

Figure 3: End-to-end semi-supervised training scheme. We use reconstruction loss for synthetic images while image-based lighting loss is applied to both real and synthetic interview images.

for the project, both in the cost of the hardware, and the greatly increased storage cost of numerous high-speed uncompressed video streams.

The project settled for recording the survivors in just a single interview lighting condition consisting of diffuse, symmetrical lighting from above. But to enable relighting in the future, each survivor was recorded in a basis of forty-one lighting conditions in several static poses in a special session toward the end of each shoot as in Figure 5. The hope was that at some point, this set of static poses in different lighting conditions, plus the interview footage in diffuse lighting, could eventually be combined through machine learning to realistically show the interview as if it had been recorded in any combination of the lighting conditions, enabling general purpose relighting. This paper presents a technique to achieve this goal, which provides a practical process for recording interview footage where the lighting can be controlled realistically after filming.

2. Related Work

Relighting virtual humans from images and video is an active research topic in computer vision and computer graphics. In this section, we summarize some of the most related work in inverse rendering, image-based relighting, and learning-based relighting methods.

Inverse Rendering. If photos of a scene can be analyzed to derive an accurate model of the scene’s geometry and materials, the model can be rendered under arbitrary new lighting using forward rendering. This *inverse rendering*

problem is a long-studied research topic in computer vision and graphics [24, 29, 41]. Many of these approaches use strong assumptions such as known illumination [15], or hand-crafted priors [2].

Unsurprisingly, relighting human bodies and faces has received particular interest recently. Many parametric models have been proposed to jointly reconstruct geometry, reflectance, and illumination of human bodies [37], faces [5, 11, 12, 14, 17, 38], eyes [3], eyelids [4], and hair [16, 42]. [22] relights videos of humans based on estimation of parametric BRDF models and wavelet-based incident illumination. [44] uses a diffuse model for the face to relight it with a radiance environment map using ratio images. [9] performs relighting by using spherical gradient illumination images to fit a cosine lobe the reflectance function. Several works estimate spatially-varying reflectance properties of a scene from either flash [27, 23] or flat-lit images [13]. [40] uses deep neural networks to estimate the parameters of a predefined geometry and reflectance model from a single image.

These parametric models are typically designed to handle specific parts of the human body. Many of these techniques rely on lightweight morphable models for geometry, a Lambertian model for skin reflectance, and a low-frequency 2nd order spherical harmonic basis for illumination. Unfortunately, these strong priors only capture low-frequency detail and do not reproduce the appearance of specular reflections and sub-surface scattering in the skin.

In contrast, we use a deep neural network to infer the subject’s reflectance field, which can be used to relight the images without explicitly modeling the geometry, material reflectance, and illumination of the images.

Image-based Relighting. When a person is recorded under a large number of individual lighting conditions, they can be accurately relit by linearly combining those one-light-at-a-time (OLAT) images with the target illumination [7]. [39] used high-speed video to capture dynamic subjects with time-multiplexed lighting conditions for relighting, but required expensive cameras, optical flow computation, and was data intensive. [8] recorded a coarser lighting basis from a multitude of viewpoints, allowing post-production control of both the lighting and the viewpoint. However, neither technique could be applied to long video segments due to the large size of the high frame rate video.

Another technique is to transfer reflectance field properties from a pre-captured subject to a target subject’s performance as in [28]. However, the quality of the lighting transfer depends on number of captured poses and the similarity in appearance of the two subjects. Relighting can also be performed by transferring local image statistics from one portrait image to the target portrait as in [32]; however, this technique does not work well for extreme lighting changes.

