

DeepFovea++: Reconstruction and Super-Resolution for Natural Foveated Rendered Videos



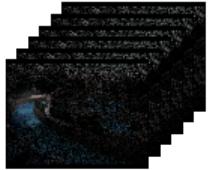
Christoph Reich Marius Memmel Jonas Henry Grebe



Fovea Rendered Video Reconstruction and Super-Resolution

Problem Setting REDS dataset [1]

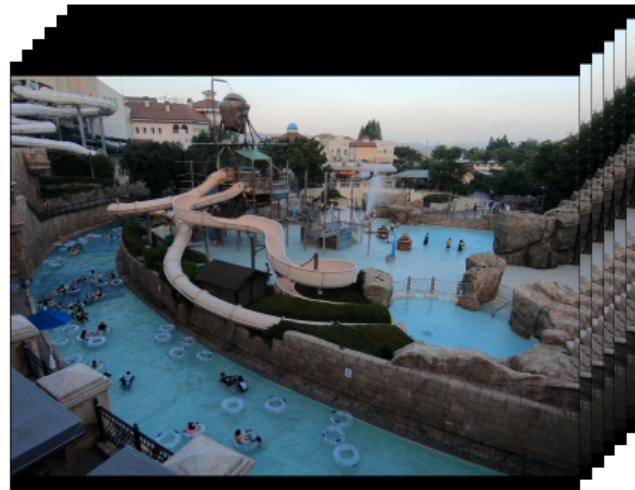
Fovea sampled input sequence



Reconstruction &
Super-Resolution



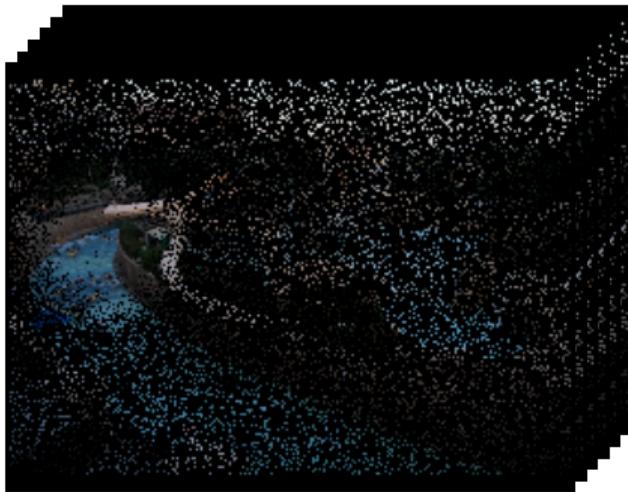
Reconstructed & upsampled sequence



Fovea Rendered Video Reconstruction and Super-Resolution

Problem Setting REDS dataset [1]

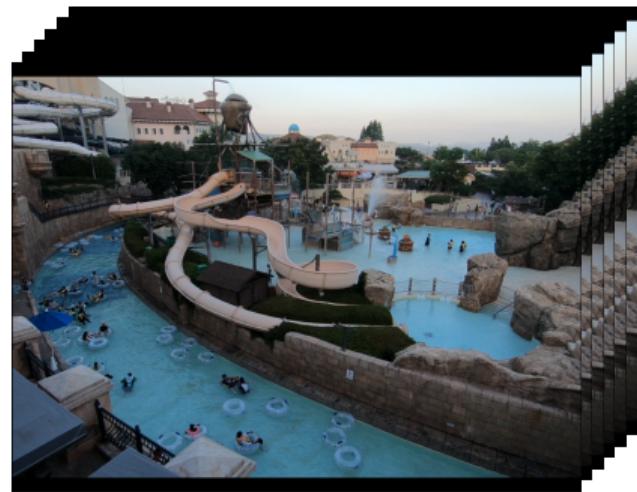
Fovea sampled input sequence



Reconstruction &
Super-Resolution



Reconstructed & upsampled sequence



DeepFovea - Kaplanyan et al. [2] (Facebook AI)

- Reconstructions of most plausible peripheral video from small portion of pixels in each frame
- Able to reconstruct video sequences (128×128)
- Recurrent U-Net reconstruction network
- Loss combination (adversarial, perceptual, optical flow)

