

# GPU-Accelerated Mobile Multi-view Style Transfer

Puneet Kohli, Saravana Gunaseelan, Jason Orozco, Yiwen Hua, Edward Li, and Nicolas Dahlquist

Leia Inc

## I. ABSTRACT

An estimated 60% of smartphones sold in 2018 were equipped with multiple rear cameras, enabling a wide variety of 3D-enabled applications such as 3D Photos. The success of 3D Photo platforms (Facebook 3D Photo, Holopix™, etc) depend on a steady influx of user generated content. These platforms must provide simple image manipulation tools to facilitate content creation, akin to traditional photo platforms. Artistic neural style transfer, propelled by recent advancements in GPU technology, is one such tool for enhancing traditional photos. However, naively extrapolating single-view neural style transfer to the multi-view scenario produces visually inconsistent results and is prohibitively slow on mobile devices. We present a GPU-accelerated multi-view style transfer pipeline which enforces style consistency between views with on-demand performance on mobile platforms. Our pipeline is modular and creates high quality depth and parallax effects from a stereoscopic image pair.

## II. INTRODUCTION

Consumers are actively seeking ways to make their traditional photography more immersive and life-like. With the rapidly increasing adoption of multi-camera smartphones, a variety of immersive applications that use depth information are now enabled at mass market scale. One of the most promising applications in this domain is *3D Photos* – which capture multiple viewpoints of a scene with a parallax effect [1]. Coupled with a Lightfield display such as those made by Leia Inc [2], 3D Photos are brought to life with added depth [3] and multi-view parallax effects [4]. Platforms such as Facebook 3D Photos and Holopix™<sup>1</sup> are making it easier for consumers to now capture, edit, and share 3D Photos on a single mobile device. Paramount to the continued success of 3D Photo platforms is the availability of image manipulation and content generation tools that consumers are accustomed to from traditional photo platforms.

To create content for 3D Photo platforms, an image containing information of multiple viewpoints (a *multi-view image*) must be generated. With advances in modern algorithms such as view synthesis [5]–[9] and in-painting [10]–[12], it is now possible to generate multi-view images from stereo pairs even on mobile devices for content creation and manipulation. However, challenges with multi-view content creation still

remain. The algorithms involved require large amounts of processing power to perform operations for multiple viewpoints along with higher storage and bandwidth requirements which scales poorly with the number of rendered views. Additionally, mobile platforms have limited computing power to support on-demand (~0.5 seconds) performance. To enable content creation at scale these challenges must be addressed along with providing users simple tools to generate high quality content.

Artistic style transfers using neural networks, pioneered by Gatys *et. al* [13], [14], produce visually pleasing content that has been the topic of research of many recent works [15]–[18] and seen widespread success in consumer-facing applications such as Microsoft Pix [19] and Prisma [20]. Given the commercial success of neural style transfers as a simple way to make even simple images look artistic, we aim to bring the effect to 3D Photo platforms. On a Lightfield display, seeing the artistic effect come to life due to the parallax and depth effect combined is a truly unique experience. An additional and hidden benefit of applying style transfer to multi-view images is that any artifacts such as edge inconsistency or poor in-painting is generally unnoticeable once an artistic style is applied. Though neural style transfer techniques work well on single images, naively applying such techniques on individual views of a multi-view image yields inconsistently styled images which results in 3D fatigue when viewed on a multi-view display [21].

Recent works have aimed to solve the consistency between multiple views for artistic style transfer [22], [23]. These approaches only consider two views and are not directly extendable to multiple views. They are also optimized for specific styles and need to be re-trained for new styles. Video style transfer [24]–[26] is another similar line of work where consistent results need to be produced in the temporal domain. These techniques are currently not computationally efficient enough for running in real-time on mobile devices. In our case, we aim to solve a generalized scenario where multiple views have an artistic style applied, generated from input stereo image and disparity pairs. Our method is agnostic to both the number of views being generated and the style transfer model being used.



Fig. 1. An example result from running our style transfer pipeline to generate four views using a stereo image and disparity pair.

<sup>1</sup>An image-sharing social network for multi-view images. <https://www.holopix.com>

To the best of our knowledge, there are no techniques that maintain multi-view style consistency while also addressing the performance concerns on mobile devices. In this work, we propose an end-to-end pipeline for generating stylized multi-view imagery given stereo image and disparity pairs that can run on-demand even on mobile computing platforms. We address the issues of both multi-view style consistency as well as performance. An example result from our pipeline is shown in Figure 1.

