

Artistic style transfer for videos

Manuel Ruder, Alexey Dosovitskiy, Thomas Brox

Department of Computer Science
University of Freiburg
{rudera, dosovits, brox}@cs.uni-freiburg.de

Abstract. In the past, manually re-drawing an image in a certain artistic style required a professional artist and a long time. Doing this for a video sequence single-handed was beyond imagination. Nowadays computers provide new possibilities. We present an approach that transfers the style from one image (for example, a painting) to a whole video sequence. We make use of recent advances in style transfer in still images and propose new initializations and loss functions applicable to videos. This allows us to generate consistent and stable stylized video sequences, even in cases with large motion and strong occlusion. We show that the proposed method clearly outperforms simpler baselines both qualitatively and quantitatively.

1 Introduction

There have recently been a lot of interesting contributions to the issue of style transfer using deep neural networks. Gatys et al. [3] proposed a novel approach using neural networks to capture the style of artistic images and transfer it to real world photographs. Their approach uses high-level feature representations of the images from hidden layers of the VGG convolutional network [10] to separate and reassemble content and style. This is done by formulating an optimization problem that, starting with white noise, searches for a new image showing similar neural activations as the *content image* and similar feature correlations (expressed by a Gram matrix) as the *style image*.

The present paper builds upon the approach from Gatys et al. [3] and extends style transfer to video sequences. Given an artistic image, we transfer its particular style of painting to the entire video. Processing each frame of the video independently leads to flickering and false discontinuities, since the solution of the style transfer task is not stable. To regularize the transfer and to preserve smooth transition between individual frames of the video, we introduce a temporal constraint that penalizes deviations between two frames. The temporal constraint takes the optical flow from the original video into account: instead of penalizing the deviations from the previous frame, we penalize deviation along the point trajectories. Disoccluded regions as well as motion boundaries are excluded from the penalizer. This allows the process to rebuild disoccluded regions and distorted motion boundaries while preserving the appearance of the rest of the image, see Fig. 1.

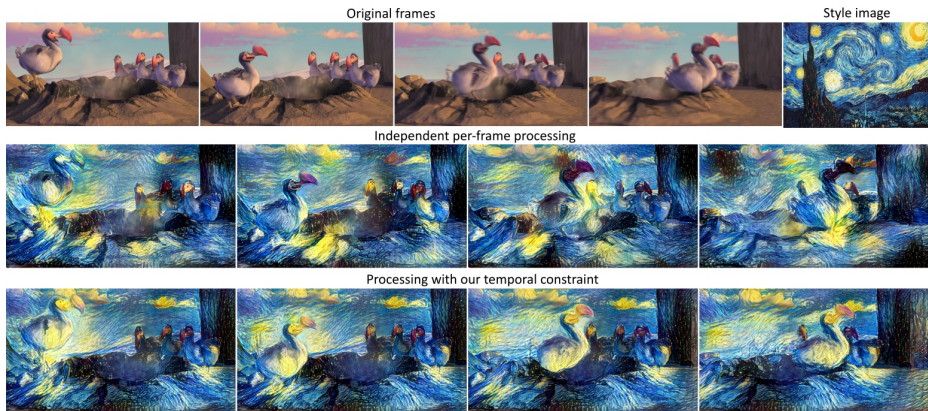


Fig. 1. Scene from *Ice Age* (2002) processed in the style of *The Starry Night*. Comparing independent per-frame processing to our time consistent approach, the latter is clearly preferable. Best observed in the supplemental video, see section 8.1.

In addition, we present two extensions of our approach. The first one aims on improving the consistency over larger periods of time. When a region that is occluded in some frame and disoccluded later gets rebuilt during the process, most likely this region will have a different appearance than before the occlusion. To solve this, we make use of long term motion estimates. This allows us to enforce consistency of the synthesized frames before and after the occlusion.

Secondly, the style transfer tends to create artifacts at the image boundaries. For static images, these artifacts are hardly visible, yet for videos with strong camera motion they move towards the center of the image and get amplified. We developed a multi-pass algorithm, which processes the video in alternating directions using both forward and backward flow. This results in a more coherent video.

We quantitatively evaluated our approach in combination with different optical flow algorithms on the Sintel benchmark. Additionally we show qualitative results on several movie shots. We were able to successfully eliminate most of the temporal artifacts and can create smooth and coherent stylized videos.

2 Related work

Style transfer using deep networks: Gatys et al. [3] showed remarkable results by using the VGG-19 deep neural network for style transfer. Their approach was taken up by various follow-up papers that, among other things, proposed different ways to represent the style within the neural network. Li et al. [5] suggested an approach to preserve local patterns of the style image. Instead of using a global representation of the style, computed as Gram matrix, they used patches of the neural activation from the style image. Nikulin et al. [7] tried the style transfer algorithm by Gatys et al. on other nets than VGG and proposed

several variations in the way the style of the image is represented to archive different goals like illumination or season transfer. However, we are not aware of any work that applies this kind of style transfer to videos.