Our approach also uses a set of lighting basis conditions

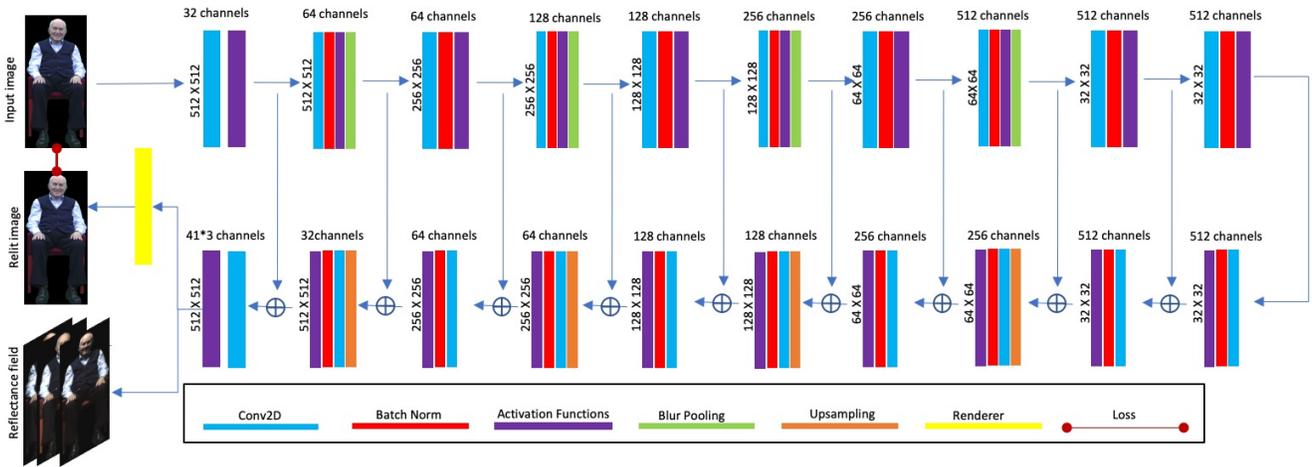


Figure 4: The architecture of our neural network. The input image is passed through a U-Net style architecture to regress to the set of OLAT images. When the ground truth is available, the network prioritizes the reconstruction loss of the OLAT imageset. Otherwise, the network is trained based on the feedback of the relit image.

to perform relighting. But instead of recording OLAT’s for every moment of the video, our neural network infers OLATs for each video frame based on exemplars from static poses, enabling dynamic performance relighting.

Learning-Based Relighting [25] trains a deep neural network to map images of a subject lit by spherical gradient basis illumination to a set of one-light-at-a-time (OLAT) images for relighting. Similar to this approach, we map the interview lighting images to a set of OLAT images. But unlike [25], we employ a semi-supervised training scheme to train the network due to the lack of ground truth in our dataset and to work on the single interview lighting condition that’s available. [34] proposes a neural network that takes a single portrait photo under any lighting environment and relights the subject with arbitrary target illumination. This network was trained on a large set of subjects individuals under a dense set of lighting conditions to predict the input illumination and perform relighting by replacing the illumination at the bottleneck of the neural network. The technique is overall successful, but the low resolution of the predicted illumination limits the quality of the relit result. [21] also uses multitudes of data from 70 subjects with dense lighting conditions to estimate a more detailed HDR lighting environment from a single portrait image. [36] presents a recent advance in Style Transfer techniques, where a video can be changed to a different style by registering to one or more keyframes in the new style. This is most often used to transfer non-photorealistic rendering styles such as a pastel drawing, but can also be used to transfer a new style of lighting. However, this technique has not been applied to create arbitrarily relightable models and requires registration from the style exemplars to the video sequence. In comparison,

our method is designed to perform realistic relighting from a single lighting condition by providing the neural network a set of reflectance field exemplars of how the subject actually should appear under OLAT lighting conditions.

