

Figure 2. Examples of the annotated regions in the challenging subsets of DTU: *Specular reflection* (row 1-2), *Shadow* (row 3-4), and *Occlusion* (row 5-6).

Net achieves high-quality reconstruction in various scenes.

To further verify the efficiency, we compare RayMVSNet against the baselines by visualizing the relationship between the overall accuracy of the reconstructed point cloud and the GPU memory consumption. As shown in Figure 4, RayMVSNet achieves state-of-the-art performance and requires less GPU memory compared to most of the baselines. This demonstrates that RayMVSNet is light weight, thanks to the mechanism of ray-based representation.

Last, we conduct experiments of replacing the MVSNet with other MVSNet variants, e.g., UCS-MVSNet, Fast-MVSNet, and CVP-MVSNet, for coarse depth estimation. We found consistent improvement of depth estimation for the alternative backbones. In particular, our method with a UCS-MVSNet backbone achieves a 0.326 overall score on the DTU dataset.

4. Explanation of quantitative comparison on Tanks & Temples

We compared RayMVSNet against the baselines on the Tanks & Temples dataset. Since the experiment is designed for evaluating the generality of the proposed method,

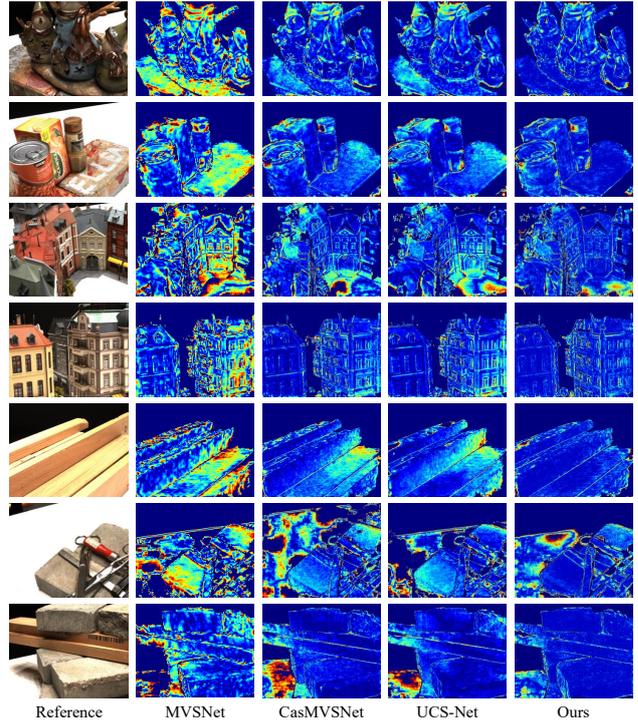


Figure 3. Visual comparison of the estimated depth map by RayMVSNet and the baselines.

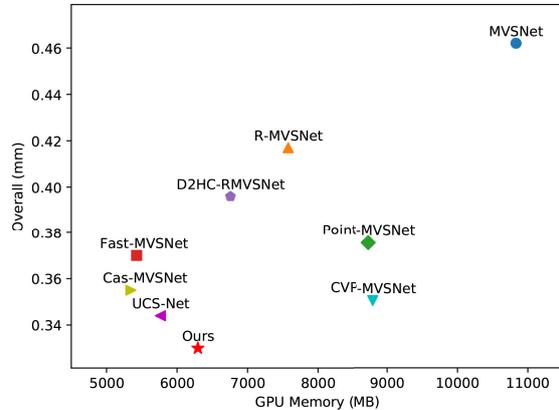


Figure 4. Visualizations of the overall reconstruction score and the GPU memory consumption. RayMVSNet achieves state-of-the-art performance and is light weight compared to most of the baselines.

we only compare RayMVSNet to existing learning-based methods that trained on the DTU dataset. Methods that have been fine-tuned on other datasets (e.g. Blended-MVS) are not considered. Those methods include Att-MVSNet [2], Vis-MVSNet [1], AA-RMVSNet [4], and EPP-MVSNet [3].

