Supplementary Material: Self-Supervised Visibility Learning for Novel View Synthesis

Yujiao Shi¹, Hongdong Li¹, Xin Yu²

¹Australian National University and ACRV ²University of Technology Sydney

yujiao.shi@anu.edu.au, hongdong.li@anu.edu.au, xin.yu@uts.edu.au

1. Robustness to Source-View Permutations

In real world scenarios, the input source views are usually unordered. Thus, we take into account source-view permutation invariance when designing our framework. To demonstrate this, we randomly permute the input source views and feed them to the same trained model, denoted as "Ours (permuted)". Quantitative results are presented in Tab. 1 and qualitative visualizations are provided in Fig. 1.

From the results, it can be seen that there is negligible difference on the synthesized images with different input view orders. This demonstrates the robustness of our method to source-view permutations.

| Table 1: Robustness of our method to different input view orders |
|--|
|--|

| | Tanks and Temples | | | | | | | | | | | | DTU | | |
|-----------------|-------------------|---------|-------|--------|-------|-------|--------|---------|-------|------------|-------|-------|--------|-------|-------|
| | Truck | | | Train | | | M60 | | | Playground | | | | | |
| | LPIPS↓ | . SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ | LPIPS↓ | . SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ |
| Ours | 0.233 | 0.708 | 21.33 | 0.386 | 0.542 | 18.81 | 0.250 | 0.732 | 19.20 | 0.245 | 0.710 | 22.12 | 0.177 | 0.721 | 19.20 |
| Ours (permuted) | 0.233 | 0.708 | 21.34 | 0.385 | 0.542 | 18.81 | 0.250 | 0.732 | 19.13 | 0.246 | 0.710 | 22.13 | 0.177 | 0.721 | 19.20 |



Figure 1: Qualitative comparison of generated images by our method with different input view orders. The four examples are the same as those in Fig. 1 in the main paper.

2. Additional Visualization on the DTU Dataset

In this section, we present more visualization examples on the DTU dataset, as shown in Fig. 2. Scenes in the DTU dataset are simpler than those in the Tanks and Temples dataset, *i.e.*, there is always a single salient object in each image. However, there are large textureless regions and some of the object surfaces (*e.g.*, the second example in Fig. 2) are reflective.

For those regions, COLMAP almost fails to predict the depths, and thus the images synthesized by FVS [2] suffer severe artifacts. EVS estimates target-view depths and source-view visibility from the source view depths. When the source view depths are inaccurate and/or disagree with each other, the error will be accumulated to the final synthesized image. In contrast, our method does not rely on the accuracy of source view depths and is able to handle textureless and reflective regions. Thus our synthesized images are much closer to ground-truth images.



Figure 2: Additional qualitative visualization of generated results on the DTU dataset with six views as input.

3. Additional Visualization on the Source-View Visibility Estimation (SVE)

In this section, we provide more visualization results on our visibility-aware aggregation mechanism, which is indicated by Eq. (2) in the main paper. We choose the four examples presented in Fig. 1 in the main paper for this visualization. The results are illustrated in Fig. 3, Fig. 4, Fig. 5 and Fig. 6, respectively. The camera movements between source and target views are provided under each input source view image in the Figures.

4. Qualitative Illustration on the Proposed Soft Ray-Casting (SRC)

The proposed soft ray-casting (SRC) mechanism is one of the key components in our framework. It converts a multi-modal surface existence probability along a viewing ray to a single-modal depth probability.

We demonstrate its necessity by removing it from our whole pipeline, denoted as "Ours w/o ray-casting" in the main paper. Fig. 7 presents the qualitative comparison. For the depth visualization, we apply a softargmax to the surface probability distribution (in "Ours w/o ray-casting") or the depth probability distribution (in "Our whole pipeline"). From Fig. 7, it can be observed that the generated depth and image by "Ours w/o ray casting" are more blurred than "Our whole pipeline". This is mainly because the surface probability distribution is always multi-modal and it does not reflect pure depth information, demonstrating the effectiveness of the proposed soft ray-casting mechanism.

5. Residual Visualization Between Images Before and After Refinement

We visualize the aggregated images by Eq.(2) in the paper and the corresponding output images after the refinement network in Fig. 8. To show the differences, we also present the residual images between the aggregated images and final output images.

