers, respectively.[9] The increase in computational time is due to graph filtering. Even though graph filtering is computationally efficient, we still have to filter all the feature maps – 428 in total – at all the scales of the backward maps. One could reduce the computational time, by, *e.g.*, filtering only $g_t(\cdot)$. We noticed however that this usually requires fine tuning of the graph filtering parameters for each pair of content and style images to achieve satisfactory results.

## 6. Discussion and Conclusion

A technical report proposing a fast method for photorealistic style transfer was made public [11] while writing this report. The authors adapt the deep network of [8] used for fast artistic style transfer: the deep decoders trained to invert different levels of the VGG-19 in [8] now use unpooling layers instead of upsampling layers. A non-photorealistic stylised image is then obtained exactly as in [8] but using these new decoders. To achieve a photorealistic result, this intermediate image is post-processed by solving a quadratic problem involving a graph regularisation term with the matting Laplacian of [32]. Our method is fundamentally different. First, the network architecture is derived by unrolling the gradient descent algorithm, which is a novel approach for style transfer. Second, photorealism is obtained by modifying the network at runtime, not by post-processing. Third, our network contains only 281 795 parameters once a style is fixed, while the networks in [8] contains more than 7 million parameters, including about 3.5 million pre-trained parameters of the VGG-19. We remark however that the method of [11] does not need retraining for each new style whereas ours does. We present in Fig. 7 results obtained with this method and ours. We notice that we achieve similar qualities to [11], whereas the color palette of the result obtained by [2] still contains colors of the content image. Finally, we also remark that no results on videos are presented in [11].

We proposed a new deep neural network architecture for fast artistic style transfer which has a key advantage: it can be restructured at runtime to incorporate, *e.g.*, the photorealistic prior proposed by Luan *et al.* [2]. This property is made possible because our network architecture is derived by unrolling the gradient descent algorithm used to minimise Gatys *et al.*'s style transfer loss: modification of this loss implies modifications of the gradient descent algorithm which can be passed on our network architecture. Finally, let us highlight that, beyond the style transfer problem, this principle could be used to derive flexible deep architectures architectures to solve other problems formulated as a minimisation of a loss function.

## References

[1] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *CVPR*, 2016.

[2] F. Luan, S. Paris, E. Shechtman, and K. Bala, "Deep photo style transfer," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4990–4998.

[3] L. A. Gatys, A. S. Ecker, and Bet, "Texture synthesis using convolutional neural networks," in *NIPS*, 2015.

[4] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *ECCV*, 2016.

[5] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky, "Texture networks: Feed-forward synthesis of textures and stylized images," 2016.

---

[9]The computation time does not include the computation of $\mathsf{L}$ nor the final guided filtering step, and was measured on an Nvidia Tesla M60 Maxwell GPU.

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

[7] V. Dumoulin, J. Shlens, and M. Kudlur, "A learned representation for artistic style," in *ICLR*, 2017.

[8] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang, "Universal style transfer via feature transforms," in *NIPS*, 2017.

[9] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *ICCV*, 2017, pp. 1501–1510.

[10] G. Ghiasi, H. Lee, M. Kudlur, V. Dumoulin, and J. Shlens, "Exploring the structure of a real-time, arbitrary neural artistic stylization network," in *BMVC*, 2017.

[11] Y. Li, M.-Y. Liu, X. Li, M.-H. Yang, and J. Kautz, "A closed-form solution to photorealistic image stylization," *arXiv:1802.06474*, 2018.

[12] E. Reinhard, M. Adhikhmin, B. Gooch, and P. Shirley, "Color transfer between images," *IEEE Computer graphics and applications*, vol. 21, no. 5, pp. 34–41, 2001.

[13] A. A. Efros and W. T. Freeman, "Image quilting for texture synthesis and transfer," in *SIGGRAPH*, 2001.

[14] A. Hertzmann, C. E. Jacobs, N. Oliver, B. Curless, and D. H. Salesin, "Image analogies," in *SIGGRAPH*, 2001.

[15] M. He, J. Liao, L. Yuan, and P. V. Sander, "Neural color transfer between images," *arXiv:1710.00756*, 2017.

[16] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *ICML*, 2010.

[17] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM journal on imaging sciences*, vol. 2, no. 1, pp. 183–202, 2009.