Detail-revealing deep video super-resolution - Tao et al. [3] (2017)

- Sub-pixel motion compensation layers (improved frame alignment)

Video restoration with enhanced deformable convolutional networks - Wang et al. [4] (2019)

- Enhanced deformable convolution (frame alignment, temporal and spatial attention)

Detail-revealing deep video super-resolution - Tao et al. [3] (2017)

- Sub-pixel motion compensation layers (improved frame alignment)

Video restoration with enhanced deformable convolutional networks - Wang et al. [4] (2019)

- Enhanced deformable convolution (frame alignment, temporal and spatial attention)

Additional Work

- Video super-resolution with convolutional neural networks - Kappeler et al. [5] (2016)
- Video super-resolution via deep draft-ensemble learning - Liao et al. [6] (2015)
- Dictionary-based multiple frame video super-resolution - Dai et al. [7] (2015)

Reconstruction Model Architecture

Method

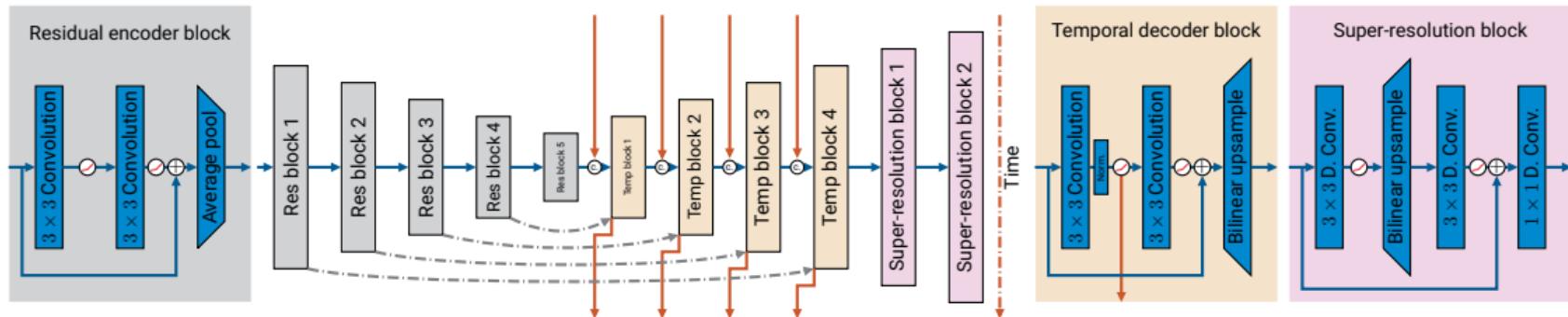


Figure: Architecture of the reconstruction network.

Loss function

$$\mathcal{L} = w_{\text{sv}} \mathcal{L}_{\text{sv}} + w_{\text{adv}} \mathcal{L}_{\text{adv}} + w_{\text{adv fft}} \mathcal{L}_{\text{adv fft}} + w_{\text{flow}} \mathcal{L}_{\text{flow}} + w_{\text{LPIPS}} \mathcal{L}_{\text{LPIPS}} \quad (1)$$

Loss function

$$\mathcal{L} = w_{\text{sv}} \mathcal{L}_{\text{sv}} + w_{\text{adv}} \mathcal{L}_{\text{adv}} + w_{\text{adv fft}} \mathcal{L}_{\text{adv fft}} + w_{\text{flow}} \mathcal{L}_{\text{flow}} + w_{\text{LPIPS}} \mathcal{L}_{\text{LPIPS}} \quad (1)$$

- Supervised loss (general and adaptive robust loss function - Barron [8])

$$\mathcal{L}_{\text{sv}} = p(x, \alpha, c) + \log Z(\alpha), \quad p(x = \sum_{i \in (chw)} \hat{\mathbf{I}} - \mathbf{I}, \alpha, c) = \frac{|\alpha-2|}{\alpha} \left(\left(\frac{(x/c)^2}{|\alpha-2|} + 1 \right)^{(\alpha/2)} - 1 \right)$$

- Adversarial loss [9] $\mathcal{L}_{\text{adv, adv fft}} = -\mathbb{E} [\log(D(\hat{\mathbf{I}}))]$