3. Method

One of the most effective ways to perform realistic relighting is to combine a dense set of basis lighting conditions (a *reflectance field*) with according to a novel lighting environment to simulate the appearance in the new lighting. However, this approach is not ideal for a dynamic performance since it requires either high-speed cameras, or requires the actor to sit still for several seconds to capture the set of OLAT images. [25] overcomes this limitation by using neural networks to regress 4D reflectance fields from just two images of a subject lit by gradient illumination. They postulate that one can also use flat-lit images to achieve similar results with less high-frequency detail. Since the method casts relighting as a supervision regression problem, it requires pairs of tracking images and their corresponding OLAT images as ground truth for training.

In the New Dimensions in Testimony project, most of the Holocaust survivors’ interview footage was captured in front of a green screen so that the virtual backgrounds can be added during post-production. However, this setup poses difficulties for achieving consistent illumination between the actors and the backgrounds in the final testimony videos and does not provide the ground truth needed for supervision training. In this paper, we use the limited OLAT data to train a neural network to infer reflectance fields from synthetically relit images. The synthetic relit images are improved by matching them with the input interview im-



Figure 5: Reflectance field: 27 of 41 one-light-at-a-time images.

ages through a differentiable renderer, enabling an end-to-end training scheme. For more training details, see Figure 3.

In this section, we describe the data acquisition process, how we relate the OLAT reflectance field exemplars with the interview footage, and how we train an end-to-end neural network to regress reflectance fields for realistic relighting.

3.1. Data Acquisition and Processing

Each Holocaust survivor was recorded over a 180-degree field of view using an array of 50 Panasonic X900MK 60fps progressive scan consumer camcorders, each four meters away framed on the subject. Toward the end of each Holocaust survivor’s lengthy interview, they were captured in several different static poses under a reflectance field lighting basis of 41 lighting conditions as in Figure 5. The lighting conditions were formed using banks of approximately 22 lights each of the 931 light sources on the 8m diameter dome [8]. This somewhat lower lighting resolution was chosen to keep the capture time shorter than what would be required to record each of the 931 lights individually and to avoid too great of a degree of underexposure so that we could use the same exposure settings as the interview lighting without touching the cameras.

Data Processing. The original resolution of our video frames images is 1920×1080 in portrait orientation. For each image, we crop the full body of the actors, and then use Grabcut [31] to mask out the background. The images are then padded and resized to 512×512 .

Synthetic Tracking Frames. We use a mirror ball image captured right after the interview session as a light probe [6]. This light probe represents the illumination of the interview session. For convenience, we convert the light probe

to a latitude-longitude format. Then we use mirror ball images captured in the OLAT session to find their projections in our target environment illumination map. By taking a weighted combination of the images in the OLAT set according to these projections, we are able to relight all the static poses of the actor.

The OLAT images are not as well exposed as the interview lighting since fewer light sources are on, and we discovered that the consumer video cameras applied a weaker level of gamma correction to the darker range of pixel values, making dark regions appear even darker, presumably as a form of noise suppression. Thus, we developed a dual gamma correction curve to linearize the image data:

$$I' = (1 - I) * I^{\gamma_1} + I * I^{\gamma_2} \quad (1)$$

where γ_1, γ_2 describe the gamma we use for the lower and upper part of the gamma curve, and we interpolate between these two curves according to the brightness of the pixel. We optimize γ_1, γ_2 so that the OLAT reflectance field exemplars, relit with the measured interview lighting condition, match the appearance of the first frame of the interview video. Though each subject is only recorded as a reflectance field in a few poses, these synthetic relit images play an important role in bringing the output of the network closer to the illumination of the input video footage.

3.2. Network Architecture

We cast the relighting problem as prediction of the reflectance field, and use these measurements to render the subject under arbitrary illumination. To be consistent with the Holocaust survivor dataset, we define a reflectance field to have 41 OLAT images. Our goal is to predict how the actor would look under 41 specified lighting conditions for every frame of dynamic performance. The structure of our neural network resembles the structure of the popular image transformation architecture with skip connections [30]. The encoder consists of ten blocks of 3×3 convolution layers each followed by a batch-normalization layer and a leaky ReLU activation function. A blur-pooling operation [43] is used at the end of the block to decrease the spatial resolution and increase the number of channels. Note that the first block of the encoder does not have a batch-normalization layer and uses a 7×7 convolution layer.

