# NEURAL STYLE TRANSFER FOR COMPUTER GAMES

**Eleftherios Ioannou, Steve Maddock**
Department of Computer Science
University of Sheffield
Sheffield, UK
{eioannou1, s.maddock}@sheffield.ac.uk

## ABSTRACT

Neural Style Transfer (NST) research has been applied to images, videos, 3D meshes and radiance fields, but its application to 3D computer games remains relatively unexplored. Whilst image and video NST systems can be used as a post-processing effect for a computer game, this results in undesired artefacts and diminished post-processing effects. Here, we present an approach for injecting depth-aware NST as part of the 3D rendering pipeline. Qualitative and quantitative experiments are used to validate our in-game stylisation framework. We demonstrate temporally consistent results of artistically stylised game scenes, outperforming state-of-the-art image and video NST methods.
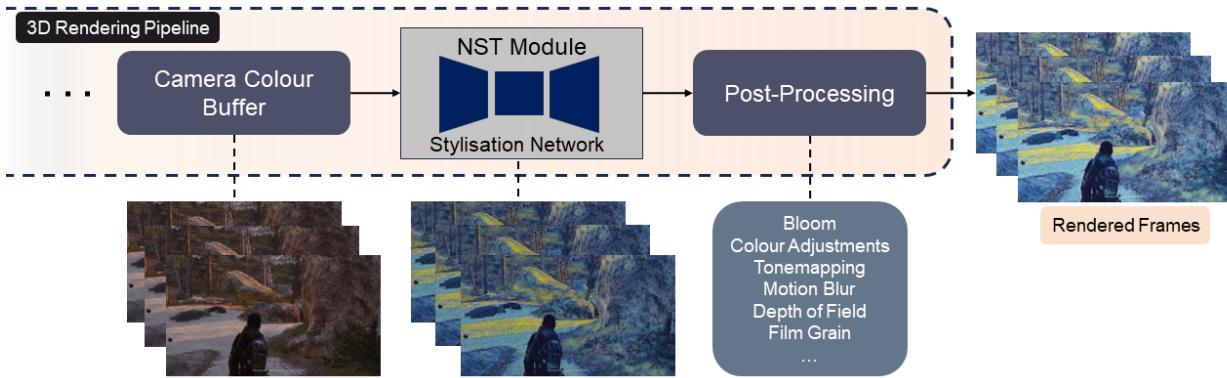
https://ioannoue.github.io/nst-for-computer-games.html



Figure 1: Our proposed framework injects NST as part of the 3D rendering pipeline.

## 1 Introduction

Neural Style Transfer (NST) refers to the process of changing the appearance of an input image based on a reference style image, whilst preserving the input image's underlying content. For example, a photograph of a landscape can be made to take on the style of a Van Gogh painting. More recently, NST has been extended to work for three-dimensional data, such as 3D meshes, point clouds, and radiance fields. Given the inherently creative and artistic nature of NST, another domain where its application holds immense potential is within the realm of 3D computer games. By integrating style transfer techniques into computer games, developers could dynamically alter the visual aesthetics of a game in real-time, and players could be given the ability to choose from an array of artistic styles, influencing the appearance of the game's world and characters according to their preferences.

However, there is limited work applying NST to 3D computer games. Whilst NST approaches have not been specifically tailored for 3D computer games, image and video NST methods can be applied at the end of the 3D computer graphics

pipeline, as a post-processing effect [1]. This essentially treats the data as a sequence of images. Here, temporal consistency across consecutive frames is the prominent challenge. Some video NST approaches utilise optical flow information and introduce a temporal consistency loss to achieve temporal stability [2, 3, 4, 5, 6], whilst other approaches rely on improving the stability of the transferred content and style features [7, 8, 9]. Nevertheless, employing such models at the post-process stage of the 3D computer graphics pipeline, results in undesired flickering effects and inconsistent stylisations.

Previous work has demonstrated that the utilisation of G-buffer data can lead to improved quality of generated stylised game scenes [10, 11]. Our work takes advantage of the intermediate data that is generated by a 3D computer graphics rendering pipeline, and proposes an approach for integrating an NST model at an earlier stage of the rendering process (Figure 1), resulting in improved, more stable artistic stylisations of game worlds. Our method retrieves data from the camera colour buffer, generates consistent stylised game frames, and writes back to the colour buffer, before post-processing. We believe this is the first work to stylise in this way. The primary contributions of our work can be summarised as follows: 1) We train a fast Stylisation network on both a real-world and a synthetic image dataset, capable of producing fast high-quality artistic stylisations; 2) We present an approach that integrates a trained stylisation network at an early stage of the rendering pipeline, avoiding the visual artefacts and inconsistencies that occur when employing stylisation as a post-effect; 3) We evaluate the results of our system qualitatively and quantitatively, demonstrating how the games community can benefit from the NST field.