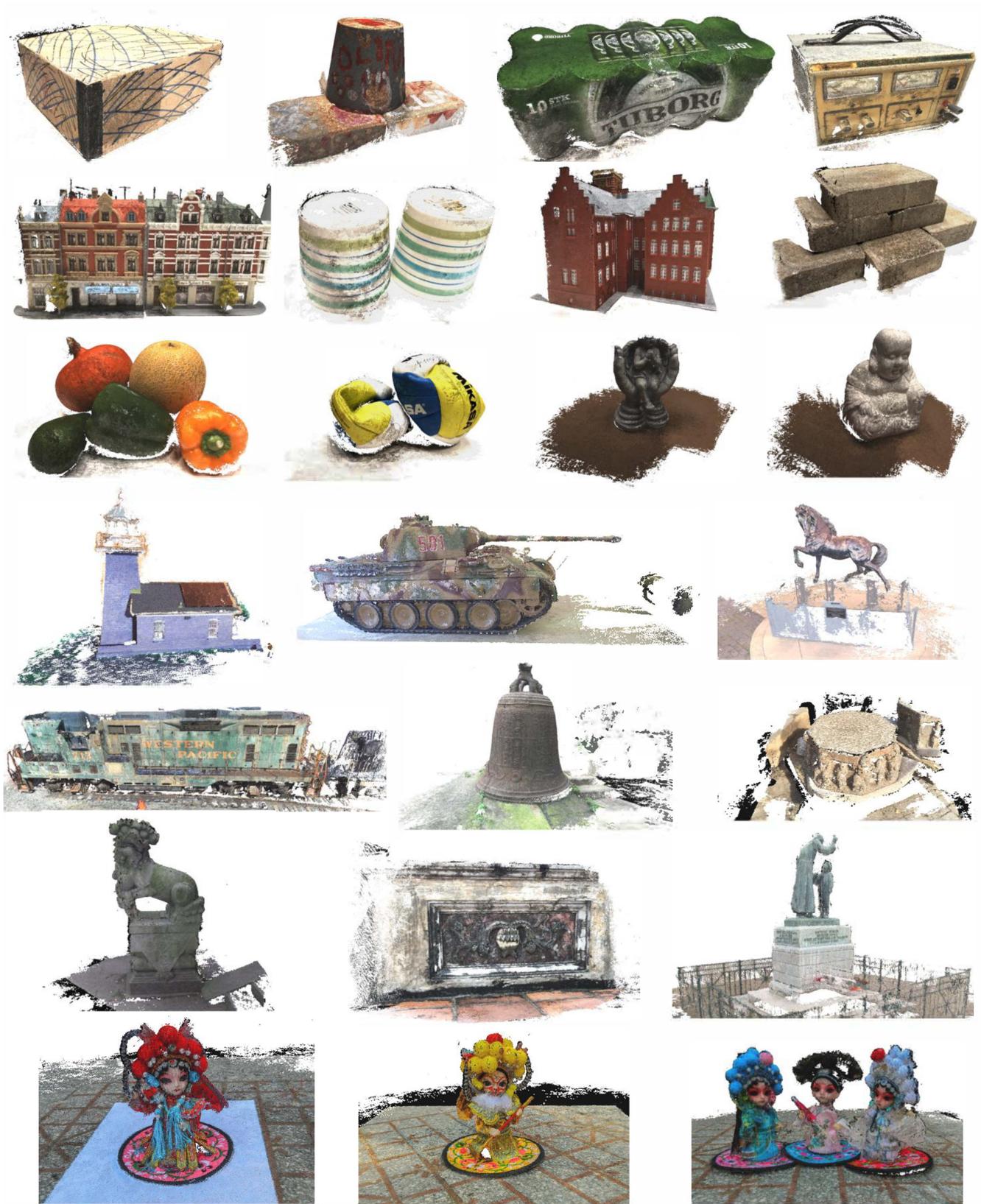


Figure 5. Visualizations of the reconstructed point cloud on DTU, Tanks & temples, and BlendedMVS.

References

- [1] Rui Chen, Songfang Han, Jing Xu, and Hao Su. Visibility-aware point-based multi-view stereo network. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3695–3708, 2020. [2](#)
- [2] Keyang Luo, Tao Guan, Lili Ju, Yuesong Wang, Zhuo Chen, and Yawei Luo. Attention-aware multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1590–1599, 2020. [2](#)
- [3] Xinjun Ma, Yue Gong, Qirui Wang, Jingwei Huang, Lei Chen, and Fan Yu. Epp-mvsnet: Epipolar-assembling based depth prediction for multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5732–5740, 2021. [2](#)
- [4] Zizhuang Wei, Qingtian Zhu, Chen Min, Yisong Chen, and Guoping Wang. Aa-rmvsnet: Adaptive aggregation recurrent multi-view stereo network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6187–6196, 2021. [2](#)