Painted animations: One common approach to create video sequences with an artistic style is to generate artificial brush strokes to repaint the scene. Different artistic styles are gained by modifying various parameters of these brush strokes, like thickness, or by using different brush placement methods. To achieve temporal consistency Litwinowicz [6] was one of the first who used optical flow. In his approach, brush strokes were generated for the first frame and then moved along the flow field. Later, this approach was refined. Hays et al. [4] proposed new stylistic parameters for the brush strokes to mimic different artistic styles. O’Donovan et al. [8] formulated an energy optimization problem for an optimal placement and shape of the brush strokes and also integrated a temporal constraint into the optimization problem by penalizing changes in shape and width of the brush strokes compared to the previous frame. These approaches are similar in spirit to what we are doing, but they are only capable of applying a restricted class of artistic styles.

3 Style transfer in still images

In this section, we briefly review the style transfer approach introduced by Gatys et al. [3]. The aim is to generate a stylized image \mathbf{x} showing the content of an image \mathbf{p} in the style of an image \mathbf{a} . Gatys et al. formulated an energy minimization problem consisting of a *content loss* and a *style loss*. The key idea is that features extracted by a convolutional network carry information about the content of the image, while the correlations of these features encode the style.

We denote by $\Phi^l(\cdot)$ the function implemented by the part of the convolutional network from input up to the layer l . The feature maps extracted by the network from the original image \mathbf{p} , the style image \mathbf{a} and the stylized image \mathbf{x} we denote by $\mathbf{P}^l = \Phi^l(\mathbf{p})$, $\mathbf{S}^l = \Phi^l(\mathbf{a})$ and $\mathbf{F}^l = \Phi^l(\mathbf{x})$ respectively. The dimensionality of these feature maps we denote by $N_l \times M_l$, where N_l is the number of filters (channels) in the layer, and M_l is the spatial dimensionality of the feature map, that is, the product of its width and height.

The content loss, denoted as $\mathcal{L}_{content}$, is simply the mean squared error between $\mathbf{P}^l \in \mathbb{R}^{N_l \times M_l}$ and $\mathbf{F}^l \in \mathbb{R}^{N_l \times M_l}$. This loss need not be restricted to only one layer. Let $L_{content}$ be the set of layers to be used for content representation, then we have:

$$\mathcal{L}_{content}(\mathbf{p}, \mathbf{x}) = \sum_{l \in L_{content}} \frac{1}{N_l M_l} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2. \quad (1)$$

The style loss is also a mean squared error, but between the correlations of the filter responses expressed by their Gram matrices $A^l \in \mathbb{R}^{N_l \times N_l}$ for the style image \mathbf{a} and $G^l \in \mathbb{R}^{N_l \times N_l}$ for the stylized image \mathbf{x} . These are computed as $A_{ij}^l = \sum_{k=1}^{M_l} S_{ik}^l S_{jk}^l$ and $G_{ij}^l = \sum_{k=1}^{M_l} F_{ik}^l F_{jk}^l$. As above, let L_{style} be the set of layers

we use to represent the style, then the style loss is given by:

$$\mathcal{L}_{style}(\mathbf{a}, \mathbf{x}) = \sum_{l \in \mathcal{L}_{style}} \frac{1}{N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2 \quad (2)$$

Overall, the loss function is given by

$$\mathcal{L}_{singleimage}(\mathbf{p}, \mathbf{a}, \mathbf{x}) = \alpha \mathcal{L}_{content}(\mathbf{p}, \mathbf{x}) + \beta \mathcal{L}_{style}(\mathbf{a}, \mathbf{x}), \quad (3)$$

with weighting factors α and β governing the importance of the two components.

The stylized image is computed by minimizing this energy with respect to \mathbf{x} using gradient-based optimization. Typically it is initialized with random Gaussian noise. However, the loss function is non-convex, therefore the optimization is prone to falling into local minima. This makes the initialization of the stylized image important, especially when applying the method to frames of a video.

4 Style transfer in videos

We use the following notation: $\mathbf{p}^{(i)}$ is the i^{th} frame of the original video, \mathbf{a} is the style image and $\mathbf{x}^{(i)}$ are the stylized frames to be generated. Furthermore, we denote by $\mathbf{x}'^{(i)}$ the initialization of the style optimization algorithm at frame i . By x_j we denote the j^{th} component of a vector \mathbf{x} .