The proposed method is highly configurable and can support different algorithms for each of the steps involved. This includes style transfer, novel view synthesis [5]–[9], and inpainting [10]–[12]. In a highly modular fashion, each component of the pipeline is independent of each other and can be replaced by an equivalent algorithm. This facilitates fine-tuning the proposed pipeline for different design constraints such as optimizing for power or quality.

In summary, the pipeline proposed in this work has the following main contributions:

- Multi-view Consistent neural style transfer and parallax effects
- GPU-Accelerated on-demand performance ( $\sim 0.5s$ ) on mobile platforms.
- Modular components for individual algorithms.

### III. METHODOLOGY

The proposed method for multi-view consistent stylizing combines existing monoscopic style transfer techniques with view synthesis in a highly modular fashion, while promoting 3D scene integrity. An overview is given in Algorithm 1 and Figure 2 showing the four steps to stylize the left view, re-project to the right viewpoint, apply a guided filter to the left and right stylized views, then synthesize any number of stylized novel views. This section explores how each module contributes to create immersive stylized 3D content.

**Algorithm 1** Multi-view consistent style transfer. Given stereo input views  $I_l, I_r$ , disparity maps  $\Delta_l, \Delta_r$ , and style guide  $G_s$ , render stylized views  $S_{1,2,\dots,n}$  at desired output viewpoints  $x_1, x_2, \dots, x_n$ .

- 1: Infer stylized left view  $S_l$  from  $I_l$  and  $G_s$  using style transfer module.
- 2: Re-project  $S_l$  to  $S_r$  at the same viewpoint as  $I_r$  using view synthesis module.
- 3: Apply guided filter module to  $S_l, S_r$  with guides  $I_l, I_r$  to produce filtered stylized views  $S'_l, S'_r$ .
- 4: **for**  $i$  in  $\{1, 2, \dots, n\}$  **do**
- 5:   Re-project  $S'_l$  and  $S'_r$  to  $x_i$  using view synthesis module and blend into  $S_i$ .
- 6: **end for**

#### A. Input

We assume that left and right views  $I_l, I_r$  are given with corresponding disparity maps  $\Delta_l, \Delta_r$ . In applications where

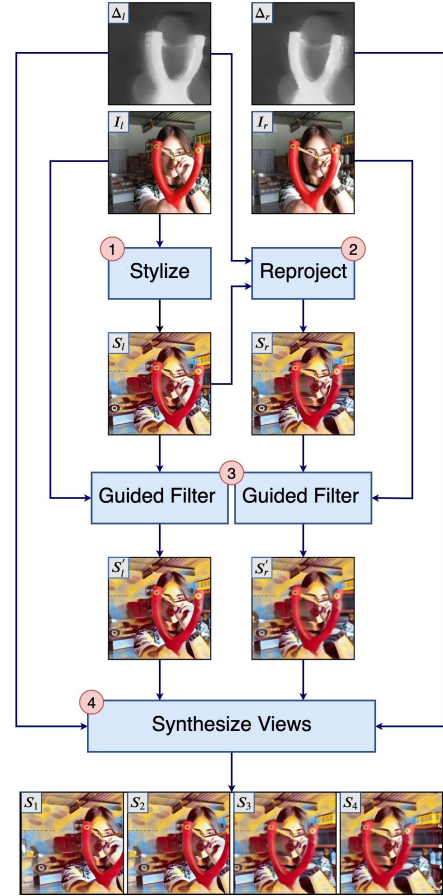


Fig. 2. Algorithm flow of our multi-view style transfer pipeline. Red numbers correspond to steps in Algorithm 1.

disparity maps are not given, a variety of well-known estimation methods can be used [27]–[30]. In our case, we use a neural network trained on stereo images from Holopix™ prior to our experiments.

#### B. Style Transfer

The proposed method stylizes the left view  $S_l$  using existing style transfer methods [15], [31]. We use a network based on [15] along with Instance Normalization [32] in our experiments.

We achieve three main benefits by applying the style transfer network to stylize only one of the input stereo views. The first is to achieve style consistency between multiple views, which is the primary motivation for our proposed method. The challenge is that style transfer networks are sensitive to small changes in the input. When stylizing each view individually, the parallax effect between neighboring views is enough to cause visual differences such as stylized features entirely appearing and disappearing from one view to the next (see Figure 4). Such inconsistencies cause 3D fatigue when viewed on a multi-view display.