## 2   Related Work

### 2.1   Image & Video NST

Gatys *et al.* [12] proposed a model that minimises a content loss and a style loss, based on features extracted from pre-trained CNNs on object classification. Since this seminal work, multiple NST approaches have emerged, proposing end-to-end systems trained on singular styles that manage to improve upon the efficiency and time required to generate one stylisation. These are capable of synthesising a stylised output with a single forward pass through the network [13, 14, 15, 16]. While some models trained to capture the style of a particular artist or a specific art genre [17, 18], and more efficient multiple-style-per-model approaches were also developed [19, 20, 21], recently, the research has shifted to developing arbitrary style transfer systems. The method of Huang and Belongie [22] suggested the use of an Adaptive Instance Normalisation layer (AdaIN) that allows transferring the channel-wise mean and variance feature statistics between the content and style feature activations, thus achieving arbitrary style transfer. Other arbitrary-style-per-model methods were also developed that improve upon the performance [23, 24, 25, 26, 27, 28, 29] or offer solutions tailored to particular challenges [30, 31]. Meta networks were also employed [32], as well as systems making use of the recently developed transformer architecture [7, 33, 34].

To alleviate the issue of temporal inconsistency across subsequent frames when video is considered, Ruder *et al.* [35, 3] employed a temporal constraint based on optical flow information. Typically, the optical flow map (calculated between two frames of the original video) is used to warp the previous stylised frame to give an estimation of the next frame. This gives a temporal loss function that can be minimised during training. Other work has subsequently improved the computation speed [2, 4] or demonstrated structure and depth-preserving qualities [36, 37]. Gao *et al.* [6] developed a fast model that incorporates multiple styles, while arbitrary video style transfer models have also been proposed [38, 8, 9, 39], some of which are extensions from image NST approaches with additional temporal considerations [40, 7, 36].

### 2.2   NST in 3D Computer Games

Although methods exist for stylising three-dimensional data and can offer 3D artists diverse options for generating or improving a game's assets, no substantial efforts have been noted for real-time in-game artistic stylisation. The image and video NST models [7, 8, 9, 6] can potentially be integrated at the end of a computer game's rendering pipeline, intercepting each rendered frame and producing a stylised version of it. An example of this has been exhibited by Unity's implementation [41] which is based on the method of Ghiasi *et al.* [23] that produces a stylised image from an input image in a single forward pass from the neural network. Multi-style in-game style transfer is achieved allowing the viewer to change the stylisation of the scene in real-time. Nonetheless, the implementation does not consider any G-buffer or 3D information. Instead, it intercepts the final rendered 2D image (using an off-screen buffer), which means it can be applied as a final 'filter' for any game. This also results in unstable stylisations and causes the intended post-process effects being diminished.

The recent approach by Richter *et al.* [10] to enhancing the photorealism of computer-generated imagery might be the first to take into account information from intermediate buffers (G-buffers) that becomes available through a game

engine's rendering process. This method, although explicitly focused on photorealistic enhancement, can be significantly impactful to style transfer algorithms that consider the stylisation of game environments. Their technique also works at the end of the rendering pipeline – the image enhancement network outputs an enhanced image given an input rendered image. However, the image enhancement network is fed with information about the geometry, materials, and lighting, extracted from intermediate rendering buffers during training.

Similarly, the image-to-image translation method proposed by Mittermueller *et al.* [11] trains a network to learn the mapping between low-poly game scenes to a synthetic dataset compiled using the Red Dead Redemption 2 (RDR2) game. The mapping takes into account intermediate data such as depth, normals, and albedo generated by conventional game rendering pipelines, for improved image domain transfer. Although the developed *EST-GAN* validates the effectiveness of G-buffer data for the generation of stylistic game scenes, it does not utilise this information in real-time and does not demonstrate any impact on the stability of sequential stylised game frames.

Integrating NST at the end of the 3D rendering pipeline is the only approach that has been suggested for the synthesis of real-time photorealistic [10] or artistic [41] game worlds. Nevertheless, the amount of post-processing that is executed, and the unpredictable camera movement and scene shifts, do not allow for coherent and robust post-process stylisations. Here we propose an approach for producing stable and aesthetically pleasing visual effects in computer games by integrating a style transfer model before the post-process stage of the 3D computer game rendering pipeline.