- Optical flow loss [10, 2] $\mathcal{L}_{\text{flow}} = \frac{1}{5} \sum_{i=1}^5 \frac{1}{c_{\text{rgb}} h_i w_i} \left\| \hat{\mathbf{I}}_i - \text{Warp} \left[\hat{\mathbf{I}}_{i+1} \right] \right\|_1$

- Perceptual loss [11] $\mathcal{L}_{\text{LPIPS}} = \frac{1}{5} \sum_{i=1}^5 \frac{1}{b c_i h_i w_i} \left\| \text{VGG}_{i,2}(\hat{\mathbf{I}}) - \text{VGG}_{i,2}(\mathbf{I}) \right\|_1$

Discriminator Architecture

Method

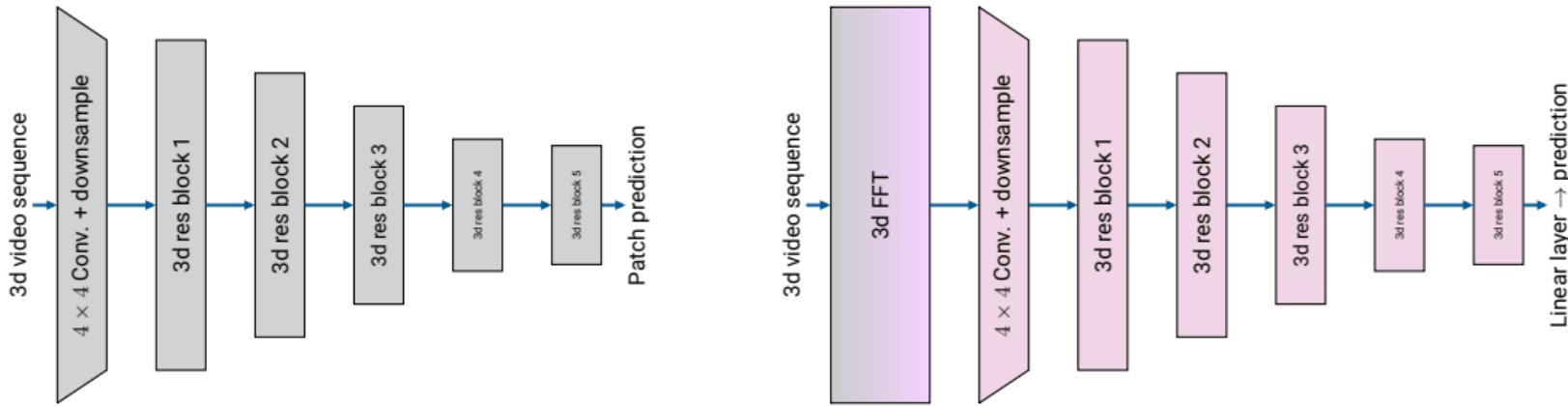


Figure: Discriminator network on the left and FFT discriminator network on the right.

REDS dataset - Nah et al. [1]