4.1 Short-term consistency by initialization

When the style transfer for consecutive frames is initialized by independent Gaussian noise, two frames of a video converge to very different local minima, resulting in a strong flickering. The most basic way to yield temporal consistency is to initialize the optimization for the frame $i+1$ with the stylized frame i . Areas that have not changed between the two frames are then initialized with the desired appearance, while the rest of the image has to be rebuilt through the optimization process.

If there is motion in the scene, this simple approach does not perform well, since moving objects are initialized incorrectly. Thus, we take the optical flow into account and initialize the optimization for the frame $i+1$ with the previous stylized image warped: $\mathbf{x}'^{(i+1)} = \omega_i^{i+1}(\mathbf{x}^{(i)})$. Here ω_i^{i+1} denotes the function that warps a given image using the optical flow field that was estimated between image $\mathbf{p}^{(i)}$ and $\mathbf{p}^{(i+1)}$. Clearly, the first frame of the stylized video $\mathbf{x}'^{(1)}$ still has to be initialized randomly.

We experimented with two state-of-the-art optical flow estimation algorithms: DeepFlow [12] and EpicFlow [9]. Both are based on Deep Matching [12]: DeepFlow combines it with a variational approach, while EpicFlow relies on edge-preserving sparse-to-dense interpolation.

4.2 Temporal consistency loss

To enforce stronger consistency between adjacent frames we additionally introduce an explicit consistency penalty to the loss function. This requires detection of disoccluded regions and motion boundaries. To detect disocclusions, we perform a forward-backward consistency check of the optical flow [11]. Let $\mathbf{w} = (u, v)$ be the optical flow in forward direction and $\hat{\mathbf{w}} = (\hat{u}, \hat{v})$ the flow in backward direction. Denote by $\tilde{\mathbf{w}}$ the forward flow warped to the second image:

$$\tilde{\mathbf{w}}(x, y) = \mathbf{w}((x, y) + \hat{\mathbf{w}}(x, y)). \quad (4)$$

In areas without disocclusion, this warped flow should be approximately the opposite of the backward flow. Therefore we mark as disocclusions those areas where the following inequality holds:

$$|\tilde{\mathbf{w}} + \hat{\mathbf{w}}|^2 > 0.01(|\tilde{\mathbf{w}}|^2 + |\hat{\mathbf{w}}|^2) + 0.5 \quad (5)$$

Motion boundaries are detected using the following inequality:

$$|\nabla \hat{u}|^2 + |\nabla \hat{v}|^2 > 0.01|\hat{\mathbf{w}}|^2 + 0.002 \quad (6)$$

Coefficients in inequalities (5) and (6) are taken from Sundaram et al. [11].

The temporal consistency loss function penalizes deviations from the warped image in regions where the optical flow is consistent and estimated with high confidence:

$$\mathcal{L}_{temporal}(\mathbf{x}, \boldsymbol{\omega}, \mathbf{c}) = \frac{1}{D} \sum_{k=1}^D c_k \cdot (x_k - \omega_k)^2. \quad (7)$$

Here $\mathbf{c} \in [0, 1]^D$ is per-pixel weighting of the loss and $D = W \times H \times C$ is the dimensionality of the image. We define the weights $\mathbf{c}^{(i-1, i)}$ between frames $i-1$ and i as follows: 0 in disoccluded regions (as detected by forward-backward consistency) and at the motion boundaries, and 1 everywhere else. Potentially weights between 0 and 1 could be used to incorporate the certainty of the optical flow prediction. The overall loss takes the form:

$$\begin{aligned} \mathcal{L}_{shortterm}(\mathbf{p}^{(i)}, \mathbf{a}, \mathbf{x}^{(i)}) &= \alpha \mathcal{L}_{content}(\mathbf{p}^{(i)}, \mathbf{x}^{(i)}) + \beta \mathcal{L}_{style}(\mathbf{a}, \mathbf{x}^{(i)}) \\ &+ \gamma \mathcal{L}_{temporal}(\mathbf{x}^{(i)}, \omega_{i-1}^i(\mathbf{x}^{(i-1)}), \mathbf{c}^{(i-1, i)}). \end{aligned} \quad (8)$$

We optimize one frame after another, thus $\mathbf{x}^{(i-1)}$ refers to the already stylized frame $i-1$.

Furthermore we experimented with the more robust absolute error instead of squared error for the temporal consistency loss; results are shown in section 8.