[18] A. Mousavi and R. G. Baraniuk, "Learning to invert: Signal recovery via deep convolutional networks," in *ICASSP*, 2017.

[19] C. Metzler, A. Mousavi, and R. Baraniuk, "Learned d-amp: Principled neural network based compressive image recovery," in *NIPS*, 2017.

[20] M. Borgerding, P. Schniter, and S. Rangan, "Amp-inspired deep networks for sparse linear inverse problems," *IEEE Trans. on Signal Processing*, 2017.

[21] J. Zhang and B. Ghanem, "Ista-net: Iterative shrinkage-thresholding algorithm inspired deep network for image compressive sensing," *arXiv preprint arXiv:1706.07929*, 2017.

[22] S. Wang, S. Fidler, and R. Urtasun, "Proximal deep structured models," in *NIPS*, 2016.

[23] R. Liu, X. Fan, S. Cheng, X. Wang, and Z. Luo, "Proximal alternating direction network: A globally converged deep unrolling framework," *arXiv preprint arXiv:1711.07653*, 2017.

[24] Y. Chen and T. Pock, "Trainable nonlinear reaction diffusion: A flexible framework for fast and effective image restoration," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1256–1272, 2017.

[25] J. Tompson, K. Schlachter, P. Sprechmann, and K. Perlin, "Accelerating eulerian fluid simulation with convolutional networks," *arXiv preprint arXiv:1607.03597*, 2016.

[26] A. Tewari, M. Zollhöfer, H. Kim, P. Garrido, F. Bernard, P. Pérez, and C. Theobalt, "Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction," in *ICCV*, 2017.

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

[28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[29] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, Y. W. Teh and M. Titterington, Eds., vol. 9. Chia Laguna Resort, Sardinia, Italy: PMLR, 13–15 May 2010, pp. 249–256. [Online]. Available: http://proceedings.mlr.press/v9/glorot10a.html

[30] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv:1412.6980*, 2014.

[31] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*. European Conference on Computer Vision, September 2014. [Online]. Available: https://www.microsoft.com/en-us/research/publication/microsoft-coco-common-objects-in-context/

[32] A. Levin, D. Lischinski, and Y. Weiss, "A closed-form solution to natural image matting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 2, pp. 228–242, 2008.

[33] D. Shuman, S. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE Signal Processing Magazine*, vol. 30, no. 3, pp. 83–98, 2013.

[34] D. K. Hammond, P. Vandergheynst, and R. Gribonval, "Wavelets on graphs via spectral graph theory," *Appl. Comput. Harmon. Anal.*, vol. 30, no. 2, pp. 129–150, 2011.

[35] E. D. Napoli, E. Polizzo, and Y. Saad, "Efficient estimation of eigenvalue counts in an interval," *Numerical Linear Algebra with Applications*, vol. 23, no. 4, pp. 674–692, 2016.

[36] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NIPS*, 2016.

[37] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv:1609.02907*, 2016.

Figure 8. *Adjusting the intensity of style transfer at runtime by varying $\alpha$ in (19) for an artistic result (top, with Content & Style 1) or a photorealistic result (bottom, with Content & Style 2).*

## Appendix A - Varying the Intensity of the Transfer

The intensity of the transfer can be modified by adjusting the parameter $\lambda_s$ in

$$\mathcal{L}(\mathsf{X}) = \lambda_c \, \mathcal{L}_c(\mathsf{X}, \mathsf{X}_c) + \lambda_s \, \mathcal{L}_s(\mathsf{X}, \mathsf{X}_s). \tag{17}$$

If one multiplies $\lambda_s$ by $\alpha > 0$, and keeps the stepsize $\mu$ constant, the gradient descent algorithm becomes

$$\mathsf{X}^{(t+1)} = \mathsf{X}^{(t)} - \mu \left[ \lambda_c \, \nabla \mathcal{L}_c(\mathsf{X}^{(t)}) + \alpha \lambda_s \, \nabla \mathcal{L}_s(\mathsf{X}^{(t)}) \right]. \tag{18}$$

This suggests that the intensity of the transfer can be adjusted at runtime as follows

$$\mathsf{X}^{(t+1)} = \mathsf{X}^{(t)} - \alpha \, g_t \left( \mathsf{X}^{(t)} \right), \tag{19}$$

in our network.