## 3 Injecting NST into the 3D rendering pipeline

### 3.1 Style Transfer Network

The network architecture is shown in Figure 2. Similarly to state-of-the-art methods [13, 42, 43], we utilise a Transformation network $f_W$, that intercepts an input image $x$ and transforms it into an output image $\hat{y}$ via the mapping $\hat{y} = f_W(x)$. To improve upon the efficiency and inference time required to generate a stylised frame given an input image, we reduce the number of residual layers and remove the ReLU activation function from the first three convolutional layers. The final configuration of our network consists of three convolutional layers followed by instance normalisation, two residual layers (composed of convolutions, instance normalisation, and ReLU), and three deconvolutional layers that upsample the input and then perform convolution. The first two deconvolutional layers are followed by instance normalisation and ReLU activation.
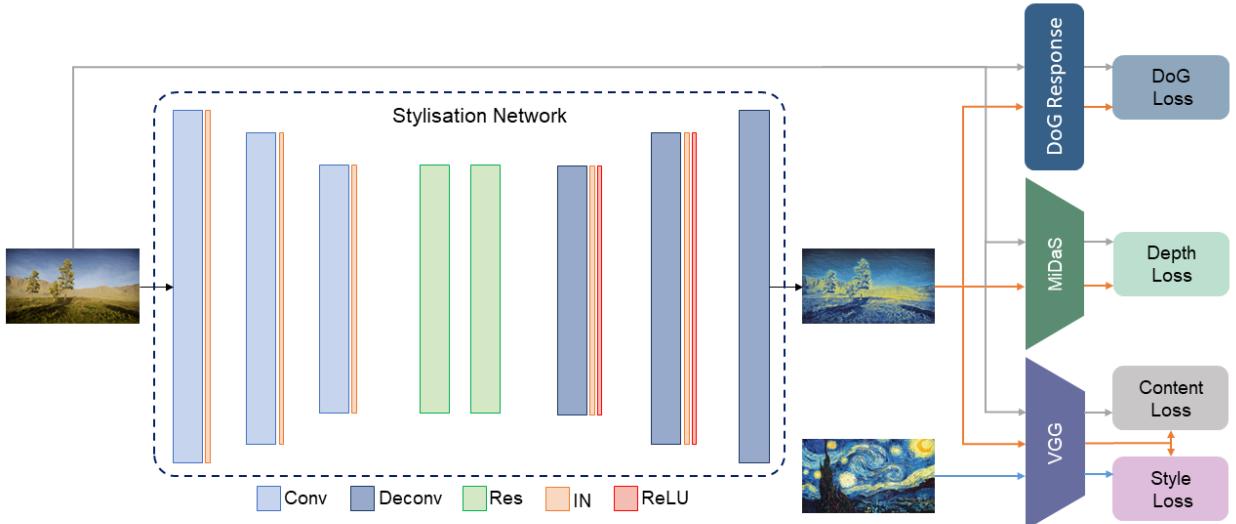


Figure 2: The Stylisation Network consists of three convolutional layers (Conv), two residual layers (Res) and three deconvolutional layers (Deconv). Instance normalisation layers (IN) and the ReLU activation function are included at the first two deconvolution layers.

### 3.1.1 Content & Style Losses

We use the perceptual loss functions introduced in the work of Johnson *et al.* [13] and employ a pre-trained image recognition network (*VGG-16* [44]) to produce feature representations of the original and transformed images. Content loss is defined as the Euclidean distance between the feature representations of the input image and the corresponding transformed image, as extracted from the $relu2\_2$ layer:

$$l_{content}^{\phi_0}(\hat{y}, x) = \frac{1}{C_j H_j W_j} \|\phi_0^j(\hat{y}) - \phi_0^j(x)\|_2^2 \tag{1}$$

where $\phi_0$ is the image classification network, $\phi_0^j$ represents the activations of the $j^{th}$ layer of $\phi_0$, and $H \times W \times C$ is the shape of the processed image.

The style is represented by features extracted from multiple layers of *VGG-16* (*J = {relu1_2, relu2_2, relu3_3, relu4_3}*). The Gram matrix $G$ is then computed to give feature correlations that can be utilised to define the style loss function. This is then defined as the squared Frobenius norm of the difference between the calculated Gram-based style representations:

$$\mathcal{L}_{style}^{\phi_0,j}(\hat{y}, y) = \|G_j^{\phi_0}(\hat{y}) - G_j^{\phi_0}(y)\|_F^2 \tag{2}$$

and it is summed up for all the layers $j$ in $J$. Here, $y$ and $\hat{y}$ refer to the original style image and the transformed image, respectively.

### 3.1.2 Depth Loss

Previous approaches that consider depth information during training [42, 43, 37] have shown improvements to the synthesised results in terms of structure retainment and depth preservation performance. As the trained stylisation network is required to be used in a game setting – and it is highly desired to sustain the depth of the stylised game frames – we utilise a depth reconstruction network (MiDaS) [45] to define a depth reconstruction loss [43, 37]:

$$\mathcal{L}_{depth}^{MiDaS}(\hat{y}, x) = \|MiDaS_1(\hat{y}) - MiDaS_1(x)\|_2^2 \tag{3}$$

### 3.1.3 Difference of Gaussians Loss

A particular issue that occurs in stylisation approaches is an undesired halo effect around distinct parts of an image. This effect is compounded by the significance that is placed on edges in human vision [46, 47] meaning that edge inconsistencies stand out. We use the Difference-of-Gaussians (DoG) operator in order to improve upon the global and local structure preservation of stylised image frames, and thus attempt to alleviate the issue of the undesired halo effect.

Inspired by the neural processing in the retina of the human eye, the DoG response is equivalent to a band-pass filter that discards most of the spatial frequencies that are present in an image. The DoG operator is derived from the concept of convolving an image with two Gaussian kernels of different standard deviations and then taking the difference between the two convolved images. This feature enhancement algorithm has been shown to produce aesthetic edge lines and has been previously utilised for image stylisation [48]. We, therefore, define a DoG loss that is based on the difference between the DoG responses (DoGR) of the original image $x$ and the corresponding stylised image $\hat{y}$:

$$\mathcal{L}_{DoG}(\hat{y}, x) = \|DoGR(\hat{y}) - DoGR(x)\|_2^2 \tag{4}$$

### 3.1.4 Training Details

The Stylisation Network is trained for 2 epochs with a batch size of 2 and a learning rate of $1 \times 10^{-3}$. The content and style weights are set to $1 \times 10^5$ and $1 \times 10^{10}$, respectively. The weight for the depth loss and the DoG loss is set to $1 \times 10^3$. The Adam optimizer [49] is employed with a learning rate of $1 \times 10^{-3}$. The setting of the hyperparameters is adopted from [43] – this maintains the optimal content-style ratio as in the implementation of Johnson *et al.* [13]. To accommodate robust stylisation of game environments and synthetic scenes, both real-world images and frames from computer-generated sources are used to train the stylisation network. The MS COCO dataset [50] is used, mixed with frames from the MPI Sintel training set [51]. The data is shuffled and all the images are resized to $360 \times 360$ during training. In order for the trained Stylisation network to be suitable for in-game stylisation, we export the trained model to the ONNX [52] format. This is supported by Unity and the Barracuda package [1].
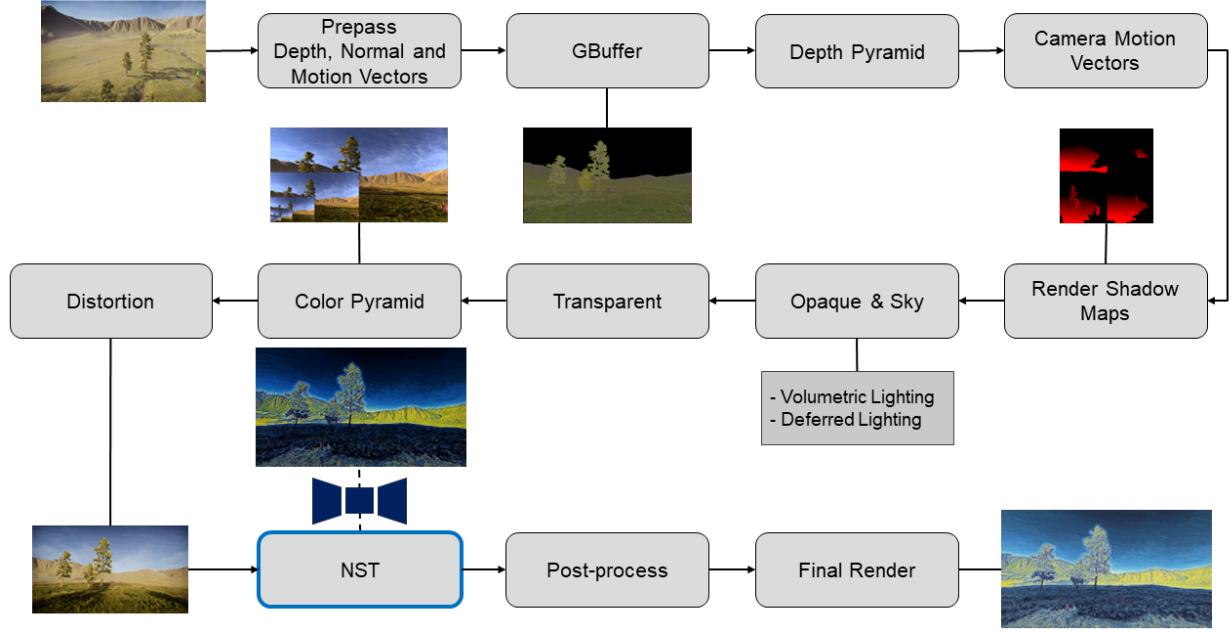
Figure 3: Overview of the modified 3D rendering pipeline. The NST module is added before the Post-process stage whilst the colour buffer information is available for read and write.

### 3.2 In-game Stylisation

To accommodate real-time in-game stylisation we use the Unity game engine and the High Definition Rendering Pipeline (HDRP) [53]. Custom Passes can be configured within Unity's rendering pipeline and can be executed at certain points during the HDRP render loop. Six injection points for a Custom Pass are offered, with a selection of buffers being available at each. The injection points are: *Before Rendering*, *After Opaque, Depth and Normal*, *Before Pre-Refraction*, *Before Transparent*, *Before Post-Process*, and *After Post Process*. To generate a stylised image, it is necessary to read and write to the colour buffer that is available after the *Opaque, Depth and Normal* stage. In order for the stylisation to affect the transparent objects in the scene, we opt to inject the custom pass before the Post-Process stage.

The overall modified Unity HDRP rendering pipeline is depicted in Figure 3. During rendering, HDRP writes colour data from visible objects (renderers) in the scene to the colour buffer. During a custom pass, a depth pyramid and a colour pyramid are created (as shown in Figure 3). The colour pyramid constitutes an iterative series of mipmaps, crafted by the HDRP, extracted from the colour buffer at a specific juncture within the rendering pipeline. The NST Module is inserted after the Distortion stage and before the Post-process stage, intercepting the colour buffer mipmap and producing an artistic stylisation for each frame (this is supported by the Barracuda package [1] that allows for neural network inference). The synthesised texture is then passed to a custom compute shader that writes the colour to the camera colour buffer. This allows for the Post-process stage to utilise the stylised frames, before the final render. Our proposed system is capable of producing stable real-time stylised frames free from undesired artefacts and flickering effects. Embedding the NST module earlier in the rendering pipeline also allows for the post-process effects (such as depth of field, bloom, or motion blur) to effectively be visible, adding to the look and feel of the game. Such effects would be diminished if the stylisation effect was applied at the final render – examples of this are shown in Figure 5.