Second, running a style transfer network incurs a significant performance cost even with competitive technology, so styliz-



Fig. 3. Comparison of stylized right image before (top row) and after (bottom row) applying guided filter. The highlighted regions show that applying a guided filter emphasizes the edges in the original 3D scene over the edges introduced by style transfer.

ing once and synthesizing many is much faster than stylizing every output view individually.

The third benefit is compatibility with any style transfer network. By not demanding modifications and retraining specifically for multi-view rendering, any network that takes a single image and returns a stylized version can be readily substituted in our pipeline, provided the quality and performance are satisfactory for single images.

### C. One In, One Out View Synthesis

Next, the stylized left view  $S_l$  is re-projected to  $S_r$  at the same viewpoint as the input right view  $I_r$  using a view synthesis module. View synthesis is an active area of research with many existing solutions [5]–[8]. We use an in-house algorithm for performing view synthesis both on the CPU and GPU. Our algorithm incorporates forward warping with a depth test and a simple in-painting technique that samples nearby regions to fill deoccluded regions.

The effect of re-projecting the stylized left view to the right viewpoint is that the style features are transported precisely to their corresponding positions in the right view. This is equivalent to having the style features present in the 3D scene at the time the stereo photo is taken.

### D. Guided Filter

The stylized views are refined by the edge-aware guided filtering methodology to produce filtered stylized left and right views  $S'_l, S'_r$  with the input views  $I_l, I_r$  serving as their respective guides. We use the CPU-based guided filter [33] in all our experiments.

When viewing images on a multi-view display, the quality of edges plays an essential role in 3D perception, however, the style transfer process tends to degrade the edges of objects in the 3D scene. By applying a guided filter to the stylized views using their corresponding un-stylized views as guides, the edges of original 3D objects are reinforced while

reducing the stylization edges, resulting in a more immersive 3D experience. See Figure 3 for comparison.

### E. Two In, Many Out View Synthesis

The final step is to synthesize the desired output views  $S_{1,2,\dots,n}$  from the filtered stylized left and right views  $S'_l, S'_r$  and given disparity maps  $\Delta_l, \Delta_r$ . This is done by repeated application of the same *one in, one out view synthesis* module described in Section III-C to re-project both  $S'_l$  and  $S'_r$  to every output viewpoint. Each output view  $S_i$  is the result of re-projecting both  $S'_l$  and  $S'_r$  to the desired viewpoint of  $S_i$ , say  $x$ , and blending based on proximity of  $x$  to the left and right viewpoints  $l$  and  $r$ .

## IV. EXPERIMENTS AND RESULTS

### A. Experimentation

We tested four approaches for generating multi-view style transferred images and compared both qualitative and quantitative results. We ran all experiments on a RED Hydrogen One [34] mobile phone that has four-view autostereoscopic display. The device runs on the Qualcomm® Snapdragon 835 Mobile Platform coupled with the Qualcomm® Adreno™ 540 GPU.

**Baseline Approach.** The first approach naively applies neural style transfer to each of the synthesized views individually. This fails to produce style consistent views due to the unstable nature of neural style transfer models [22] and also does not produce on-demand results. A supplementary observation shows that style transfer is considerably slower than view synthesis. The overall performance scales poorly with the number of views.

**Approach 2.** Neural style-transfer is first applied to each of the stereoscopic input views and then novel view synthesis is performed using the stylized pair and the original disparity maps as inputs. Although this runs significantly faster than the *baseline* method, the rendered views produce undesirable ghosting artifacts and an overall inconsistent styling between the stereoscopic pair leading to 3D fatigue.

**Approach 3.** Here the style inconsistency problem left unsolved from the previous methods is tackled. Neural style is applied only to the input left image to create the stylized left image. Novel views are synthesized only from this stylized left image. View synthesis is simultaneously performed using both the original naturalistic left and right images. This naturalistic multi-view image is used as a guide for a guided filter pass on the stylized multi-view image. The resulting multi-view image is sharpened to reduce blurring artifacts. This method produces consistently styled views with relatively sharp edges.

However, the drawback of this approach is that there is a limited depth effect due to using only the left image for the styled novel view synthesis. Additionally, as the edges do not line up perfectly in the guided filtering step, ghosting artifacts



are produced around the edges in the output views.