The decoder follows a similar structure with ten blocks of bilinear upsampling followed by a convolution layer. At the end of each decoder block, we use skip connections to concatenate the network features with their corresponding activations in the encoder. All convolution layers are followed by a ReLU activation except for the last convolution layer where a sigmoid activation is used. At the end of the decoder is a differentiable renderer that takes as input a whole set of OLAT images to render a subject under a new



(a) Reference (b) Sun et al.[34] (c) Ours

Figure 6: Comparison with Single Image Portrait Relighting. Our result has greater lighting detail and looks much closer to the reference lighting.

calibrated illumination condition. For network details, see Figure 4.

3.3. Loss Function

Our model is trained through the minimization of a weighted combination of two loss functions. A reconstruction loss minimizes errors between the set of OLAT images in the dataset and set of OLAT images predicted by the network. The second loss is an image-based relighting loss that minimizes the errors between the input image and the rendered image lit with the predicted reflectance field. The backgrounds are masked out in all loss calculations.

Reconstruction Loss. This loss ensures the accurate inference of the network by matching the network prediction with the ground truth. Since the per-pixel photometric loss often leads to blurry output images, we choose to minimize the loss in feature space with a perceptual loss. Letting $VGG^{(i)}(I)$ be the activations of the i th layer of a VGG network [33], the reconstruction loss is defined as:

$$L_{rec} = \sum_{i=1}^N \sum_{j=1}^M \|VGG^{(j)}(I_{pred}^i) - VGG^{(j)}(I_{gt}^i)\|_2 \quad (2)$$

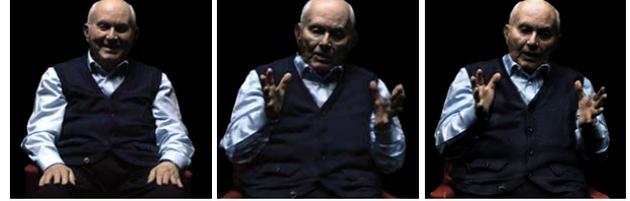
where N is the number of images in a complete OLAT set, and M is the number of VGG layers to be used.

Image-Based Relighting Rendering Loss This self-supervision loss makes the network more robust to unseen poses of the actor in the training set. Given a predicted reflectance field $R(\theta, \phi, x, y)$ and the calibrated interview lighting environment L_i , we can relight the actor as follows:

$$I_{relit} = \sum_{\theta, \phi} R_{x,y}(\theta, \phi) L_i(\theta, \phi) \quad (3)$$

where $R_{x,y}(\theta, \phi)$ represents how much light is reflected toward the camera by pixel (x, y) as a result of illumination from direction (θ, ϕ) . Matching this relit image I_{relit} and the input image I gives us the rendering loss:

$$L_{render} = \sum_{j=1}^M \|VGG^{(j)}(I_{relit}) - VGG^{(j)}(I)\|_2 \quad (4)$$



(a) Reference (b) Texler et al. [35] (c) Ours

Figure 7: Comparison with Style Transfer based relighting. Our method reproduces more convincing shadows and high-lights.

The full objective is the weighted combination of the two loss functions:

$$L = \lambda_1 L_{rec} + \lambda_2 L_{render} \quad (5)$$

Implementation details. We use two sets of data to train the network. The first set consists of six poses with groundtruth OLAT images and the six corresponding relit images showing the reflectance field exemplars under the simulated interview lighting condition. The second set consists of 100 frames of the target video. We train the first set with both a reconstruction loss L_{rec} and a rendering loss L_{render} for 100 epochs, and then we train the second set for 4 epochs with only the rendering loss before going back to supervised training. The training process continues until we reach 1040 epochs. We use the ADAM optimizer [20] with $\beta_1 = 0.9, \beta_2 = 0.999$ and a learning rate of 0.001.