## 4 Results & Discussion

Our NST system is embedded in the rendering pipeline, intercepting each G-buffer colour frame and producing a stylised version that is then passed through the Post-process stage. Figure 4 demonstrates some final rendered frames for different reference style images and for different game scenes. Our method is capable of producing robust stylisations even for complicated scenes with difficult lighting. The final renders do not suffer from undesired artefacts or flickering

effects, while the halo effect around the objects in the scene is significantly reduced. The following subsections further demonstrate temporally consistent stylisations of frames from various open-source games [54, 55, 56, 57] and compare the results against state-of-the-art methods in image and video style transfer. Videos of our results and comparisons to state-of-the-art methods are included on the project's website.



Figure 4: Our approach for different style images and different game scenes. Original content images are above stylised frames. Adjacent frames show temporal stability.

## 4.1 Qualitative Results

Qualitative comparisons against four state-of-the-art methods – AdaAttN [7], CSBNet [9], MCCNet [8], and FVMST [6] – are shown in Figure 5. This includes stylisations for two consecutive frames for two different game scenes and two different style images.

Figure 5 shows that AdaAttN [7] preserves much of the content information, however, the stylisation effect is not very visible – the yellow colour that is eminent in the style image is absent from the stylised frames (Figure 5(b)). The video NST methods, CSBNet [9] and MCCNet [8], reproduce the style image more faithfully than AdaAttN, but they create undesired artefacts such as the yellow halo around the trees that are visually distinguished from the background. FVMST [6] also captures the style quite well, but generates a white artefact encircling the mountain's background edge and it produces a sudden shift of the sky colour (from light, it turns to dark blue/black) that is not visible in the original frames. Our approach reproduces the style image faithfully and eliminates the undesired effects that are visible in the results of the state-of-the-art methods. The structure of the original frames is preserved comparably to the stylisations of AdaAttN but with higher stylisation intensity and better preservation of the luminance and lighting of the scene.

To demonstrate the effectiveness of our approach, we also include close-ups of consecutive frames of a game scene that includes a moving 3D object in a complex background. When looking at the zoomed-in cut-outs in Figure 5(a), the halo effect around the object's edge is more noticeable in the generated frames of AdaAttN [7], CSBNet [9], and MCCNet [8], while the stylisations of FVMST [6] create a disturbing white blob. In addition, the close-ups provide strong evidence of the capability of our approach to retaining the luminance in the scene and the game's post-effects. The prominent depth-of-field effect in the original frames is completely ignored when the stylisation is performed at the final render using state-of-the-art methods that enhance the details on the background. Our system makes the 3D object stand out and preserves the lighting and the game's overall look and feel.

## 4.2 Quantitative Results

For quantitative comparisons, frames are extracted from 4 different games and 12 different gameplays, including indoor and outdoor scenes, featuring moving objects and complicated lighting. This results in an evaluation dataset of 2100
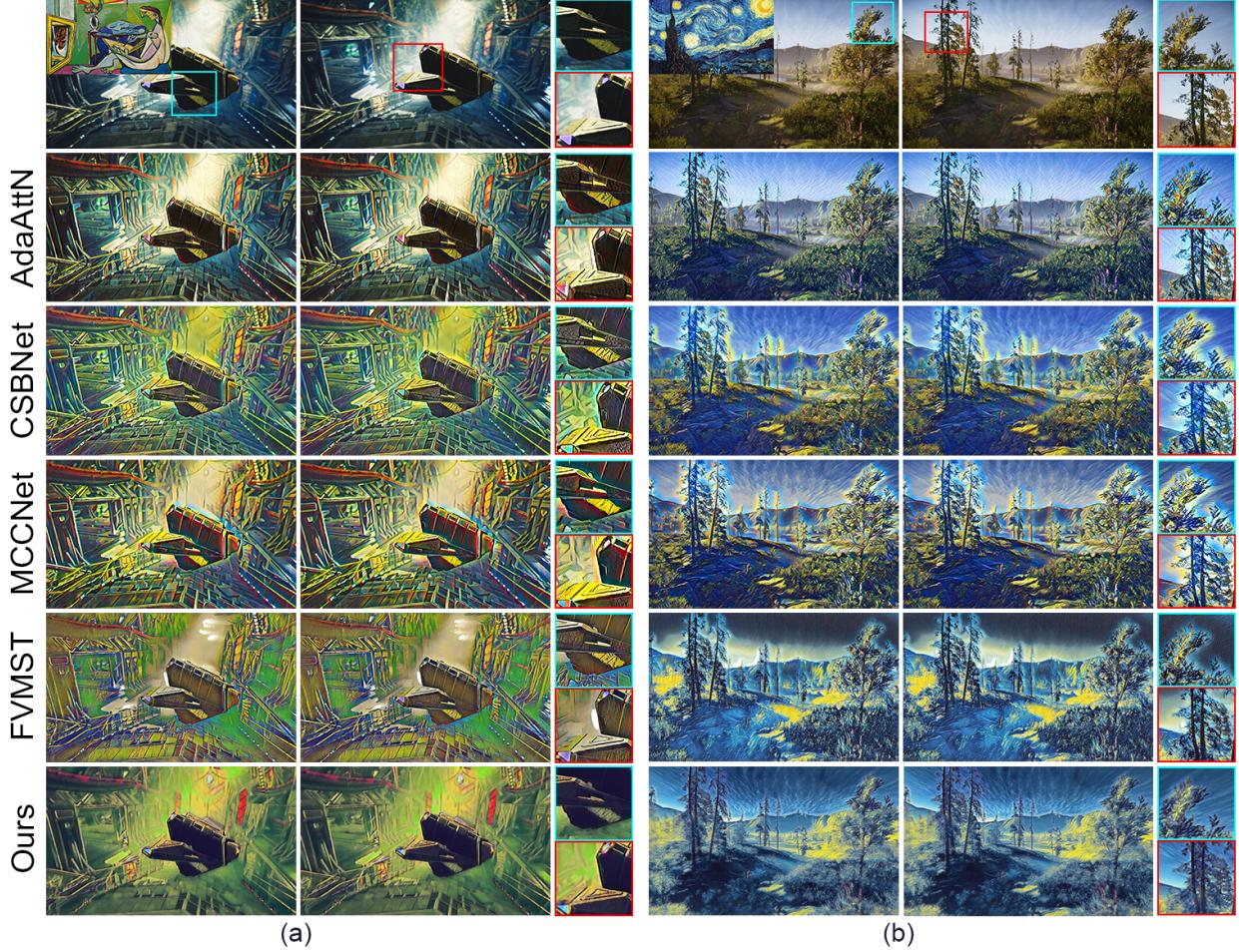
Figure 5: Comparison against state-of-the-art approaches. Top row: original frames, with the style image top left; two adjacent frames from two different game scenes are shown. We provide zoomed-in cut-outs on the right of each two-frame sequence, for better comparisons. Our method produces robust stylisations that capture the style image more efficiently and preserve content and luminance information of the scene more effectively in comparison with the state-of-the-art approaches.

