

# RoSAMDepth: Robust Self-supervised Depth Estimation Leveraging Segment Anything Model

## Supplementary Material

### 001 1. Methods Details

#### 002 1.1. SAM Mask Generation

003 SAM often produces a hierarchical set of masks for a single  
004 object, detecting not only the entire object but also its con-  
005 stituent parts (e.g., wheels, doors, windows, lights, and even  
006 bumpers for a car, as shown in Fig. 1). While detailed, this  
007 hierarchy introduces redundancy and potential noise for op-  
008 erations requiring object-level information. To address this,  
009 we adopt a simple yet effective strategy. We assume that in  
010 overlapping masks, the mask with the largest spatial extent  
011 represents the object instance, while smaller nested masks  
012 represent its parts.

013 Formally, let  $M_{sam}^{raw} \in \{0, 1\}^{K \times H \times W}$  denote the set of  
014  $K$  binary masks produced by SAM for an input image  $I$ ,  
015 where  $M_{sam,i}^{raw} \in \{0, 1\}^{H \times W}$  is the  $i$ -th mask. Specifically,  
016 we iterate through the masks with an area below a threshold  
017 (set to  $10^5$  pixels, approx.  $200 \times 500$ ). If a small mask is sig-  
018 nificantly covered by a larger mask (overlap ratio  $> 0.75$ ), it  
019 is considered a sub-part and is merged into the larger parent  
020 mask. As shown in Fig. 1, the resulting unified segmenta-  
021 tions serve as robust object-level priors, ensuring that the  
022 network achieves consistent and accurate depth predictions  
023 across the entire object instance.

#### 024 1.2. Training Framework

025 We build upon the Syn2Real-Depth framework [6] to in-  
026 vestigate how object-level priors from SAM [3] can en-  
027 hance robust self-supervised depth estimation in diverse  
028 real-world scenes. Syn2Real-Depth serves as an ideal base-  
029 line due to its real adaptation stage designed for diverse  
030 real data. The framework employs ManyDepth [5] as the  
031 depth backbone, originally defining three domain-specific  
032 models:  $\Phi_{day}$  (trained on real daytime data),  $\Phi_{syn}$  (trained  
033 on diverse synthetic data), and  $\Phi_{real}$  (trained on diverse  
034 real data). The pose network follows the standard archi-  
035 tecture described in [2, 5]. Since our contribution focuses  
036 on integrating SAM-based guidance into the real adaptation  
037 stage, we simplify the notation to better reflect our train-  
038 ing paradigm. Specifically, we designate  $\Phi_{syn}$  as the initial  
039 teacher  $\Phi_t$ , and  $\Phi_{real}$  as the student  $\Phi_s$ , which is trained on  
040 real data with object-level priors.

#### 041 1.3. Implementation Details

042 Our model is trained on a single NVIDIA 4090D GPU, tak-  
043 ing about 10 hours. Regarding the hyperparameters, we set  
044 the temperature  $T = 1.0$  and the number of feature scales

$S = 3$  for  $L_{src}$ . For the temperature-controlled sharpening  
sigmoid function  $S_\kappa(\cdot)$ , the sharpening temperature is set to  
 $\kappa = 15.0$ . In the object-level reliability estimation (ORE)  
strategy, we employ a scaling factor  $\beta = 1.0$  and a bias  
term  $\epsilon = 1.0$ . Finally, the balancing weights for the total  
objective  $L_{total}$  are set to  $\lambda_1 = 0.05$  and  $\lambda_2 = 0.2$ .

#### 051 1.4. Supervision Details

052 To ensure robust performance under adverse conditions,  
053 our training on real data inherits the auxiliary losses from  
054 Syn2Real-Depth [6]. Specifically, we incorporate the ex-  
055 ternal loss term  $L_{ext}$ , which comprises the K-L divergence  
056 ( $L_{dis}^{KL}$ ), the cost volume learning loss ( $L_{cv}$ ), and the su-  
057 pervision for pose ( $L_T$ ), following the exact formulation  
058 in [6]. We retain these components unchanged to lever-  
059 age their proven efficacy in the real data. This allows us  
060 to isolate the source of performance gains, demonstrating  
061 that our proposed object-level integration provides additive  
062 improvements over a strong baseline.

### 063 2. More Qualitative Examples

064 As supplementary to the main paper, we provide more qual-  
065 itative results on nuScenes [1] (day-clear, night, day-rain),  
066 Robotcar [4] (daytime, nighttime), as shown in Figs. 2  
067 and 3.