References

- Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In Proceedings of the IEEE International Conference on Computer Vision, pages 7781–7790, 2019. 2
- [2] Gernot Riegler and Vladlen Koltun. Free view synthesis. In European Conference on Computer Vision, pages 623–640. Springer, 2020. 2

 $Rot = [0.90^{\circ}, 1.73^{\circ}, -0.71^{\circ}]$

Trans = [-25.56cm, -2.50cm, 3.32cm]

 $Rot = [-0.95^{\circ}, -4.10^{\circ}, 1.67^{\circ}]$

 $Rot = [16.40^\circ, 1.85^\circ, 14.23^\circ]$ Trans = [42.91cm, -8.60cm, 121.41cm]

 $Rot = [15.83^{\circ}, 2.58^{\circ}, 13.22^{\circ}]$ Trans = [8.29cm, -1.48cm, 157cm]

 $Rot = [0.10^{\circ}, -2.71^{\circ}, 0.55^{\circ}]$ Trans = [-21.16cm, -0.34cm, 17.46cm]

Trans = [-24.79 cm, 2.41 cm, 198.15 cm]

Figure 4: Qualitative illustration on our visibility-aware aggregation (Train).

 $Rot = [18.20^{\circ}, 1.54^{\circ}, 15.74^{\circ}]$

(a) Input Source-View Images

(b) Warped Source-View Images

Figure 3: Qualitative illustration on our visibility-aware aggregation (Truck).

 $Rot = [-6.46^{\circ}, -8.52^{\circ}, -2.26^{\circ}]$ Trans = [-51.16cm, 22.23cm, 100.32cm]

 $Rot = [-1.34^{\circ}, -11.69^{\circ}, 1.28^{\circ}]$ Trans = [-23.63cm, -3.57cm, 57.15cm]

 $Rot = [-6.45^{\circ}, -11.41^{\circ}, -2.27^{\circ}]$ Trans = [-68.72cm, 23.02cm, 93.47cm]

 $Rot = [-0.41^{\circ}, -1.08^{\circ}, 0.94^{\circ}]$ Trans = [24.33cm, 30.47cm, -4.96cm]

(a) Input Source-View Images

(c) Target-View Ground Truth

 $Rot = [137.92^{\circ}, 5.13^{\circ}, -139.39^{\circ}]$ Trans = [56.66cm, 16.41cm, 211.42cm]

 $Rot = [137.13^{\circ}, 1.47^{\circ}, -138.50^{\circ}]$ Trans = [-36.27 cm, 28.48 cm, 120.76 cm]

 $Rot = [142.98^{\circ}, 13.25^{\circ}, -144.97^{\circ}]$ ${\rm Trans} = [0.29 {\rm cm}, 15.12 {\rm cm}, 210.44 {\rm cm}]$

 $Rot = [145.35^{\circ}, 16.37^{\circ}, -147.98^{\circ}]$ Trans = [-106.24 cm, -14.62 cm, 73.63 cm]

 $Rot = [28.33^{\circ}, -8.87^{\circ}, -28.85^{\circ}]$ ${\rm Trans} = [-33.01 {\rm cm}, -1.91 {\rm cm}, 8.70 {\rm cm}]$

 $Rot = [109.67^{\circ}, -7.22^{\circ}, -110.39^{\circ}]$ Trans = [53.16cm, 23.28cm, 128.66cm]

(a) Input Source-View Images

(b) Warped Source-View Images

Figure 5: Qualitative illustration on our visibility-aware aggregation (M60).

 $Rot = [-0.82^{\circ}, -2.05^{\circ}, 034^{\circ}]$

 $Rot = [-0.68^{\circ}, -1.72^{\circ}, 0.21^{\circ}]$

 $Rot = [-0.83^{\circ}, -4.84^{\circ}, 0.17^{\circ}]$ Trans = [15.00 cm, -6.87 cm, 171.10 cm]

 $Rot = [-0.62^{\circ}, -1.74^{\circ}, 0.08^{\circ}]$ Trans = [-25.01cm, -7.28cm, 166.32cm]

 $Rot = [-0.55^{\circ}, -2.53^{\circ}, -0.09^{\circ}]$ Trans = [-54.80 cm, -7.45 cm, 165.68 cm]

 $Rot = [-0.69^{\circ}, -10.34^{\circ}, -0.33^{\circ}]$ Trans = [-15.70cm, -6.60cm, 173.51cm]

(a) Input Source-View Images

(c) Target-View Ground Truth

Figure 6: Qualitative illustration on our visibility-aware aggregation (Playground).

(b) Warped Source-View Images

ground Truck (a) Generated depth (the first row) and image (the second row) by "Ours w/o ray-casting".

Train M60 (b) Generated depth (the first row) and image (the second row) by "Our whole pipeline".

Figure 7: Qualitative comparison between "Ours w/o ray-casting" and "Our whole pipeline".

Aggregated ImagesOutput ImagesResidual ImagesGround TruthsFigure 8: Visualization of aggregated images, output images after refinement, residual images between the aggregated images and the
output images, and the ground truths. The four examples are from scenes "Truck", "Train", "M60" and "Playground", respectively. Their
corresponding input images are presented in Fig.3, Fig.4, Fig.5, and Fig.6, respectively.