

# NUWA-3D: Learning 3D Photography Videos via Self-supervised Diffusion on Single Images

## (Supplementary Material)

Xiaodong Wang<sup>1\*</sup>, Chenfei Wu<sup>2\*</sup>, Shengming Yin<sup>2</sup>, Minheng Ni<sup>2</sup>, Jianfeng Wang<sup>3</sup>, Linjie Li<sup>3</sup>, Zhengyuan Yang<sup>3</sup>, Fan Yang<sup>2</sup>, Lijuan Wang<sup>3</sup>, Zicheng Liu<sup>3</sup>, Yuejian Fang<sup>1†</sup>, Nan Duan<sup>2†</sup>

<sup>1</sup>Peking University <sup>2</sup>Microsoft Research Asia <sup>3</sup>Microsoft Azure AI  
{wangxiaodong21s@stu, fangyj@ss}.pku.edu.cn, {chewu, v-sheyin, t-mni, jianfw, Lindsey.Li, zhengyang, fanyang, lijuanw, zliu, nanduan}@microsoft.com,

## Appendix

### A.1 Metric Introduction

**For Novel View Synthesis:**

**LPIPS:** Learned Perceptual Image Patch Similarity (LPIPS) [Zhang *et al.*, 2018] measures image similarity based on deep features extracted by a learned network, and it better coincides with human judgment.

**PSNR:** Peak Signal-to-Noise Ratio (PSNR) measures image similarity with the pixel-wise independence assumption.

**SSIM:** Structural Similarity Index Measure (SSIM) is a perception-based method and measures image similarity by considering structural information.

**For Image Outpainting:**

**FID:** Fréchet Inception Distance (FID) [Heusel *et al.*, 2017] is used to measure the diversity and fidelity of the generated image.

**IS:** Inception Score (IS) [Salimans *et al.*, 2016] is a common metric to measure the fidelity of the generated image.

**CLIP-SIM:** CLIP Similarity Score (CLIP-SIM) [Radford *et al.*, 2021] is used to measure the semantic consistency between the generated image and text.

### A.2 Image filtering on MSCOCO

MSCOCO-2017 [Caesar *et al.*, 2018] contains 172 classes: 80 thing classes, 91 stuff classes, and 1 unlabeled class. We focus on thing classes and utilize pixel-level annotations. Specifically, we regard the regions belonging to thing classes as the input object regions, and models are supposed to outpaint the remaining regions. For each image, we utilize the pixel-level annotations to construct the binary mask to indicate the known objects and unknown regions, and captions are regarded as the text prompts. Finally, we filter the training set and the validation set to get an outpainting dataset containing 117266 training images and 4952 test images.

### A.3 More Samples of Out-animation

Figure. 1 shows more samples of our proposed out-animation task. Our method can handle a wide variety of scenes from the open domain.

### A.4 More Samples on MSCOCO

**Out-animated Samples.** Figure. 2. (a) presents the out-animated results using the processed input objects from MSCOCO.

**Outpainted Samples.** Figure. 2. (b) presents the out-painted results using the processed input objects from MSCOCO.

## References

- [Caesar *et al.*, 2018] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1209–1218, 2018.
- [Heusel *et al.*, 2017] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [Radford *et al.*, 2021] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, and Jack Clark. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [Salimans *et al.*, 2016] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.
- [Shih *et al.*, 2020] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8028–8038, 2020.
- [Zhang *et al.*, 2018] 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.

\*Equal contribution.

†Corresponding author.