We show the effect of varying $\alpha$ in Fig. 8. Results are presented for artistic and photorealistic style transfers. These results show that we still have a degree of freedom to adjust the intensity of the style, even though the network was trained with $\alpha = 1$.

Let us highlight that all the results presented in this work are generated with $\alpha = 1$ in absence of graph filtering and with $\alpha = 1.2$ in presence of graph filtering in the network. This permits us to correct a loss in color saturation that we observe in most results when graph filtering is introduced in the network.

## Appendix B - Effect of the Guided Filter

As in [2] and [11], we authorise ourselves to post-process the obtained stylized images with a guided filter (using the content image as guide). Let us note that we did not use the code provided by [2] or [11] but used the OpenCV implementation of the guided filter. We discuss in this section the effect of this post-processing. In particular, we show that this post-processing permits us to remove remaining artefacts, such as compression artefacts, or enhances some details smoothed during style transfer. However, let us emphasize that the effect of the guided filter is limited to the improvement of these details and that the graph filtering introduced at runtime almost yields final results.
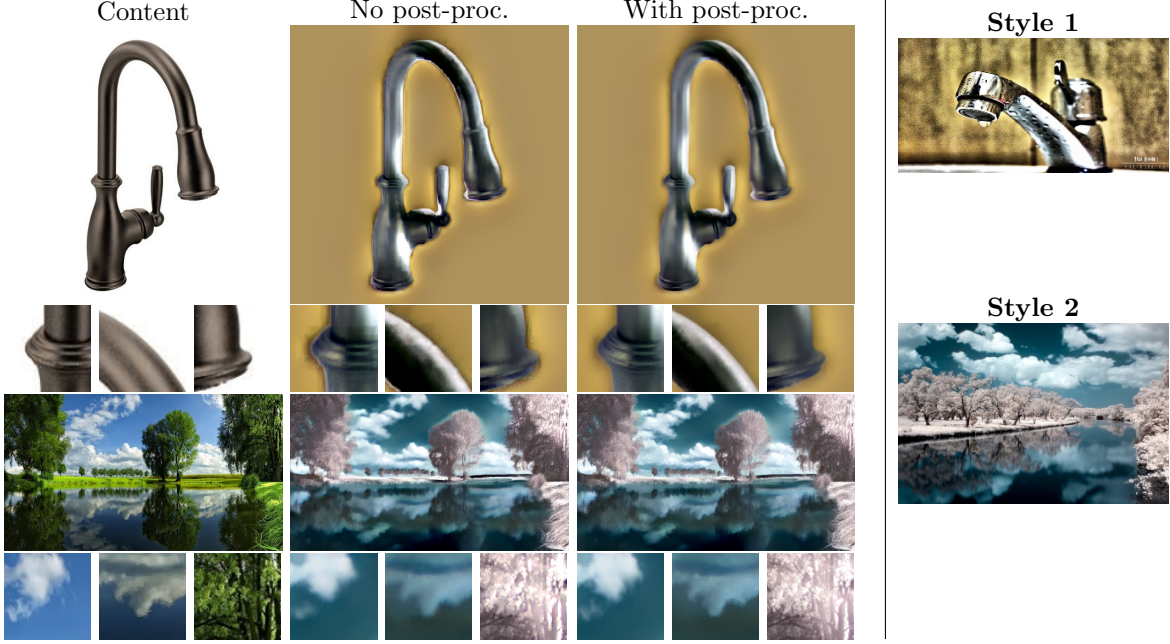
Figure 9. *Effect of post-processing with a guided filter. First column: Content images. Second: Results without guided filtering. Third: Results with guided filtering. Fourth: Style images. Please refer to Appendix B for an interpretation of the results.*

In the first example of Fig. 9, we see that the guided filter removes artefacts that appear near the boundary of the object. These artefacts are most certainly due to image compression which perturbed the construction of the matting Laplacian, hence their presence after graph filtering. In the last example of Fig. 9, we notice that the guided filter is able to sharpen some details and edges.