4. Evaluation

We evaluate our technique by relighting several hundred frames of interview footage and comparing to relit images made with that of [34] and [35]. We do not have ground truth relighting for each frame of the video to compare to, so we employ a user study to evaluate our method against prior works. Finally, we show how our method is able to realistically relight the dynamic performance of the subject with arbitrary poses and motions.

4.1. Single Image Portrait Relighting

We first compared our method with a state-of-the-art lighting estimation and relighting for portrait photos [34]. Their neural network was trained on a dataset of numerous synthetically relit portrait images of 18 individuals from pre-captured OLAT data. From Figure 6, we can see that our method performs much more believable relighting, as the single image portrait relighting result only reproduces the low-frequency components of the novel illumination. Note also that we cropped our method’s result down to just the face to match the output capability of the Single Image Portrait Relighting network, whereas our model is able to relight more of the body as shown in Figure 7.



Figure 8: **Relighting results** - Row 1: Input interview videos. Row 2,3: OLAT predictions on two patterns. Row 4,5: Relighting results with two HDRI lighting environments: Grace Cathedral and Pisa Courtyard. See more examples in our video.

4.2. Relighting as Style Transfer

We next compared our approach with the state-of-the-art style transfer technique of [35] that takes several keyframes to use as style and transfer the styles or relighting from those keyframes to the video. As we can see from Figure 7, the shading on the inner palms of the actor is not supposed to be in shadow, but since the provided keyframes do not cover this pose, [35] predicts the wrong shading in this area. In contrast, thanks to self-supervised learning, our network is able to recover a more reasonable rendition of the shading

one would expect for this pose. For side by side comparison, see our supplementary video.

4.3. User Study

We conducted a user study to evaluate which relighting technique produced preferable results. We showed users a reference image of the subject under one of the OLAT conditions, and then short video clips of the subject’s interview re-lit by that condition using our approach, Single Image Portrait Relighting [34], and Style Transfer based

relighting [35]. We finally asked users two questions: 1) Which video clip looks more like the reference image, and 2) Which video clip looks better? From 61 responses, all users answered both questions with the same answer: 52 chose the video clip rendered with our approach, and 9 chose the video rendered with [35], while none chose [34]. This showed a clear preference for our approach.

4.4. Relighting Dynamic Performance

We perform relighting for interview footage of three Holocaust survivors. The first survivor was recorded in 2012, while the other two survivors were recorded in 2015. In 2012, the OLAT set consists of 41 patterns while there are 146 patterns used in 2015. Because our method is not restricted to any OLAT patterns, it can generalize to the new setup as long as the diffuse lighting condition from above is guaranteed. For consistency, we choose evenly distributed 41 out of 146 OLAT patterns to train our network. It is important to note that none of the evaluated interview videos are used to train the neural network. As we can see from Figure 8, our network is able to predict convincing reflectance fields for novel poses in the interview videos, enabling it to realistically place these interviews in any lighting environment.

5. Future Work

In this project, we made use of both the diffusely-lit interview footage and the reflectance field exemplars of each subject, but we only used a single one of the available viewpoints in the data. It seems possible that even better relighting results could be obtained by leveraging some or all of the views of the subject from the other cameras' positions, even though these other views are also recorded in the same diffuse interview lighting. The reason is that the multiple viewpoints carry additional information about the subject's three-dimensional shape, and knowing the subject's 3D shape is also useful for predicting their shading and shadowing under new lighting conditions. For future work, it would be of interest to use the 50 viewpoints available to reconstruct a 3D model such as a Neural Radiance Field [26] for each frame of the interview footage and to leverage these models during training so that the network is better able to learn how shape and the appearance under novel illumination are connected. However, at this time, such reconstruction techniques might be prohibitively expensive to run on hours of video material.