frames (9 gameplays $\times$ 200 frames and 3 gameplays $\times$ 100 frames), and we evaluate utilising 10 different style images. The average results are reported in Table 1.

To quantitatively gauge the performance of our method in video stability and temporal coherence we utilise the warping error that is calculated as the difference between a warped next frame (using optic flow) and the original next frame. *FlowNetS* [58] is used to compute the optical flow of the original videos. In addition, we employ the LPIPS (Learned Perceptual Image Patch Similarity) metric [59] to measure the average perceptual distances between the adjacent frames in order to verify the smoothness of the stylised game sequences. The results show that our approach is superior to the state-of-the-art methods in generating temporally consistent in-game stylisations.

Perceptual metrics are employed to quantitatively assess the stylisation quality. SSIM [60] and Content error ($\mathcal{L}_c$) [12] are used to evaluate the effectiveness of the methods in retaining content information; SIFID [61] and Style error [12] are used to evaluate style performance. Our system manages to preserve content adequately. Whilst our algorithm's effectiveness in reproducing the style image is sufficient, some stylisation qualities are lost when the post-process effects are performed on top of the stylisations. In order to retain the intended post-effects applied to a game, certain aspects of the style likeness to the original image are traded off. Arguably, this compromise can be deemed desirable in a game setting and, as has been demonstrated, this trade-off leads to more consistent and temporally stable stylisations.

7

| Method | Warping Error ↓ | LPIPS Error ↓ | SSIM ↑ | SIFID ↓ | $\mathcal{L}_c$ ↓ | $\mathcal{L}_s$ ↓ |
|---|---|---|---|---|---|---|
| AdaAttN [7] | 1.6477 | 0.3217 | **0.7820** | 1.6115 | **0.4945** | 1.0391 |
| CSBNet [9] | 1.7458 | 0.3908 | 0.6370 | 2.2468 | 0.8674 | 1.0053 |
| MCCNet [8] | 1.6519 | 0.3547 | 0.6637 | 1.5555 | 0.8065 | 1.0042 |
| FVMST [6] | 1.8524 | 0.3215 | 0.5855 | 2.2529 | 0.7834 | 1.0077 |
| Ours (image) | 1.6764 | 0.3602 | 0.6740 | **1.2063** | 0.6532 | **0.9808** |
| Ours (in-game) | **1.5798** | **0.2930** | 0.6057 | 1.8679 | 0.7830 | 1.0612 |

Table 1: Quantitative results. Warping Error and LPIPS error (both in the form $\times 10$) capture the smoothness of the generated video. SSIM and $\mathcal{L}_c$ relate to content preservation, and SIFID and $\mathcal{L}_s$ quantify the style performance. Results are given for our NST system injected in the game's rendering pipeline (game) and for the NST network applied as a post-effect (image). Bold values in Warping Error and LPIPS Error indicate our (in-game) approach is best at preserving temporal consistency.
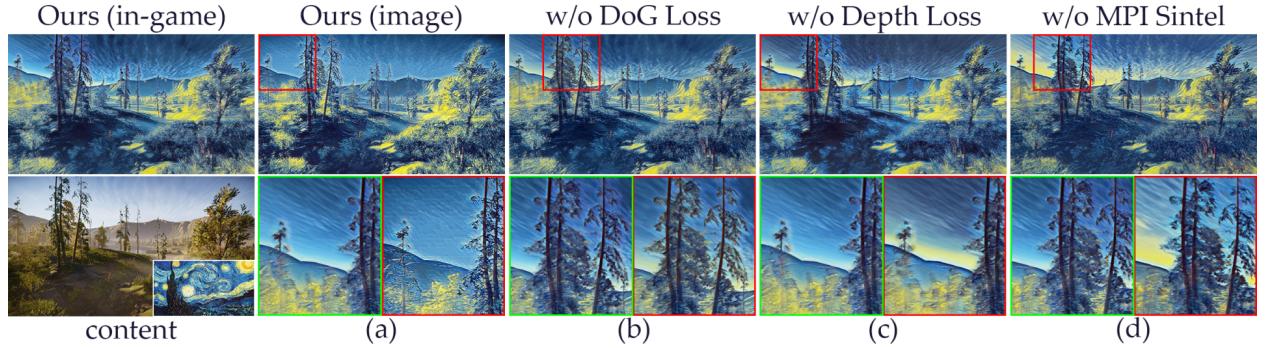
## 4.3 Ablation Study