## Appendix C - Additional Style Transfer Results

We present additional photorealistic style tranfer results in Fig. 10 and compare with those obtained in [2]. In the first two results in Fig. 10, we notice that the color palette of the results obtained in [2] still contains colors present in the original content $X_c$ (blue of the sky in the first result, and yellow of the wall in the second). On the contrary, our results appear more faithful to the color palette of the style images. In the fourth result, we remark that both methods achieve similar qualities. In the last two examples, the color palette of the results obtained in [2] seems richer and more faithful to the color palette of the style images than in our results. Nevertheless, the results of [2] suffers from artefacts (see, *e.g.*, the sky in the last example). The local affinity constraint imposed via the matting Laplacian $L$ is enforced more strongly in our case thanks to the projection onto span($U_k$), instead of just a simpler regularisation with $L$ in [2]. This explains the reduced number of artefacts but also, sometimes, the limited color palettes in our results.

Finally, we present additional artistic style tranfer results in Fig. 11 to illustrate the fact that the method works on several styles. The results that we obtain are similar to those obtained with the method of [1].

## Appendix D - Limitations of the Method

In this section, we discuss some style transfer results which are not satisfying. Examples of such results are presented in Fig. 12

A typical case where the method fails is when we start from a content image essentially monochromatic and use a colorful style image (see the first two examples in Fig. 12). We notice that the color
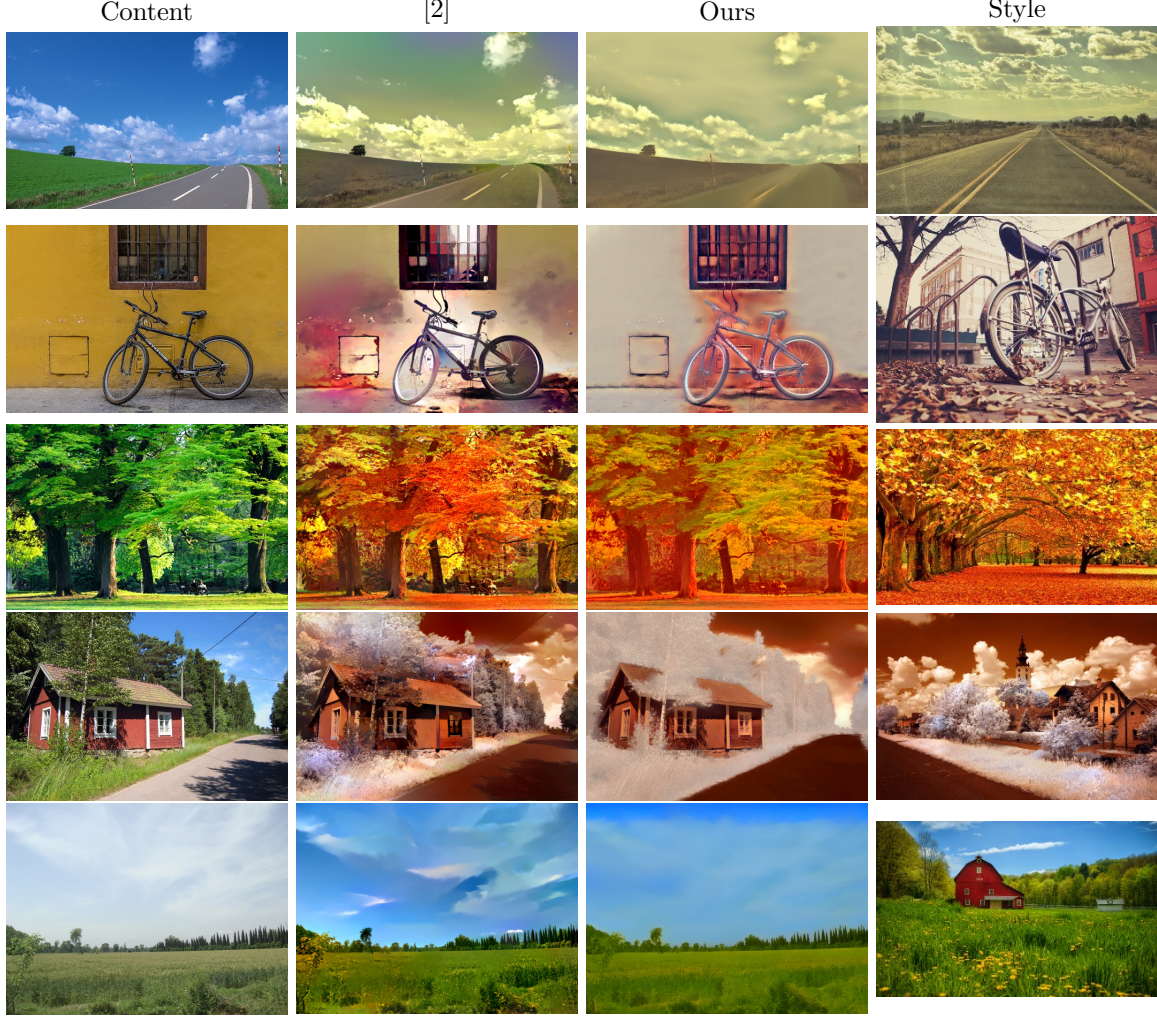
Content      [2]      Ours      Style