6. Conclusion

In this paper, we presented a deep learning-based video relighting technique that takes diffusely lit video and a set of reflectance field exemplars of the same subject as input. We designed this technique to work with the data available

from the Holocaust survivor interviews recorded in 2014 in the New Dimensions in Testimony project and showed how we can realistically render the Holocaust survivor interview footage in novel lighting conditions. The technique suggests that this approach could be used to obtain high-quality relighting of new interview footage, assuming that the subjects can also be recorded under a variety of directional lighting conditions in a number of static poses. This provides the relighting network with subject-specific information for how to relight the video than just the single interview lighting condition alone.

7. Acknowledgement

We wish to thank the SHOAH foundation for their effort in making the project New Dimension in Testimony possible and for sharing precious data of the three Holocaust survivors Pinchas Gutter, Aaron Elster, Eva Schloss. We thank Kathleen Haase and Christina Trejo for their coordination. We thank Yajie Zhao and Mingming He for their insightful suggestions. This research was sponsored by the Army Research Office and was accomplished under Cooperative Agreement Number W911NF-20-2-0053, and sponsored by the U.S. Army Research Laboratory (ARL) under contract number W911NF-14-D-0005, the CONIX Research Center, one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program sponsored by DARPA and in part by the ONR YIP grant N00014-17-S-FO14. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute.

References

- [1] Ron Artstein, David Traum, Oleg Alexander, Anton Leuski, Andrew Jones, Kallirroi Georgila, Paul Debevec, William Swartout, Heather Maio, and Stephen Smith. Time-offset interaction with a holocaust survivor. In *Proceedings of the 19th International Conference on Intelligent User Interfaces, IUI '14*, page 163–168, New York, NY, USA, 2014. Association for Computing Machinery.
- [2] Jonathan T. Barron and Jitendra Malik. Shape, illumination, and reflectance from shading. *TPAMI*, 2015.
- [3] Pascal Bérard, Derek Bradley, Markus Gross, and Thabo Beeler. Lightweight eye capture using a parametric model. *ACM Trans. Graph.*, 35(4), July 2016.
- [4] Amit Bermano, Thabo Beeler, Yeara Kozlov, Derek Bradley, Bernd Bickel, and Markus Gross. Detailed spatio-temporal reconstruction of eyelids. *ACM Trans. Graph.*, 34(4), July 2015.
- [5] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Tech-*