4.3 Long-term consistency

The short-term model has the following limitation: when some areas are occluded in some frame and disoccluded later, these areas will likely change their

appearance in the stylized video. This can be counteracted by also making use of long-term motion, i.e. not only penalizing deviations from the previous frame, but also from temporally more distant frames. Let J denote the set of indices each frame should take into account, relative to the frame number. E.g. $J = \{1, 2, 4\}$ means frame i takes frames $i-1$, $i-2$ and $i-4$ into account. Then, the loss function with long-term consistency is given by:

$$\begin{aligned} \mathcal{L}_{longterm}(\mathbf{p}^{(i)}, \mathbf{a}, \mathbf{x}^{(i)}) &= \alpha \mathcal{L}_{content}(\mathbf{p}^{(i)}, \mathbf{x}^{(i)}) + \beta \mathcal{L}_{style}(\mathbf{a}, \mathbf{x}^{(i)}) \\ &+ \gamma \sum_{j \in J: i-j \geq 1} \mathcal{L}_{temporal}(\mathbf{x}^{(i)}, \omega_{i-j}^i(\mathbf{x}^{(i-j)}), \mathbf{c}_{long}^{(i-j,i)}) \end{aligned} \quad (9)$$

It is essential how the weights $\mathbf{c}_{long}^{(i-j,i)}$ are computed. Let $\mathbf{c}^{(i-j,i)}$ be the weights for the flow between image $i-j$ and i , as defined for the short-term model. The long-term weights $\mathbf{c}_{long}^{(i-j,i)}$ are computed as follows:

$$\mathbf{c}_{long}^{(i-j,i)} = \max(\mathbf{c}^{(i-j,i)} - \sum_{k \in J: i-k > i-j} \mathbf{c}^{(i-k,i)}, \mathbf{0}), \quad (10)$$

where \max is taken element-wise. This means, we first apply the usual short-term constraint. For pixels in disoccluded regions we look into the past until we find a frame in which these have consistent correspondences. An advantage over simply using $\mathbf{c}^{(i-j,i)}$ is that each pixel is connected only to the closest possible frame from the past. Since the optical flow computed over more frames is more erroneous than over fewer frames, this results in nicer videos. An empirical comparison of $\mathbf{c}^{(i-j,i)}$ and $\mathbf{c}_{long}^{(i-j,i)}$ is shown in the supplementary video (see section 8.1).

4.4 Multi-pass algorithm

We found that the output image tends to have less contrast and is less diverse near image boundaries than in other areas of the image. For mostly static videos this effect is hardly visible. However, in cases of strong camera motion the areas from image boundaries move towards other parts of the image, which leads to a lower image quality over time when combined with our temporal constraint. Therefore, we developed a multi-pass algorithm which processes the whole sequence in multiple passes and alternates between the forward and backward direction. Every pass consists of a relatively low number of iterations without full convergence. At the beginning, we process every frame independently. After that, we blend frames with non-disoccluded parts of previous frames warped according to the optical flow, then run the optimization algorithm for some iterations initialized with this blend. We repeat this blending and optimization to convergence.

Let $\mathbf{x}'^{(i)(j)}$ be the initialization of frame i in pass j and $\mathbf{x}^{(i)(j)}$ the corresponding output after some iterations of the optimization algorithm. When processed

in forward direction, the initialization of frame i is created as follows:

$$\mathbf{x}'^{(i)(j)} = \begin{cases} \mathbf{x}^{(i)(j-1)} & \text{if } i = 1, \\ \delta \mathbf{c}^{(i-1,i)} \circ \omega_{i-1}^i(\mathbf{x}^{(i-1)(j)}) + (\bar{\delta} \mathbf{1} + \delta \bar{\mathbf{c}}^{(i-1,i)}) \circ \mathbf{x}^{(i)(j-1)} & \text{else.} \end{cases} \quad (11)$$

Here \circ denotes element-wise vector multiplication, δ and $\bar{\delta} = 1 - \delta$ are the blend factors, $\mathbf{1}$ is a vector of all ones, and $\bar{\mathbf{c}} = \mathbf{1} - \mathbf{c}$.

Analogously, the initialization for a backward direction pass is:

$$\mathbf{x}'^{(i)(j)} = \begin{cases} \mathbf{x}^{(i)(j-1)} & \text{if } i = N_{\text{frames}} \\ \delta \mathbf{c}^{(i+1,i)} \circ \omega_{i+1}^i(\mathbf{x}^{(i+1)(j)}) + (\bar{\delta} \mathbf{1} + \delta \bar{\mathbf{c}}^{(i+1,i)}) \circ \mathbf{x}^{(i)(j-1)} & \text{else} \end{cases} \quad (12)$$

The multi-pass algorithm can be combined with the temporal consistency loss described above. We achieved good results when we disabled the temporal consistency loss in several initial passes and enabled it in later passes only after the images had stabilized.