**Selected Approach.** Our novel approach outlined in Section III builds upon ideas from all of the previous approaches and succeeds in producing multi-view consistent stylized images with on-demand performance when GPU-accelerated on mobile devices.

### B. Quantitative Evaluation

Our evaluation metric for comparing the various approaches is the time taken to produce a multi-view image with neural style transfer applied from an input stereoscopic image pair and their associated disparity maps. We run each of the compared methods on the *MPI Sintel* stereo test dataset [35] at a 1024x436 resolution.

TABLE I  
MEAN RUNTIME COMPARISON FOR RENDERING A STYLIZED MULTI-VIEW IMAGE WITH 4, 8, AND 16 VIEWS FROM THE *MPI Sintel* STEREO DATASET.

Method	Time taken (ms)		
	4 Views	8 Views	16 Views
Baseline <sub>CPU</sub>	8352	16682	33358
Baseline <sub>GPU</sub>	1405	2832	5768
Approach <sub>2</sub>	843	849	858
Approach <sub>3</sub>	746	995	1213
Ours <sub>CPU</sub>	2311	2394	2576
<b>Ours<sub>GPU</sub></b>	<b>561</b>	<b>567</b>	<b>576</b>

Table I shows the run time comparison of our various methods on the *MPI Sintel* stereo dataset. The baseline method is run both in an end-to-end CPU and a GPU-accelerated environment to highlight the performance advantages of using the GPU. As can be seen, our novel pipeline is the fastest method and scales linearly with the number of views generated. The style transfer model used in our experiments is from [15].

Additionally, we deployed images with our stylization technique to the Holopix™ platform and compared overall engagement. We observed that images posted to Holopix™ with style transfer applied perform over 300% better than the platform average. In a control post when the same image was posted twice in succession, once with style transfer applied, and once without, the version of the image with style transfer performed over 20% better.

### C. Qualitative Evaluation

Figure 4 shows visual results from the baseline approach vs our final approach for a stereoscopic image pair taken from a stereo camera on the RED Hydrogen One [34] device. When viewing the results on a multi-view display, the baseline methods makes it difficult for the eyes to focus on due

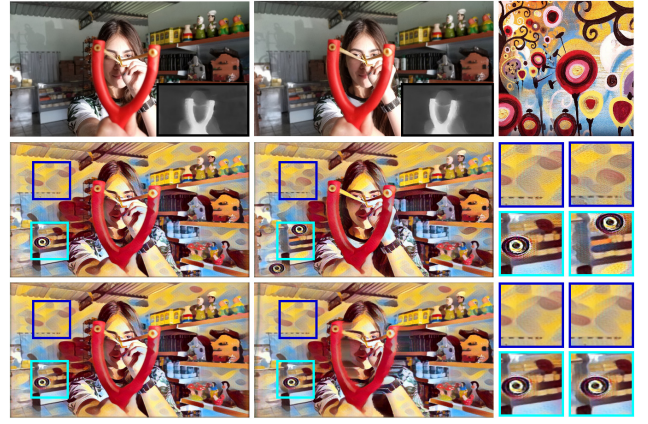


Fig. 4. Top row: original left image and disparity (left), original right image and disparity (center), style guide (right). Comparison of stylized left and right image using baseline approach (middle) and our approach (bottom). Baseline approach produces inconsistent styling between left and right image leading to 3D fatigue and our approach eliminates such inconsistency as seen in the highlighted region.

to inconsistencies causing 3D fatigue. In comparison, our approach has consistent views making the experience more comfortable to view on a multi-view display.

We show additional results in Figure 5 where we generate various images from stereo pairs and a different style. It is interesting to note that between the first view of each result and the last view, objects have consistent stylization despite having different positions, rotations, or occlusions.

## V. CONCLUSION

In this work, we propose an end-to-end pipeline for generating stylized multi-view imagery given stereo image and disparity pairs. Our experiments show that the pipeline can achieve on-demand (~0.5 seconds) results on mobile GPUs, tested up to sixteen views. The results are style consistent between each of the views and produce a highly immersive parallax effect when viewed on a lightfield display. Finally, the pipeline is modular and can support any style transfer, view synthesis, or in-painting algorithm. To the best of our knowledge, this is the first published work for producing style consistent multi-view imagery from a stereoscopic image and disparity pair with on-demand performance on mobile GPUs. In a future work, we will run our pipeline on the Jetson TX2 GPU and the Tegra X1 GPU on an NVIDIA SHIELD and we expect to see similar or better quantitative results.