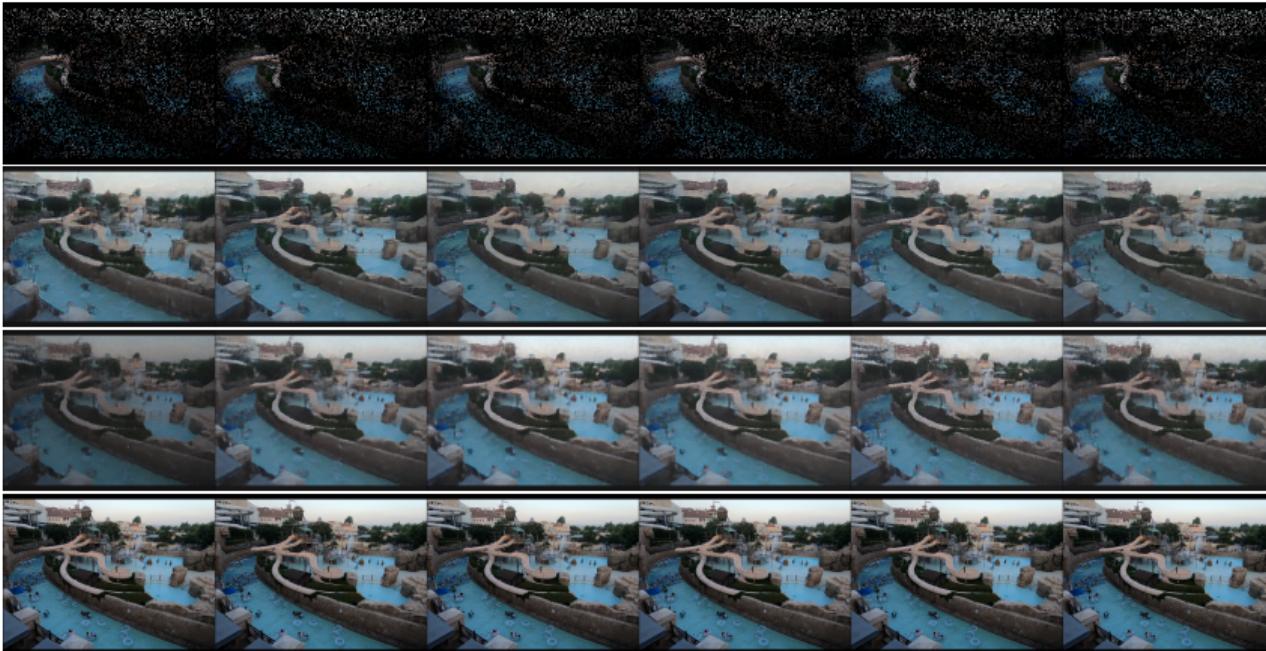
- 300 sequences of 100 high-quality natural RGB video frames each
- Resolution 720×1280

Fovea Sampling

- Low chance of pixels not being masked out in or closer to focus point
- Higher chance of being masked out moving away from focus point
- Mask generated on downsampled image (relates to approximately 19% of the information in the low-res image and to 1.1% of pixels compared to the high-res image)