- niques, SIGGRAPH '99, page 187–194, USA, 1999. ACM Press/Addison-Wesley Publishing Co.
- [6] Paul Debevec. Rendering synthetic objects into real scenes: Bridging traditional and image-based graphics with global illumination and high dynamic range photography. In *Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '98, page 189–198, New York, NY, USA, 1998. Association for Computing Machinery.
- [7] Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar. Acquiring the reflectance field of a human face. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '00, page 145–156, USA, 2000. ACM Press/Addison-Wesley Publishing Co.
- [8] Per Einarsson, Charles-Felix Chabert, Andrew Jones, Wan-Chun Ma, Bruce Lamond, Tim Hawkins, Mark Bolas, Sebastian Sylwan, and Paul Debevec. Relighting Human Locomotion with Flowed Reflectance Fields. In *Eurographics Symposium on Rendering (2006)*, June 2006.
- [9] Graham Fyffe. Cosine lobe based relighting from gradient illumination photographs. In *SIGGRAPH '09: Posters*, SIGGRAPH '09, New York, NY, USA, 2009. Association for Computing Machinery.
- [10] Graham Fyffe, Tim Hawkins, Chris Watts, Wan-Chun Ma, and Paul Debevec. Comprehensive facial performance capture. *Comput. Graph. Forum*, 30:425–434, 04 2011.
- [11] Pablo Garrido, Levi Valgaert, Chenglei Wu, and Christian Theobalt. Reconstructing detailed dynamic face geometry from monocular video. *ACM Trans. Graph.*, 32(6), Nov. 2013.
- [12] Pablo Garrido, Michael Zollhöfer, Dan Casas, Levi Valgaerts, Kiran Varanasi, Patrick Pérez, and Christian Theobalt. Reconstruction of personalized 3d face rigs from monocular video. *ACM Trans. Graph.*, 35(3), May 2016.
- [13] Paulo Gotardo, Jérémy Riviere, Derek Bradley, Abhijeet Ghosh, and Thabo Beeler. Practical dynamic facial appearance modeling and acquisition. *ACM Trans. Graph.*, 37(6), Dec. 2018.
- [14] Tim Hawkins, Andreas Wenger, Chris Tchou, Andrew Gardner, Fredrik Goransson, and Paul Debevec. Animatable Facial Reflectance Fields. In *Eurographics Symposium on Rendering*, Norkoping, Sweden, 2004.
- [15] B. K.P. Horn. Shape from shading: A method for obtaining the shape of a smooth opaque object from one view. Technical report, USA, 1970.
- [16] Liwen Hu, Chongyang Ma, Linjie Luo, and Hao Li. Single-view hair modeling using a hairstyle database. *ACM Trans. Graph.*, 34(4), July 2015.
- [17] Alexandru Eugen Ichim, Sofien Bouaziz, and Mark Pauly. Dynamic 3d avatar creation from hand-held video input. *ACM Trans. Graph.*, 34(4), July 2015.
- [18] A. Jones, K. Nagano, J. Busch, X. Yu, H. Peng, J. Barreto, O. Alexander, M. Bolas, P. Debevec, and J. Unger. Time-offset conversations on a life-sized automultiscopic projector array. In *2016 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 927–935, 2016.
- [19] Andrew Jones, Jonas Unger, Koki Nagano, Jay Busch, Xueming Yu, Hsuan-Yueh Peng, Oleg Alexander, Mark Bolas, and Paul Debevec. An automultiscopic projector array for interactive digital humans. In *ACM SIGGRAPH 2015 Emerging Technologies*, SIGGRAPH '15, New York, NY, USA, 2015. Association for Computing Machinery.
- [20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [21] Chloe LeGendre, Wan-Chun Ma, Rohit Pandey, Sean Fanello, Christoph Rhemann, Jason Dourgarian, Jay Busch, and Paul Debevec. Learning illumination from diverse portraits. In *SIGGRAPH Asia 2020 Technical Communications*, SA '20, New York, NY, USA, 2020. Association for Computing Machinery.
- [22] Guannan Li, Chenglei Wu, Carsten Stoll, Yebin Liu, Kiran Varanasi, Qionghai Dai, and Christian Theobalt. Capturing Relightable Human Performances under General Uncontrolled Illumination. *Computer Graphics Forum*, 2013.
- [23] Zhengqin Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. Learning to reconstruct shape and spatially-varying reflectance from a single image. In *SIGGRAPH Asia 2018 Technical Papers*, page 269. ACM, 2018.
- [24] S. Lombardi and K. Nishino. Reflectance and illumination recovery in the wild. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(1):129–141, 2016.
- [25] Abhimitra Meka, Christian Haene, Rohit Pandey, Michael Zollhoefer, Sean Fanello, Graham Fyffe, Adarsh Kowdle, Xueming Yu, Jay Busch, Jason Dourgarian, Peter Denny, Sofien Bouaziz, Peter Lincoln, Matt Whalen, Geoff Harvey, Jonathan Taylor, Shahram Izadi, Andrea Tagliasacchi, Paul Debevec, Christian Theobalt, Julien Valentin, and Christoph Rhemann. Deep reflectance fields - high-quality facial reflectance field inference from color gradient illumination. volume 38, July 2019.
- [26] Ben Mildenhall, Pratul Srinivasan, Matthew Tancik, Jonathan Barron, Ravi Ramamoorthi, and Ren Ng. *NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis*, pages 405–421. 11 2020.
- [27] Giljoo Nam, Joo Ho Lee, Diego Gutierrez, and Min H. Kim. Practical svbrdf acquisition of 3d objects with unstructured flash photography. *ACM Trans. Graph.*, 37(6), Dec. 2018.
- [28] Pieter Peers, Naoki Tamura, Wojciech Matusik, and Paul Debevec. Post-production facial performance relighting using reflectance transfer. *ACM Transactions on Graphics*, 26(3), July 2007.
- [29] Ravi Ramamoorthi and Pat Hanrahan. A signal-processing framework for inverse rendering. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '01, page 117–128, New York, NY, USA, 2001. Association for Computing Machinery.
- [30] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, volume 9351 of LNCS, pages 234–241. Springer, 2015. (available on arXiv:1505.04597 [cs.CV]).