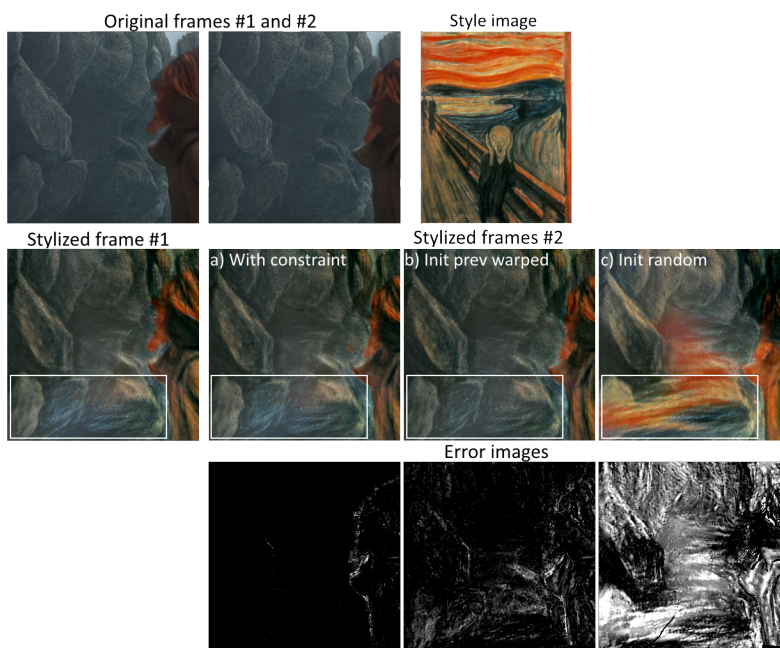


Fig. 2. Close-up of a scene from Sintel, combined with *The Scream* painting. **a)** With temporal constraint **b)** Initialized with previous image warped, but without the constraint **c)** Initialized randomly. The marked regions show most visible differences. *Error images* show the contrast-enhanced absolute difference between frame #1 and frame #2 warped back using ground truth optical flow, as used in our evaluation. The effect of the temporal constraint is very clear in the error images and in the corresponding video.

5 Experiments

In this section, we briefly describe implementation details and present experimental results produced with different versions of our algorithm. While we did our best to make the paper self-contained, it is not possible to demonstrate effects like video flickering in still images. We therefore advise the readers to watch the supplementary video, which is available at https://youtu.be/vQk_Sf17kSc.

5.1 Implementation details

Our implementation¹ is based on the Torch [2] implementation called *neural-style*². We used the following layers of the VGG-19 network [10] for computing the losses: *relu4_2* for the content and *relu1_1, relu2_1, relu3_1, relu4_1, relu5_1* for the style. The energy function was minimized using L-BFGS. For precise evaluation we incorporated the following strict stopping criterion: the optimization was considered converged if the loss did not change by more than 0.01% during 50 iterations. This typically resulted in roughly 2000 to 3000 iterations for the first frame and roughly 400 to 800 iterations for subsequent frames when optimizing with our temporal constraint, depending on the amount of motion and the complexity of the style image. Using a convergence threshold of 0.1% cuts the number of iterations and the running time in half, and we found it still produces reasonable results in most cases. However, we used the stronger criterion in our experiments for the sake of accuracy.

For videos of resolution 350×450 we used weights $\alpha = 1$ and $\beta = 20$ for the content and style losses, respectively (default values from *neural-style*), and weight $\gamma = 200$ for the temporal losses. However, the weights should be adjusted if the video resolution is different. We provide the details in section 7.2.

For our multi-pass algorithm, we used 100 iterations per pass and set $\delta = 0.5$, but we needed at least 10 passes for good results, so this algorithm needs more computation time than our previous approaches.

We used DeepMatching, DeepFlow and EpicFlow implementations provided by the authors of these methods. We used the "improved-settings" flag in DeepMatching and the default settings for DeepFlow and EpicFlow.

Runtime For the relaxed convergence threshold of 0.1% with random initialization the optimization process needed on average roughly eight to ten minutes per frame at a resolution of 1024×436 on an Nvidia Titan X GPU. When initialized with the warped previous frame and combined with our temporal loss, the optimization converges 2 to 3 times faster, three minutes on average. Optical flow computation runs on a CPU and takes roughly 3 minutes per frame pair (forward and backward flow together), therefore it can be performed in parallel with the style transfer. Hence, our modified algorithm is roughly 3 times faster than naive per-frame processing, while providing temporally consistent output videos.