Results

Qualitative results - with and without reset

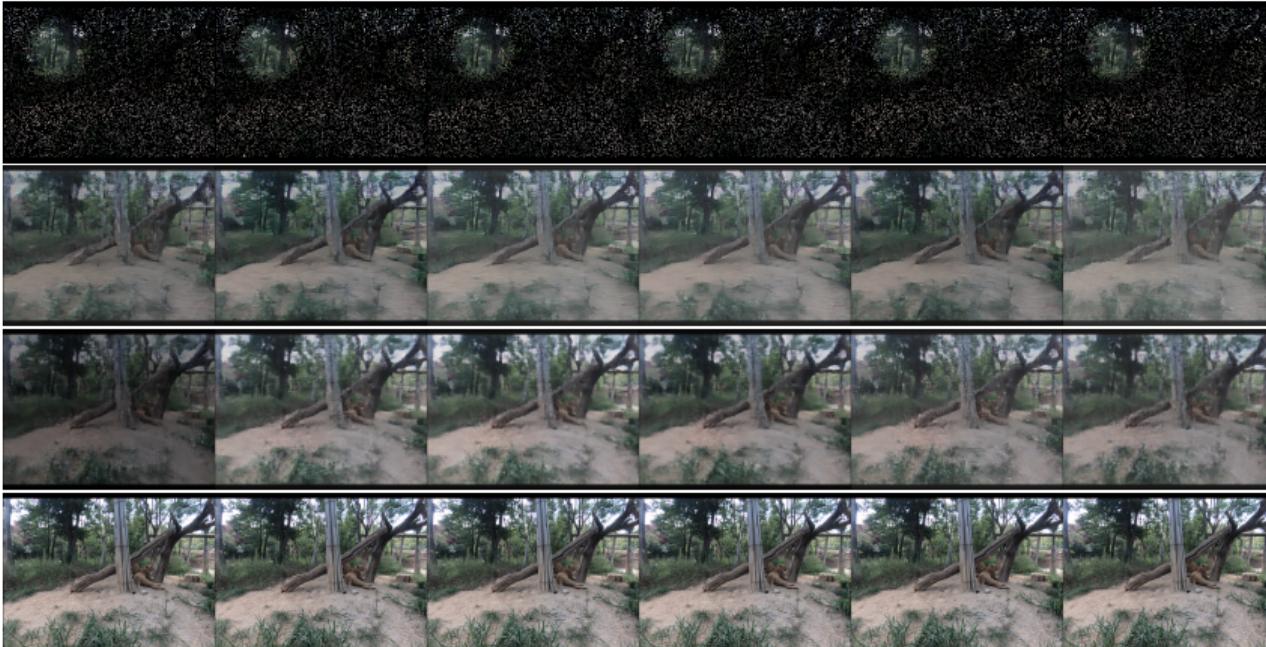


Results

Qualitative results - with and without reset

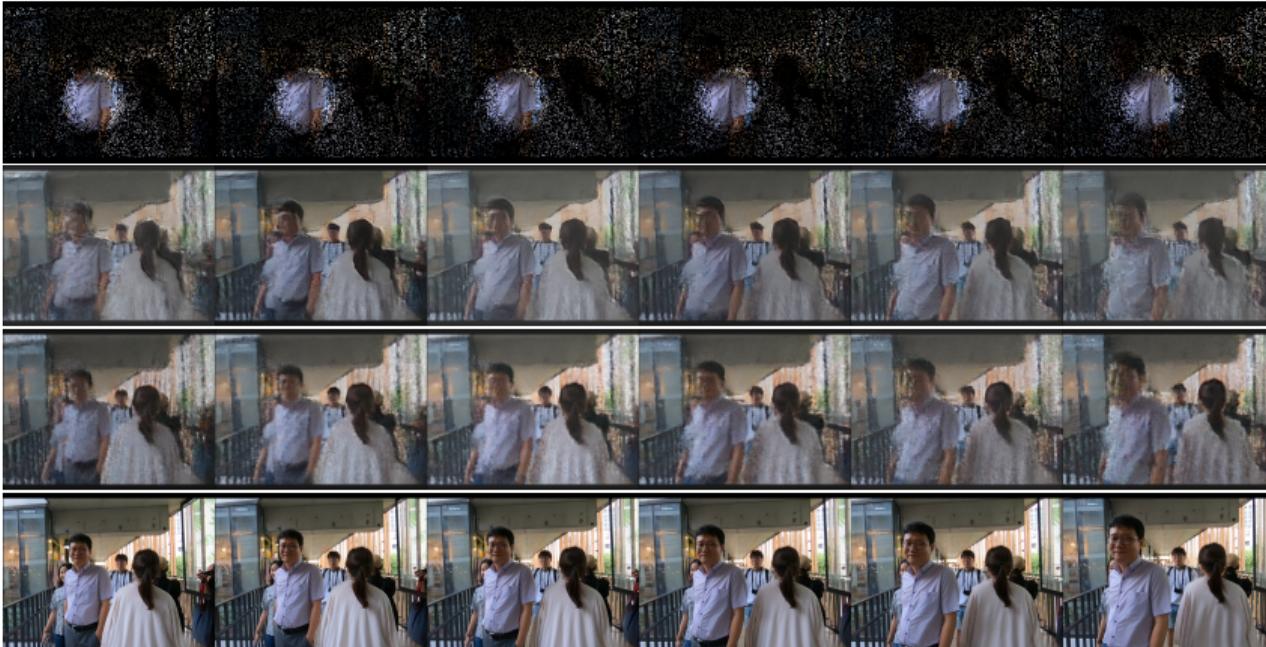


TECHNISCHE
UNIVERSITÄT
DARMSTADT



Results

Qualitative results - with and without reset



Results

Quantitative results

Reset	L1 ↓	L2 ↓	Peak signal-to-noise ratio PSNR ↑	Structural similarity SSIM ↑
✓	0.0701	0.0117	22.6681	0.9116
✗	0.061	0.009	23.8755	0.9290

$$\text{PSNR} = 10 \log_{10} \left(\frac{\max \left\{ \hat{\mathbf{I}} \right\}^2}{\text{L2} (\hat{\mathbf{I}}, \mathbf{I})} \right)$$
$$\text{SSIM} = \frac{4\mathbb{E} [\hat{\mathbf{I}}] \mathbb{E} [\mathbf{I}] \text{Cov} [\hat{\mathbf{I}}, \mathbf{I}]}{\left(\mathbb{E} [\hat{\mathbf{I}}]^2 + \mathbb{E} [\mathbf{I}]^2 \right) (\text{Var} [\hat{\mathbf{I}}] + \text{Var} [\mathbf{I}])}.$$