Figure 6: Ablation study on the impact of the different components of our system. Each column shows a zoomed-in comparison between our method (in-game) trained with all components (green) and our stylisation network (a) applied as a post-effect, (b) trained without DoG Loss (c) trained without Depth Loss, and (d) trained without the MPI Sintel data (red).

Figure 6 demonstrates example results of our approach under different configurations. In-game stylisation has a significant impact on temporal coherence, in comparison to stylising each rendered frame as a post-effect (Section 4.2). Here, we show that the latter also produces less appealing stylisations, with much of the content information being discarded. In addition, training without the DoG Loss results in more visible halos around the objects in the scene, whereas training with DoG Loss leads to generated frames with enhanced object stylisation and reduced boundary artefacts. The same applies to Depth Loss, as our method synthesises visibly improved results when depth is considered. The inclusion of the MPI Sintel dataset also has an impact on the performance – the stylisation network trained only on the MS COCO dataset neglects the synthetic nature of the game and struggles to generate frames that retain the content adequately, producing undesired effects.

## 4.4 Limitations

To demonstrate the effectiveness of applying NST as part of the rendering pipeline of a computer game, we have trained a single-style-per-network model. Future work could experiment with arbitrary-style-per-model networks [22, 7, 8] which would provide the user with the option to upload and use their own reference style image. Another important consideration in applying NST in a game setting is running time. We reduced the number of residual layers and removed activation from the initial convolution layers to improve upon the inference time of the trained network which requires approximately 0.9 seconds to stylise an image of size $512 \times 512$. When injecting stylisation in the rendering pipeline the frame rate of a game running in Unity at Full HD resolution drops to $\sim$10fps. Utilising a more lightweight network architecture (arbitrary style transfer networks report better inference time, e.g., AdaIN [22]: 0.065 seconds) could result in stylised game environments running at higher frame rates.

## 5    Conclusion

We have proposed a novel approach for injecting NST into a computer graphics rendering pipeline. Our NST framework is capable of producing coherent and temporally stable stylised frames in computer games. Our NST module intercepts frames from the colour buffer and synthesises artistic stylisations that are then written back to the camera colour buffer. Robust stylisations are achieved without interfering with the applied post-process effects. We demonstrate qualitative and quantitative results that reveal a promising new avenue for integrating NST within game development processes.

## Acknowledgements

## References

[1] Unity Technologies. Introduction to barracuda: Barracuda: 3.0.1, 2023.

[2] Haozhi Huang, Hao Wang, Wenhan Luo, Lin Ma, Wenhao Jiang, Xiaolong Zhu, Zhifeng Li, and Wei Liu. Real-time neural style transfer for videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 783–791, 2017.

[3] Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. Artistic style transfer for videos and spherical images. *CoRR*, 2017.

[4] Chang Gao, Derun Gu, Fangjun Zhang, and Yizhou Yu. Reconet: Real-time coherent video style transfer network, 2018.

[5] Tian Qi Chen and Mark Schmidt. Fast patch-based style transfer of arbitrary style. *arXiv preprint arXiv:1612.04337*, 2016.

[6] Wei Gao, Yijun Li, Yihang Yin, and Ming-Hsuan Yang. Fast video multi-style transfer. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3222–3230, 2020.

[7] Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Meiling Wang, Xin Li, Zhengxing Sun, Qian Li, and Errui Ding. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6649–6658, 2021.

[8] Yingying Deng, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, and Changsheng Xu. Arbitrary video style transfer via multi-channel correlation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2):1210–1217, May 2021.

[9] Haofei Lu and Zhizhong Wang. Universal video style transfer via crystallization, separation, and blending. In *Proc. Int. Joint Conf. on Artif. Intell.(IJCAI)*, volume 36, pages 4957–4965, 2022.

[10] Stephan R Richter, Hassan Abu AlHaija, and Vladlen Koltun. Enhancing photorealism enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):1700–1715, 2022.

[11] Martina Mittermueller, Zhanxiang Ye, and Helmut Hlavacs. EST-GAN: Enhancing style transfer gans with intermediate game render passes. In *2022 IEEE Conference on Games (CoG)*, pages 25–32, 2022.

[12] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016.

[13] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.

[14] Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, and Victor S Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *ICML*, volume 1, page 4, 2016.

[15] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6924–6932, 2017.

[16] Chuan Li and Michael Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*, pages 702–716. Springer, 2016.

[17] Artsiom Sanakoyeu, Dmytro Kotovenko, Sabine Lang, and Bjorn Ommer. A style-aware content loss for real-time hd style transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 698–714, 2018.