- [31] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. "grabcut": Interactive foreground extraction using iterated graph cuts. In *ACM SIGGRAPH 2004 Papers*, SIGGRAPH '04, page 309–314, New York, NY, USA, 2004. Association for Computing Machinery.
- [32] YiChang Shih, Sylvain Paris, Connelly Barnes, William T. Freeman, and Frédo Durand. Style transfer for headshot portraits. *ACM Trans. Graph.*, 33(4), July 2014.
- [33] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [34] Tiancheng Sun, Jonathan T. Barron, Yun-Ta Tsai, Zexiang Xu, Xueming Yu, Graham Fyffe, Christoph Rhemann, Jay Busch, Paul Debevec, and Ravi Ramamoorthi. Single image portrait relighting. *ACM Trans. Graph.*, 38(4), July 2019.
- [35] Ondřej Texler, David Futschik, Jakub Fišer, Michal Lukáč, Jingwan Lu, Eli Shechtman, and Daniel Sýkora. Arbitrary style transfer using neurally-guided patch-based synthesis. *Computers & Graphics*, 87:62–71, 2020.
- [36] Ondřej Texler, David Futschik, Michal Kučera, Ondřej Jamriška, Šárka Sochorová, Menglei Chai, Sergey Tulyakov, and Daniel Sýkora. Interactive video stylization using few-shot patch-based training. *ACM Transactions on Graphics*, 39(4):73, 2020.
- [37] Christian Theobalt, Naveed Ahmed, Hendrik Lensch, Marcus Magnor, and Hans-Peter Seidel. Seeing people in different light - joint shape, motion, and reflectance capture. *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, 13(4):663–674, Jul 2007.
- [38] J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Nießner. Face2Face: Real-time Face Capture and Reenactment of RGB Videos. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2016.
- [39] Andreas Wenger, Andrew Gardner, Chris Tchou, Jonas Unger, Tim Hawkins, and Paul Debevec. Performance relighting and reflectance transformation with time-multiplexed illumination. *ACM Trans. Graph.*, 24(3):756–764, July 2005.
- [40] Shugo Yamaguchi, Shunsuke Saito, Koki Nagano, Yajie Zhao, Weikai Chen, Kyle Olszewski, Shigeo Morishima, and Hao Li. High-fidelity facial reflectance and geometry inference from an unconstrained image. *ACM Trans. Graph.*, 37(4), July 2018.
- [41] Yizhou Yu, Paul Debevec, Jitendra Malik, and Tim Hawkins. Inverse global illumination: Recovering reflectance models of real scenes from photographs. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '99, page 215–224, USA, 1999. ACM Press/Addison-Wesley Publishing Co.
- [42] Meng Zhang, Menglei Chai, Hongzhi Wu, Hao Yang, and Kun Zhou. A data-driven approach to four-view image-based hair modeling. *ACM Trans. Graph.*, 36(4), July 2017.
- [43] Richard Zhang. Making convolutional networks shift-invariant again. In *ICML*, 2019.
- [44] Zhen Wen, Zicheng Liu, and T. S. Huang. Face relighting with radiance environment maps. In *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, volume 2, pages II–158, 2003.