Figure 10. *Photorealistic style transfer results. First column: Content images. Second: Results obtained by* [2]. *Third: Our results. Fourth: Style images. Please refer to Appendix C for comments about these results.*

in our results lacks saturation unlike the results in [2] (even though we do not judge these results really photorealistic). One of the reasons explaining this lack of saturation is the presence of graph filtering itself. Indeed, the matting Laplacian enforces the transformation to be locally affine; the content image being essentially monochromatic, our graph filtering thus essentially averages locally all the colors. In several examples using $\alpha > 1$ permits us to attenuate this effect. Nevertheless, in the first two examples in Fig. 12, this lack of color saturation is also visible in the absence of graph filtering. The reason for this behaviour needs further investigation. Some hypotheses are: Not enough monochromatic images in the training datasets; Not enough iteration layers in our network; Not enough feature channels in the forward maps.

Finally, in some cases, the network stylises too differently parts of the content image that were initially similar. This is for example quite noticeable on the houses of the two last results in Fig. 12 where extra windows appear in the stylised images obtained without graph filtering. The presence of graph filtering in the network in not enough to correct this effect. This filtering averages the style across the image, giving the impression that a veil was added in front of the content image. The reason for this effect was not clearly identified. A likely reason is that the style weight $\lambda_s$ in the training
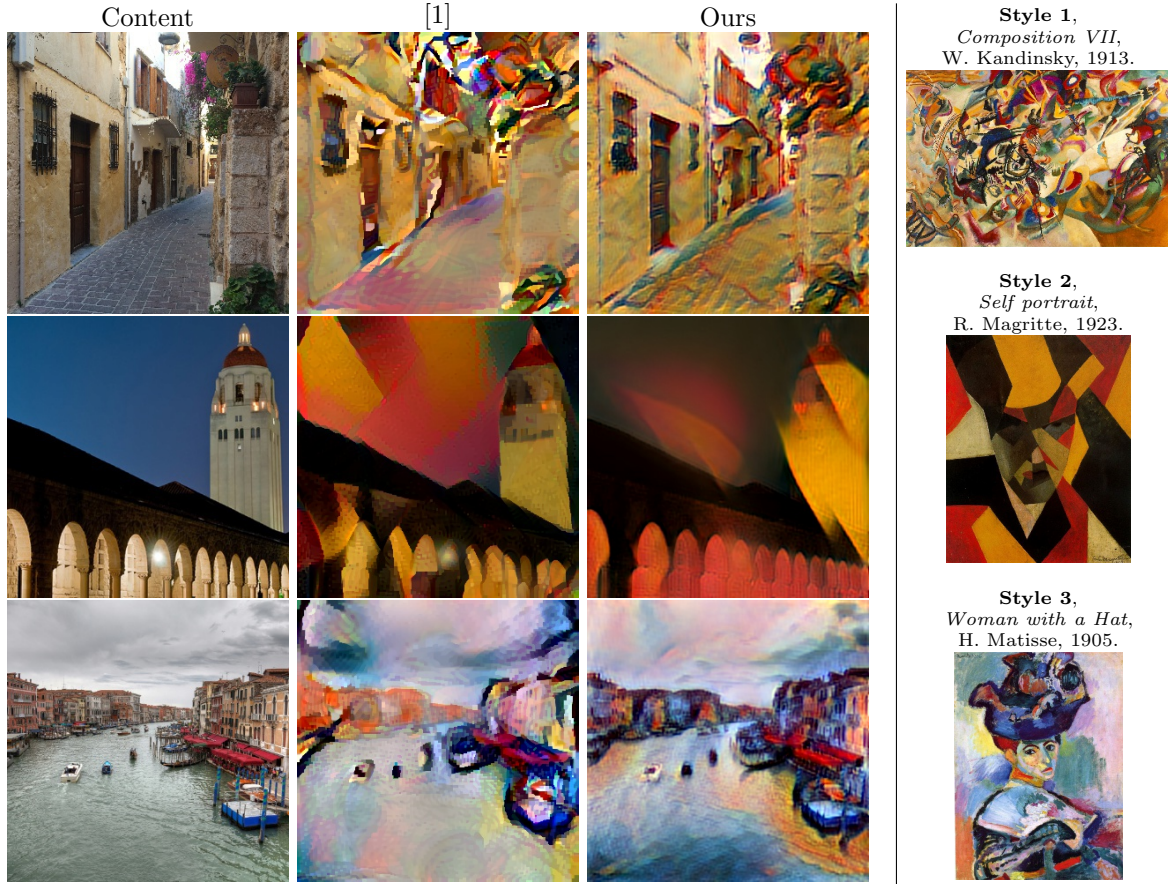
Figure 11. *Artistic style transfer results. First column: Content images. Second: Results obtained with the method of* [1]. *Third: Our results. Fourth: Style images.*