[18] Dmytro Kotovenko, Artsiom Sanakoyeu, Sabine Lang, and Bjorn Ommer. Content and style disentanglement for artistic style transfer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4422–4431, 2019.

[19] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. In *International Conference on Learning Representations*, 2017.

[20] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1897–1906, 2017.

[21] Hang Zhang and Kristin Dana. Multi-style generative network for real-time transfer. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.

[22] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017.

[23] Golnaz Ghiasi, Honglak Lee, Manjunath Kudlur, Vincent Dumoulin, and Jonathon Shlens. Exploring the structure of a real-time, arbitrary neural artistic stylization network. *arXiv preprint arXiv:1705.06830*, 2017.

[24] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. *Advances in neural information processing systems*, 30, 2017.

[25] Zheng Xu, Michael Wilber, Chen Fang, Aaron Hertzmann, and Hailin Jin. Learning from multi-domain artistic images for arbitrary style transfer. *arXiv preprint arXiv:1805.09987*, 2018.

[26] Jing Huo, Shiyin Jin, Wenbin Li, Jing Wu, Yu-Kun Lai, Yinghuan Shi, and Yang Gao. Manifold alignment for semantically aligned style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14861–14869, 2021.

[27] Dae Young Park and Kwang Hee Lee. Arbitrary style transfer with style-attentional networks. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5880–5888, 2019.

[28] Jan Svoboda, Asha Anoosheh, Christian Osendorfer, and Jonathan Masci. Two-stage peer-regularized feature recombination for arbitrary image style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13816–13825, 2020.

[29] Jie An, Haoyi Xiong, Jun Huan, and Jiebo Luo. Ultrafast photorealistic style transfer via neural architecture search. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):10443–10450, Apr. 2020.

[30] Zhiyuan Hu, Jia Jia, Bei Liu, Yaohua Bu, and Jianlong Fu. Aesthetic-aware image style transfer. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 3320–3329, 2020.

[31] Xiao-Chang Liu, Yong-Liang Yang, and Peter Hall. Learning to warp for style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3702–3711, 2021.

[32] Falong Shen, Shuicheng Yan, and Gang Zeng. Neural style transfer via meta networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8061–8069, 2018.

[33] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11326–11336, 2022.

[34] Xuan Luo, Zhen Han, Lingkang Yang, and Lingling Zhang. Consistent style transfer. *arXiv preprint arXiv:2201.02233*, 2022.

[35] Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. Artistic style transfer for videos. In Bodo Rosenhahn and Bjoern Andres, editors, *Pattern Recognition*, pages 26–36, Cham, 2016. Springer International Publishing.

[36] Shiguang Liu and Ting Zhu. Structure-guided arbitrary style transfer for artistic image and video. *IEEE Transactions on Multimedia*, 2021.

[37] Eleftherios Ioannou and Steve Maddock. Depth-aware neural style transfer for videos. *Computers*, 12(4):69, 2023.

[38] Wenjing Wang, Shuai Yang, Jizheng Xu, and Jiaying Liu. Consistent video style transfer via relaxation and regularization. *IEEE Transactions on Image Processing*, 29:9125–9139, 2020.

[39] Zijie Wu, Zhen Zhu, Junping Du, and Xiang Bai. Ccpl: Contrastive coherence preserving loss for versatile style transfer. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVI*, pages 189–206. Springer, 2022.

[40] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast image and video style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3809–3817, 2019.

[41] Thomas Deliot, Florent Guinier, and Kenneth Vanhoey. Real-time style transfer in unity using deep neural networks, 2020.

[42] Xiao-Chang Liu, Ming-Ming Cheng, Yu-Kun Lai, and Paul L Rosin. Depth-aware neural style transfer. In *Proceedings of the Symposium on Non-Photorealistic Animation and Rendering*, pages 1–10, 2017.

[43] Eleftherios Ioannou and Steve Maddock. Depth-aware neural style transfer using instance normalization. In *Computer Graphics & Visual Computing (CGVC) 2022*. Eurographics Digital Library, 2022.

[44] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015.

[45] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.

[46] Stephen E Palmer. *Vision science: Photons to phenomenology*. MIT press, 1999.

[47] David Marr and Ellen Hildreth. Theory of edge detection. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 207(1167):187–217, 1980.

[48] Holger Winnemöller, Jan Eric Kyprianidis, and Sven C Olsen. Xdog: An extended difference-of-gaussians compendium including advanced image stylization. *Computers & Graphics*, 36(6):740–753, 2012.

[49] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[50] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[51] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), editor, *European Conf. on Computer Vision (ECCV)*, Part IV, LNCS 7577, pages 611–625. Springer-Verlag, October 2012.

[52] ONNX. Open neural network exchange, 2019.

[53] Unity Technologies. High definition render pipeline overview: High definition rp: 12.1.12, 2021.

[54] Unity Technologies. Unity terrain - hdrp demo scene, 2022.

[55] Unity Technologies. Unity-technologies/fontainebleaudemo: Fontainebleau demo, 2022.

[56] POLYGONAUTIC. Seed hunter, 2020.

[57] Unity Technologies. Book of the dead: Environment: Hdrp: Tutorial projects, 2023.

[58] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2462–2470, 2017.

[59] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.

[60] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.

[61] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. SinGAN: Learning a generative model from a single natural image. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4570–4580, 2019.