¹ GitHub: <https://github.com/manuelruder/artistic-videos>

² GitHub: <https://github.com/jcjohnson/neural-style>

5.2 Short-term consistency

We evaluated our short-term temporal loss on 5 diverse scenes from the MPI Sintel Dataset [1], with 20 to 50 frames of resolution 1024×436 pixels per scene, and 6 famous paintings (shown in section 7.1) as style images. The Sintel dataset provides ground truth optical flow and ground truth occlusion areas, which allows a quantitative study. We warped each stylized frame i back with the ground truth flow and computed the difference with the stylized frame $i - 1$ in non-disoccluded regions. We use the mean square of this difference (that is, the mean squared error) as a quantitative performance measure.

On this benchmark we compared several approaches: our short-term consistency loss with DeepFlow and EpicFlow, as well as three different initializations without the temporal loss: random noise, the previous stylized frame and the previous stylized frame warped with DeepFlow. We set $\alpha = 1$, $\beta = 100$, $\gamma = 400$.

A qualitative comparison is shown in Fig. 2. Quantitative results are in Table 2. The most straightforward approach, processing every frame independently, performed roughly an order of magnitude worse than our more sophisticated methods. In most cases, the temporal penalty significantly improved the results. The *ambush* scenes are exceptions, since they contain very large motion and the erroneous optical flow impairs the temporal constraint. Interestingly, on average DeepFlow performed slightly better than EpicFlow in our experiments.

Table 1. Short-term consistency benchmark results. Mean squared error of different methods on 5 video sequences, averaged over 6 styles, is shown. Pixel values in images were between 0 and 1.

	alley_2	ambush_5	ambush_6	bandage_2	market_6
DeepFlow	0.00061	0.0062	0.012	0.00084	0.0035
EpicFlow	0.00073	0.0068	0.014	0.00080	0.0032
Init prev warped	0.0016	0.0063	0.012	0.0015	0.0049
Init prev	0.010	0.018	0.028	0.0041	0.014
Init random	0.019	0.027	0.037	0.018	0.023

5.3 Long-term consistency and multi-pass algorithm

The short-term consistency benchmark presented above cannot evaluate the long-term consistency of videos (since we do not have long-term ground truth flow available) and their visual quality (this can only be judged by humans). We therefore excluded the long-term penalty and the multi-pass approach from the quantitative comparison and present only qualitative results. Please see the supplementary video for more results.

Fig. 3 shows a scene from Miss Marple where a person walks through the scene. Without our long-term consistency model, the background looks very different after the person passes by. The long-term consistency model keeps the background unchanged. Fig. 4 shows another scene from Miss Marple with fast camera motion. The multi-pass algorithm avoids the artifacts introduced by the basic algorithm.

6 Conclusion

We presented a set of techniques for style transfer in videos: suitable initialization, a loss function that enforces short-term temporal consistency of the stylized video, a loss function for long-term consistency, and a multi-pass approach. As a consequence, we can produce stable and visually appealing stylized videos even in the presence of fast motion and strong occlusion.

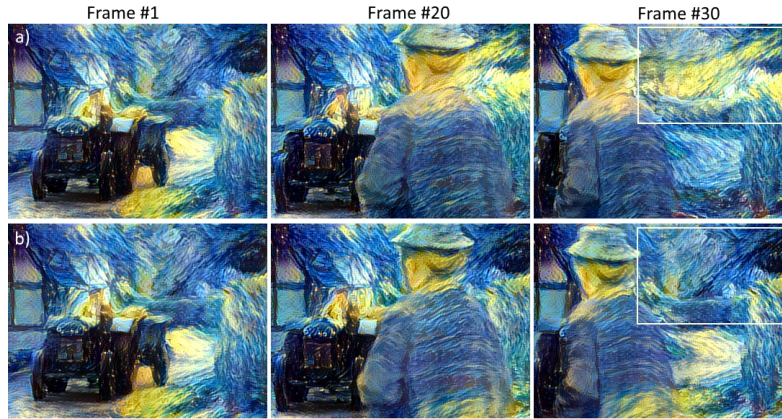


Fig. 3. Scene from Miss Marple, combined with The Starry Night painting. **a)** Short-term consistency only. **b)** Long-term consistency with $J = \{1, 10, 20, 40\}$. Corresponding video is linked in section 8.1.



Fig. 4. The multi-pass algorithm applied to a scene from Miss Marple. With the default method, the image becomes notably brighter and loses contrast, while the multi-pass algorithm yields a more consistent image quality over time. Corresponding video is linked in section 8.1.