## REFERENCES

- [1] J. Kopf, S. Alsian, F. Ge, Y. Chong, K. Matzen, O. Quigley, J. Patterson, J. Tirado, S. Wu, and M. F. Cohen, "Practical 3d photography," 2019.
- [2] D. Fattal, Z. Peng, T. Tran, S. Vo, M. Fiorentino, J. Brug, and R. G. Beusoleil, "A multi-directional backlight for a wide-angle, glasses-free three-dimensional display," *Nature*, vol. 495, no. 7441, p. 348, 2013.
- [3] I. P. Howard, B. J. Rogers, et al., *Binocular vision and stereopsis*. Oxford University Press, USA, 1995.
- [4] N. Dodgson, J. Moore, and S. Lang, "Multi-view autostereoscopic 3d display," 11 2003.
- [5] D. Scharstein, "Stereo vision for view synthesis," in *Proceedings CVPR IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 852–858, IEEE, 1996.

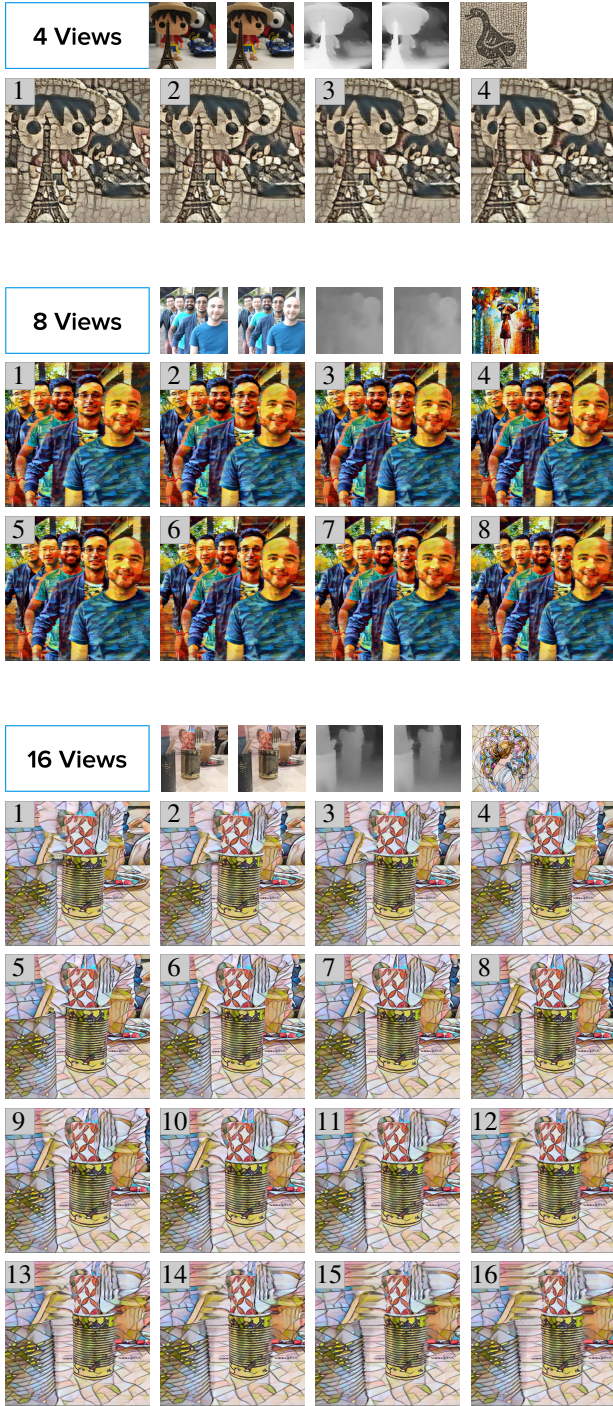


Fig. 5. Example results from running our style transfer pipeline and generating 4, 8, and 16 views respectively. The top row of each section shows the original stereo pairs, disparity maps, and style guide used. Each subsequent row shows the generated views from left-most to right-most viewpoint.