loss was too high for these styles. We noticed that modifying $\alpha$ in (19) at runtime is not sufficient to correct this hypothetical poor choice of $\lambda_s$.

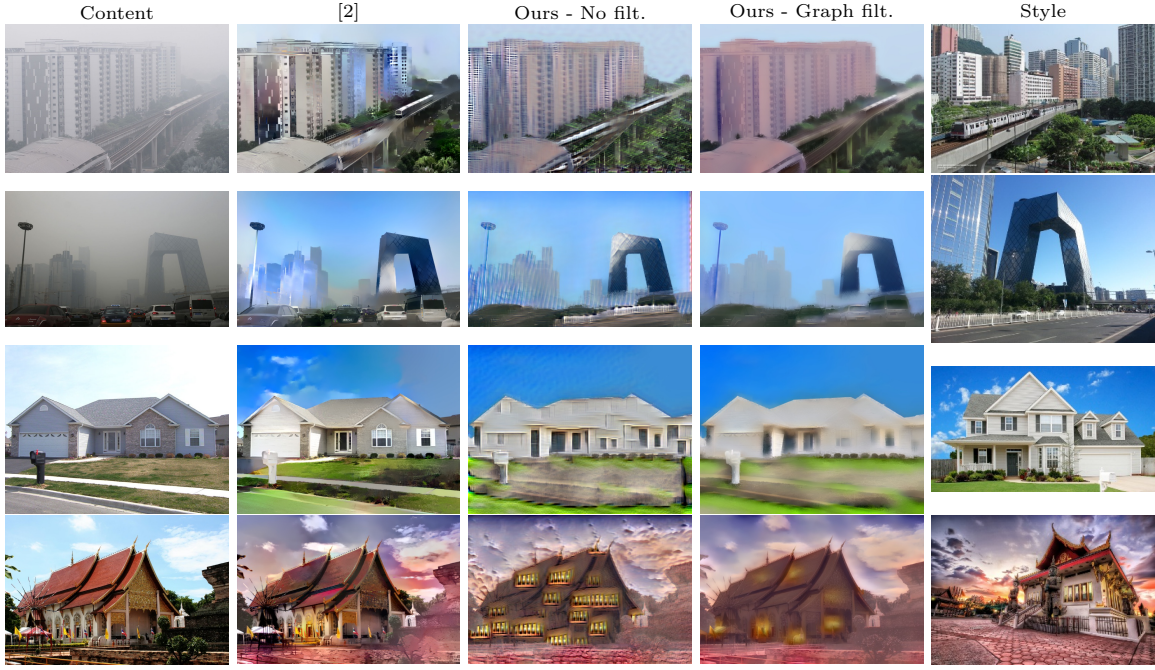| Content | [2] | Ours - No filt. | Ours - Graph filt. | Style |
|---------|-----|-----------------|--------------------|-------|



Figure 12. *Examples of results where style transfer failed. First column: Content images. Second: Results obtained by [2]. Third: Our results without graph filtering at runtime. Fourth: Our results with graph filtering at runtime. Fifth: Style images. Please refer to Appendix D for comments.*