References

1. Butler, D.J., Wulff, J., Stanley, G.B., Black, M.J.: A naturalistic open source movie for optical flow evaluation. In: A. Fitzgibbon et al. (Eds.) (ed.) European Conf. on Computer Vision (ECCV). pp. 611–625. Part IV, LNCS 7577, Springer-Verlag (Oct 2012)
2. Collobert, R., Kavukcuoglu, K., Farabet, C.: Torch7: A matlab-like environment for machine learning. In: BigLearn, NIPS Workshop (2011)
3. Gatys, L.A., Ecker, A.S., Bethge, M.: A neural algorithm of artistic style. CoRR abs/1508.06576 (2015), <http://arxiv.org/abs/1508.06576>
4. Hays, J., Essa, I.: Image and video based painterly animation. In: Proceedings of the 3rd International Symposium on Non-photorealistic Animation and Rendering. pp. 113–120. NPAR '04, ACM, New York, NY, USA (2004), <http://doi.acm.org/10.1145/987657.987676>
5. Li, C., Wand, M.: Combining markov random fields and convolutional neural networks for image synthesis. CoRR abs/1601.04589 (2016), <http://arxiv.org/abs/1601.04589>
6. Litwinowicz, P.: Processing images and video for an impressionist effect. In: Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques. pp. 407–414. SIGGRAPH '97, ACM Press/Addison-Wesley Publishing Co., New York, NY, USA (1997), <http://dx.doi.org/10.1145/258734.258893>
7. Nikulin, Y., Novak, R.: Exploring the neural algorithm of artistic style. CoRR abs/1602.07188 (2016), <http://arxiv.org/abs/1602.07188>
8. O'Donovan, P., Hertzmann, A.: Anipaint: Interactive painterly animation from video. IEEE Transactions on Visualization and Computer Graphics 18(3), 475–487 (2012)
9. Revaud, J., Weinzaepfel, P., Harchaoui, Z., Schmid, C.: EpicFlow: Edge-Preserving Interpolation of Correspondences for Optical Flow. In: CVPR 2015 - IEEE Conference on Computer Vision & Pattern Recognition. Boston, United States (Jun 2015), <https://hal.inria.fr/hal-01142656>
10. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. CoRR abs/1409.1556 (2014), <http://arxiv.org/abs/1409.1556>
11. Sundaram, N., Brox, T., Keutzer, K.: Dense point trajectories by gpu-accelerated large displacement optical flow (Sept 2010), <http://lmb.informatik.uni-freiburg.de/Publications/2010/Bro10e>
12. Weinzaepfel, P., Revaud, J., Harchaoui, Z., Schmid, C.: DeepFlow: Large displacement optical flow with deep matching. In: ICCV 2013 - IEEE International Conference on Computer Vision. pp. 1385–1392. IEEE, Sydney, Australia (Dec 2013), <https://hal.inria.fr/hal-00873592>

Supplementary material

7 Additional details of experimental setup

7.1 Style images

Style images we used for benchmark experiments on Sintel are shown in Figure 5.

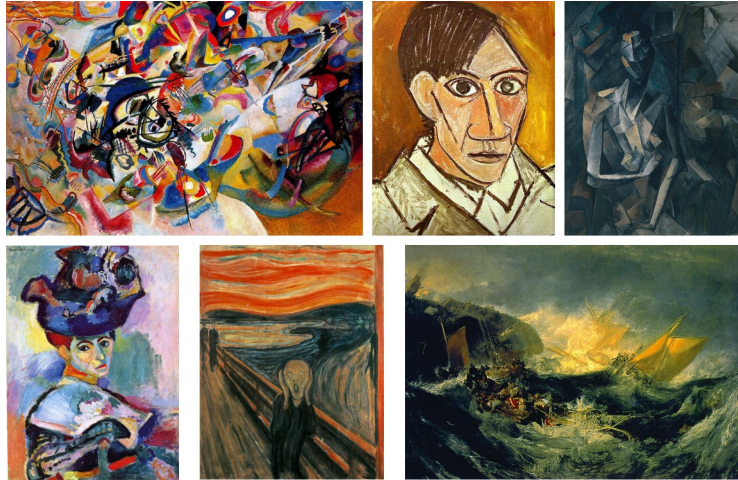


Fig. 5. Styles used for experiments on Sintel. Left to right, top to bottom: "Composition VII" by Wassily Kandinsky (1913), "Self-Portrait" by Pablo Picasso (1907), "Seated female nude" by Pablo Picasso (1910), "Woman with a Hat" by Henri Matisse (1905), "The Scream" by Edvard Munch (1893), "Shipwreck" by William Turner (1805).

7.2 Weighting of the loss components

As mentioned in the main paper, for best results the weights α , β and γ of different components of the loss function have to be adjusted depending on the resolution of the video. The settings we used for different resolutions are shown in Table 2.

Table 2. Weights of the loss function components for different input resolutions.