- [6] S. Avidan and A. Shashua, "Novel view synthesis in tensor space," in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 1034–1040, IEEE, 1997.
- [7] S. Avidan and A. Shashua, "Novel view synthesis by cascading trilinear tensors," *IEEE transactions on visualization and computer graphics*, vol. 4, no. 4, pp. 293–306, 1998.
- [8] N. Martin and S. Roy, "Fast view interpolation from stereo: Simpler can be better," *Proceedings of 3DPTV*, vol. 8, 2008.
- [9] T. Zhou, R. Tucker, J. Flynn, G. Fyffe, and N. Snavely, "Stereo magnification: Learning view synthesis using multiplane images," in *SIGGRAPH*, 2018.
- [10] S. Ravi, P. Pasupathi, S. Muthukumar, and N. Krishnan, "Image inpainting techniques-a survey and analysis," in *2013 9th International Conference on Innovations in Information Technology (IIT)*, pp. 36–41, IEEE, 2013.
- [11] K.-J. Oh, S. Yea, and Y.-S. Ho, "Hole filling method using depth based in-painting for view synthesis in free viewpoint television and 3-d video," in *2009 Picture Coding Symposium*, pp. 1–4, IEEE, 2009.
- [12] A. Telea, "An image inpainting technique based on the fast marching method," *Journal of graphics tools*, vol. 9, no. 1, pp. 23–34, 2004.
- [13] L. A. Gatys, A. S. Ecker, and M. Bethge, "A neural algorithm of artistic style," *arXiv preprint arXiv:1508.06576*, 2015.
- [14] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2414–2423, 2016.
- [15] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European Conference on Computer Vision*, 2016.
- [16] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. S. Lempitsky, "Texture networks: Feed-forward synthesis of textures and stylized images," in *ICML*, vol. 1, p. 4, 2016.
- [17] V. Dumoulin, J. Shlens, and M. Kudlur, "A learned representation for artistic style," *arXiv preprint arXiv:1610.07629*, 2016.
- [18] L. Sheng, Z. Lin, J. Shao, and X. Wang, "Avatar-net: Multi-scale zero-shot style transfer by feature decoration," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8242–8250, 2018.
- [19] G. Hua, "Ai with creative eyes amplifies the artistic sense of everyone," Jul 2017.
- [20] "Prisma labs: Turn memories into art using artificial intelligence," 2016.
- [21] F. Kooi and A. Toet, "Visual comfort of binocular and 3-d displays," *Displays*, vol. 25, pp. 99–108, 08 2004.
- [22] D. Chen, L. Yuan, J. Liao, N. Yu, and G. Hua, "Stereoscopic neural style transfer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6654–6663, 2018.
- [23] X. Gong, H. Huang, L. Ma, F. Shen, W. Liu, and T. Zhang, "Neural stereoscopic image style transfer," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 54–69, 2018.
- [24] M. Ruder, A. Dosovitskiy, and T. Brox, "Artistic style transfer for videos," in *German Conference on Pattern Recognition*, pp. 26–36, Springer, 2016.
- [25] D. Chen, J. Liao, L. Yuan, N. Yu, and G. Hua, "Coherent online video style transfer," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [26] H. Huang, H. Wang, W. Luo, L. Ma, W. Jiang, X. Zhu, Z. Li, and W. Liu, "Real-time neural style transfer for videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 783–791, 2017.
- [27] W. Hoff and N. Ahuja, "Surfaces from stereo: Integrating feature matching, disparity estimation, and contour detection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 11, no. 2, pp. 121–136, 1989.
- [28] S. Wanner and B. Goldluecke, "Variational light field analysis for disparity estimation and super-resolution," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 3, pp. 606–619, 2013.
- [29] N. Mayer, E. Ilg, P. Hausser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox, "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4040–4048, 2016.
- [30] W. Luo, A. G. Schwing, and R. Urtasun, "Efficient deep learning for stereo matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5695–5703, 2016.

- [31] H. Zhang and K. Dana, "Multi-style generative network for real-time transfer," 2017.
- [32] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv preprint arXiv:1607.08022*, 2016.
- [33] K. He, J. Sun, and X. Tang, "Guided image filtering," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, pp. 1397–1409, June 2013.
- [34] "Red hydrogen."
- [35] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *European Conf. on Computer Vision (ECCV)* (A. Fitzgibbon et al. (Eds.), ed.), Part IV, LNCS 7577, pp. 611–625, Springer-Verlag, Oct. 2012.