Figure 1. Example of the SAM masks generation.

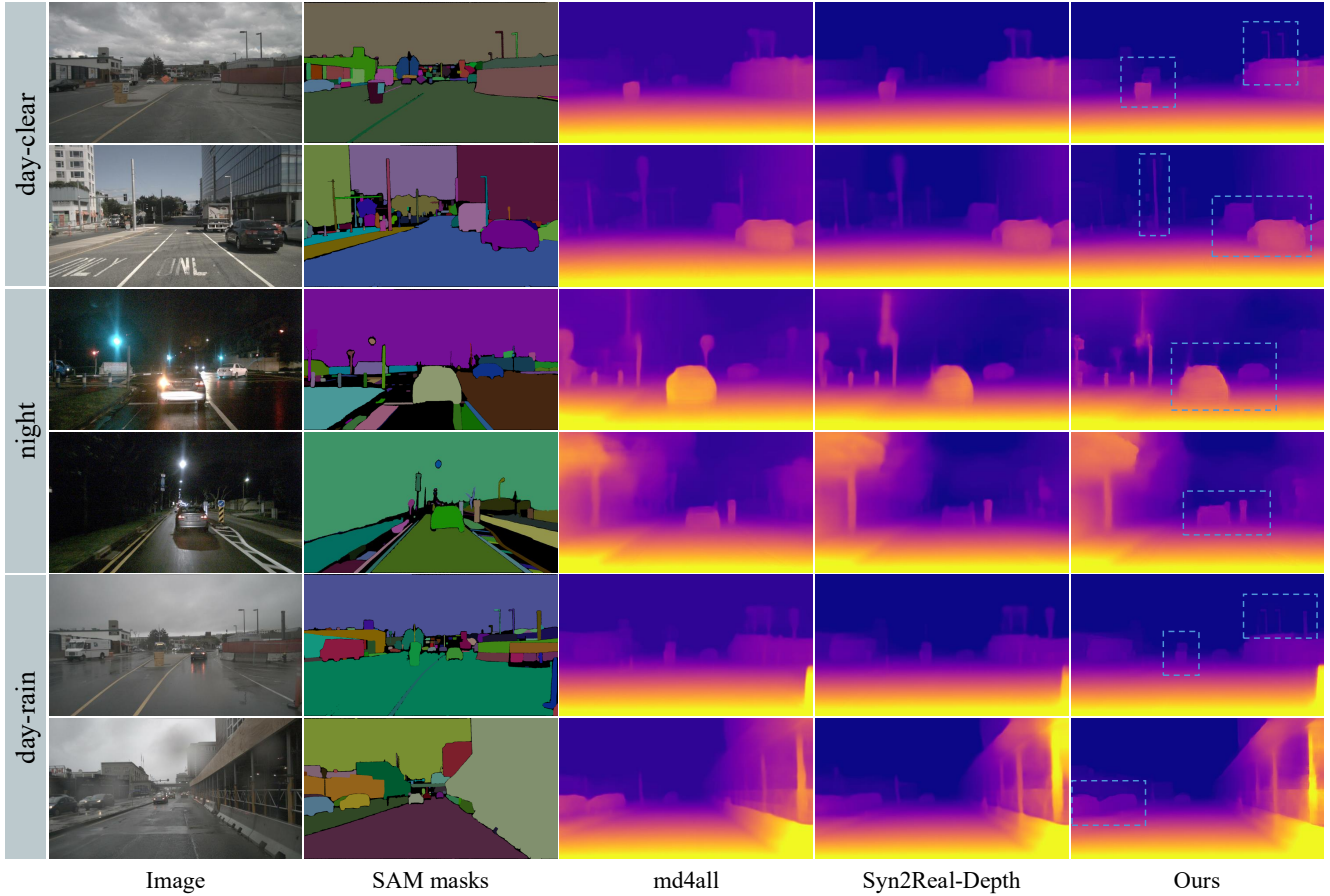


Figure 2. More qualitative results on nuScenes [1]