Conclusion

- Solid deep learning baseline to a novel problem
- Good result on the REDS dataset [1] which can be used as a benchmark
- Fast inference (reconstruction network 2.3M parameters)

Conclusion

- Solid deep learning baseline to a novel problem
- Good result on the REDS dataset [1] which can be used as a benchmark
- Fast inference (reconstruction network 2.3M parameters)

Possible Future Research

- Handle high memory consumption at training time
- Optimize the architecture of the reconstruction network

- [1] S. Nah, S. Baik, S. Hong, G. Moon, S. Son, R. Timofte, and K. M. Lee, "Ntire 2019 challenge on video deblurring and super-resolution: Dataset and study," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- [2] A. S. Kaplanyan, A. Sochenov, T. Leimkühler, M. Okunev, T. Goodall, and G. Rufo, "Deepfovea: neural reconstruction for foveated rendering and video compression using learned statistics of natural videos," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–13, 2019.
- [3] X. Tao, H. Gao, R. Liao, J. Wang, and J. Jia, "Detail-revealing deep video super-resolution," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 4472–4480.
- [4] X. Wang, K. C. Chan, K. Yu, C. Dong, and C. Change Loy, "Edvr: Video restoration with enhanced deformable convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.

- [5] A. Kappeler, S. Yoo, Q. Dai, and A. K. Katsaggelos, "Video super-resolution with convolutional neural networks," *IEEE Transactions on Computational Imaging*, vol. 2, no. 2, pp. 109–122, 2016.
- [6] R. Liao, X. Tao, R. Li, Z. Ma, and J. Jia, "Video super-resolution via deep draft-ensemble learning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 531–539.
- [7] Q. Dai, S. Yoo, A. Kappeler, and A. K. Katsaggelos, "Dictionary-based multiple frame video super-resolution," in *2015 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2015, pp. 83–87.
- [8] J. T. Barron, "A general and adaptive robust loss function," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4331–4339.

- [9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [10] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, "Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8934–8943.
- [11] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

Code Availability & Additional Resources



TECHNISCHE
UNIVERSITÄT
DARMSTADT



DeepFovea++: Reconstruction and Super-Resolution for Natural Foveated Rendered Videos

Christoph Reich
TU Darmstadt
christoph.reich@robot.tu-darmstadt.de
Marina Memmel
TU Darmstadt
marina.memmel@robot.tu-darmstadt.de
Jens-Henry Grebe
TU Darmstadt
jens-henry.grebe@robot.tu-darmstadt.de



Figure 1: Results of our proposed DeepFovea++ (best setting) framework. The first sampled input sequence of the resolution (102×256) image frames can be seen on the top sequence. It is shown in the middle, and the corresponding label sequence is shown at the bottom.

Abstract

Image super-resolution is a well-known problem in the field of computer vision. Recently, researchers extended the problem of super-resolution to videos and showed amazing results. On the other hand deep learning based methods have also been increasingly drawn more popularity since the DeepFovea publication of Facebook AI. Even though DeepFovea showed outstanding results, it was limited to the reconstruction of images of 128×128 pixels. We extend the proposed DeepFovea architecture to handle foveated video reconstruction and super-resolution ($102 \times 256 \rightarrow 768 \times 1024$) of scenes. Our proposed architecture, DeepFovea++, follows a two-stage approach. First, we propose a recursive U-Net architecture, and afterwards, the

desired super-resolution is learned by bidirectional convolution. We tested our DeepFovea++ architecture on the challenging BRDS dataset. The code is available at https://github.com/ChristophReich1996/DeepFoveaPP_for_Video_Reconstruction_and_Paper_Resolution.

1. Introduction

A full immersion with virtual reality requires a very high image resolution, a low latency, and a high frame rate. However, the human visual system (HVS) and the BrainLab Labs approached this problem by making use of the fact that human perception has different levels of visual acuity across the retina. This allows us to use a recursive U-Net architecture, and afterwards,

https://github.com/ChristophReich1996/DeepFoveaPP_for_Video_Reconstruction_and_Super_Resolution