	350×450	768×432	1024×436
α (content)	1	1	1
β (style)	20	40	100
γ (temporal)	200	200	400

8 Additional experiments

8.1 Supplementary video

A supplementary video, available at https://youtu.be/vQk_Sf17kSc, shows moving sequences corresponding to figures from this paper, plus a number of additional results:

- Results of the basic algorithm on different sequences from Sintel with different styles
- Additional comparison of the basic and the multi-pass algorithm
- Additional comparison of the basic and the long-term algorithm
- Comparison of "naive" (c) and "advanced" (c_{long}) weighting schemes for long-term consistency
- Results of the algorithm on a number of diverse videos with different style images

8.2 Robust loss function for temporal consistency

We tried using the more robust absolute error instead of squared error for the temporal consistency loss. The weight for the temporal consistency was doubled in this case. Results are shown in Figure 6. While in some cases (left example in the figure) absolute error leads to slightly improved results, in other cases (right example in the figure) it causes large fluctuations. We therefore stick with mean squared error in all our experiments.

8.3 Effect of errors in optical flow estimation

The quality of results produced by our algorithm strongly depends on the quality of optical flow estimation. This is illustrated Figure 7. When the optical flow is correct (top right region of the image), the method manages to repair the artifacts introduced by warping in the disoccluded region. However, erroneous optical flow (tip of the sword in the bottom right) leads to degraded performance. Optimization process partially compensates the errors (sword edges get sharp), but cannot fully recover.

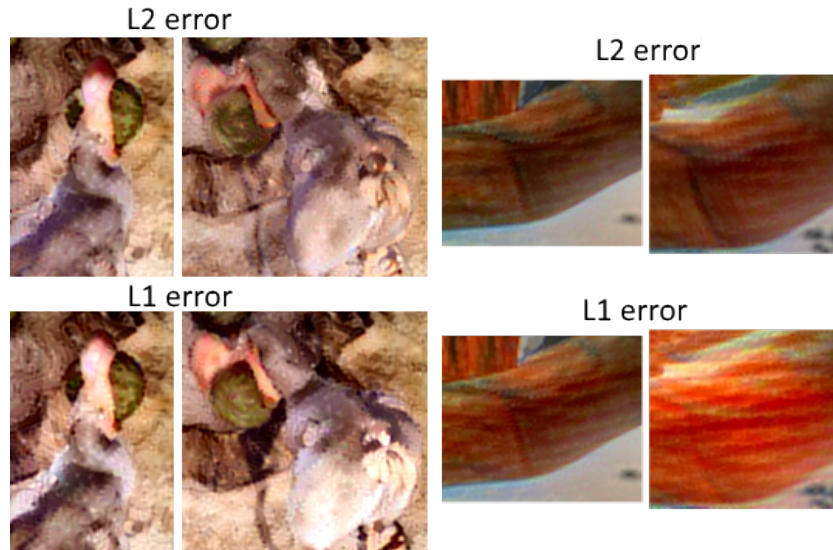


Fig. 6. *Left:* Scene from Ice Age (2002) where an absolute error function works better, because the movement of the bird wasn't captured correctly by the optical flow. *Right:* Extreme case from Sintel movie where a squared error is far superior.

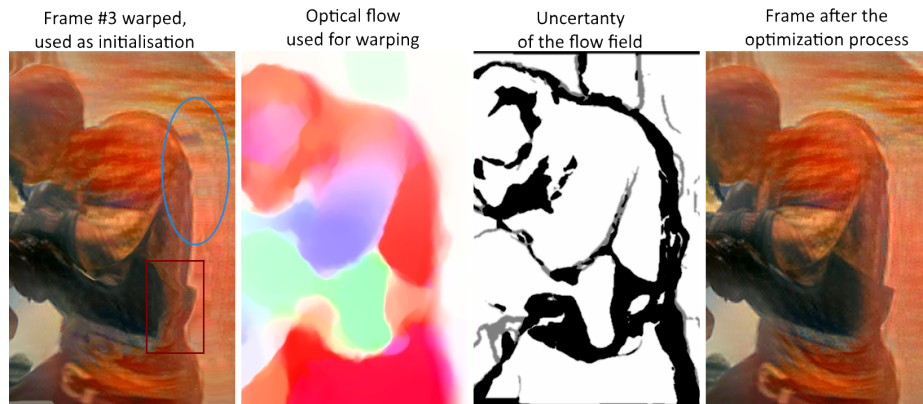


Fig. 7. Scene from the Sintel video showing how the algorithm deals with optical flow errors (red rectangle) and disocclusions (blue circle). Both artifacts are somehow repaired in the optimization process due to the exclusion of uncertain areas from our temporal constrain. Still, optical flow errors lead to imperfect results. The third image shows the uncertainty of the flow field in black and motion boundaries in gray.