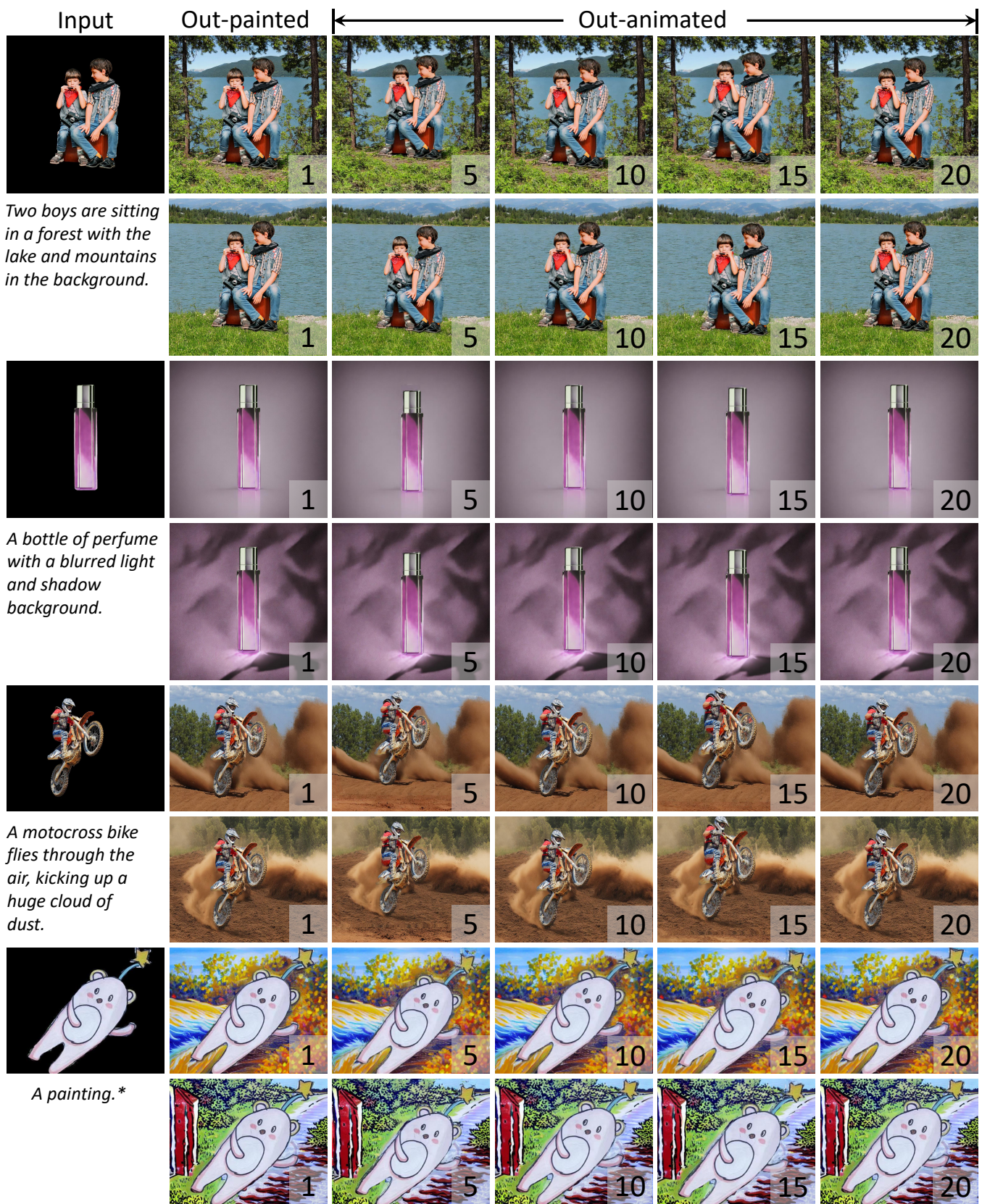
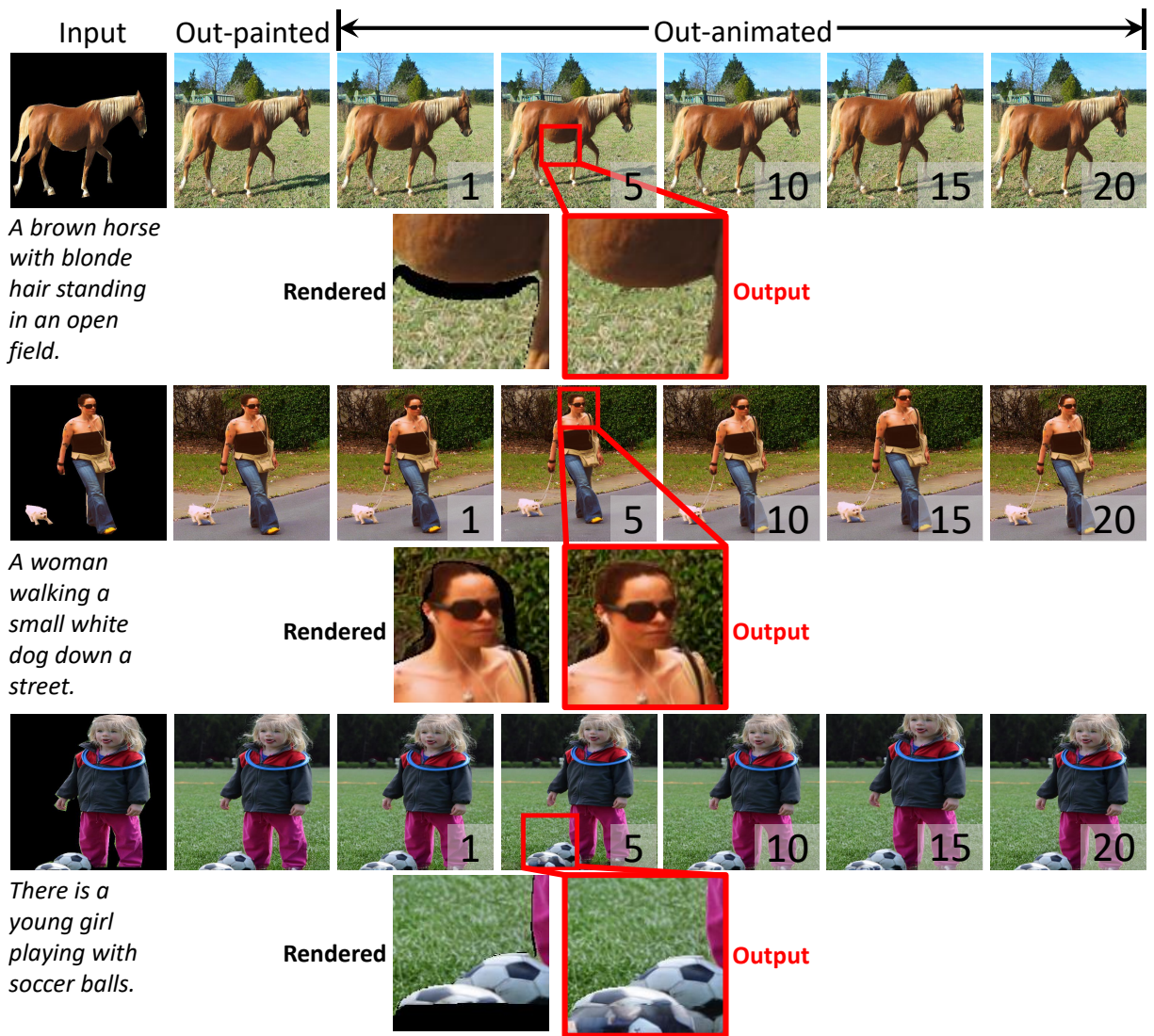
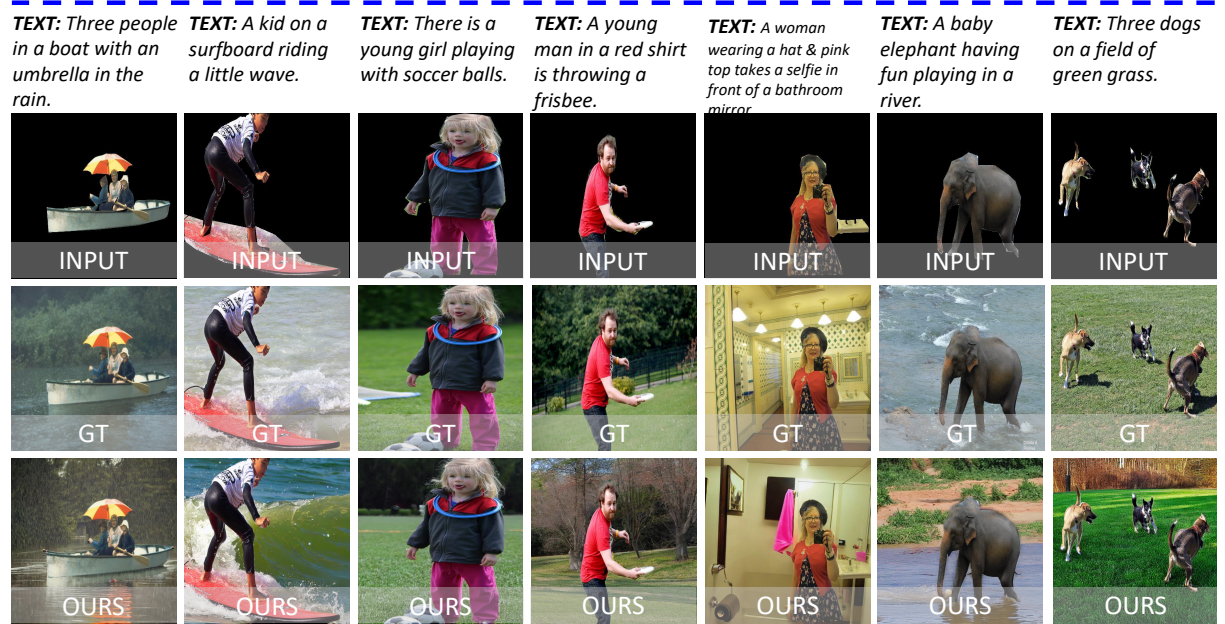


Figure 1: Illustration of our proposed out-animation task. We first outpaint the input into complete and suitable scenes based on its content and prompts, and then generate subsequent frames with 3D effects to form the out-animated videos. Our method can handle a wide variety of scenes from the open domain, such as people, objects, advertised goods, paintings, etc. (\*This input comes from an artwork of a young artist Geng Jiahao, in 2022 ANOBO “A Drop of Water with One World” exhibition. The number indicates the frame number in a 3D video.)



(a) More Samples for 3D Photography



(b) More Samples for Image Outpainting

Figure 2: More samples. (a) Out-animation sample results from MSCOCO. The 3D cycle effect follows 3d-photo [Shih *et al.*, 2020]. (b) Outpainted sample results from MSCOCO. From top to bottom: text prompt, input objects, ground truth, our prediction.