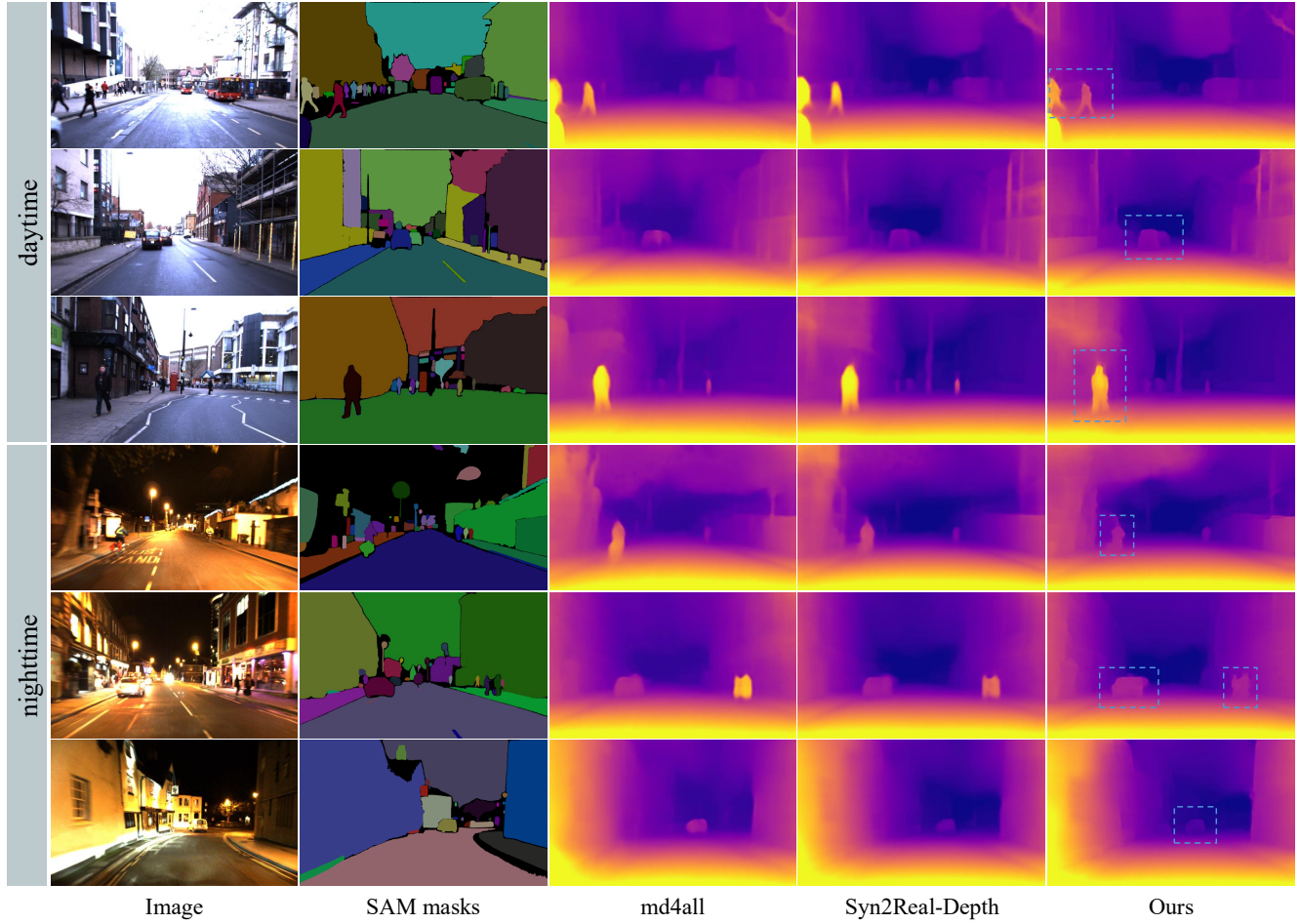


Figure 3. More qualitative results on Robotcar [4]

068

**References**069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098

- [1] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 1, 2
- [2] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3828–3838, 2019. 1
- [3] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 1
- [4] Will Maddern, Geoffrey Pascoe, Chris Linegar, and Paul Newman. 1 year, 1000 km: The oxford robotcar dataset. *The International Journal of Robotics Research*, 36(1):3–15, 2017. 1, 3
- [5] Jamie Watson, Oisín Mac Aodha, Victor Prisacariu, Gabriel Brostow, and Michael Firman. The temporal opportunist: Self-supervised multi-frame monocular depth. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1164–1174, 2021. 1
- [6] Weilong Yan, Ming Li, Haipeng Li, Shuwei Shao, and Robby T Tan. Synthetic-to-real self-supervised robust depth estimation via learning with motion and structure priors. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 21880–21890, 2025. 1