

# STEPS: Joint Self-supervised Nighttime Image Enhancement and Depth Estimation

Yupeng Zheng<sup>1,2</sup>, Chengliang Zhong<sup>3,5</sup>, Pengfei Li<sup>3</sup>, Huan-ang Gao<sup>3</sup>, Yuhang Zheng<sup>2</sup>,  
Bu Jin<sup>1,2</sup>, Ling Wang<sup>3,5</sup>, Hao Zhao<sup>4</sup>, Guyue Zhou<sup>2</sup>, Qichao Zhang<sup>1</sup> and Dongbin Zhao<sup>1</sup>

**Abstract**—Self-supervised depth estimation draws a lot of attention recently as it can promote the 3D sensing capabilities of self-driving vehicles. However, it intrinsically relies upon the photometric consistency assumption, which hardly holds during nighttime. Although various supervised nighttime image enhancement methods have been proposed, their generalization performance in challenging driving scenarios is not satisfactory. To this end, we propose the first method that jointly learns a nighttime image enhancer and a depth estimator, without using ground truth for either task. Our method tightly entangles two self-supervised tasks using a newly proposed uncertain pixel masking strategy. This strategy originates from the observation that nighttime images not only suffer from underexposed regions but also from overexposed regions. By fitting a bridge-shaped curve to the illumination map distribution, both regions are suppressed and two tasks are bridged naturally. We benchmark the method on two established datasets: nuScenes and RobotCar and demonstrate state-of-the-art performance on both of them. Detailed ablations also reveal the mechanism of our proposal. Last but not least, to mitigate the problem of sparse ground truth of existing datasets, we provide a new photo-realistically enhanced nighttime dataset based upon CARLA. It brings meaningful new challenges to the community. Codes, data, and models are available at <https://github.com/ucaszyp/STEPS>.

## I. INTRODUCTION

Pointcloud-based sensing algorithms are of great significance to computer vision society [2]–[13]. However, they are widely considered too expensive for autonomous vehicles. In this regard, image-based depth estimation has drawn a lot of attention from the robotics community [14]–[19] due to low hardware cost. Amongst learning-based depth estimation methods, self-supervised formulations using image sequences [20], [21] are quite appealing as they do not require paired RGB-D data and open up the opportunity for online adaptation [22]. And with the efforts of [23]–[26], the performance of self-supervised depth estimation on KITTI [27], Cityscapes [28], and DDAD [24] datasets is already comparable to supervised methods. However, these studies all focus on daytime image sequences where inputs are well-lit and the photometric consistency assumption generally holds. Self-driving vehicles need to run robustly during nighttime and unfortunately the photometric consistency assumption hardly holds in this challenging scenario.

<sup>1</sup>Institute of Automation, Chinese Academy of Sciences, China, {zhengyupeng2022, jinbu2022, zhangqichao2014, dongbin.zhao}@ia.ac.cn.

<sup>2</sup>Institute for AI Industry Research (AIR), Tsinghua University, China.

<sup>3</sup>Department of Computer Science and Technology, Tsinghua University, China, {gha20, li-pf22}@mails.tsinghua.edu.cn.

<sup>4</sup>Intel Labs China, Peking University, China, hao.zhao@intel.com

<sup>5</sup>Xi’an Research Institute of High-Tech, China.

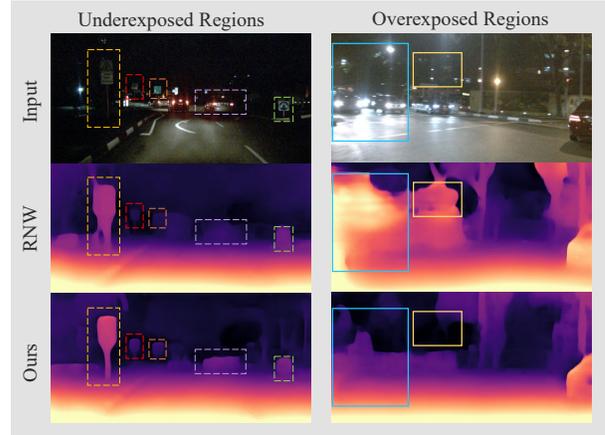


Fig. 1. The first row is the input image of two scenes, and the colored dotted and solid boxes mark the objects in underexposed regions (left) and overexposed regions (right). The last two rows are the depth maps predicted by RNW [1] and our method respectively, indicating the effectiveness of image enhancement and uncertain pixel masking in underexposed and overexposed regions.

A natural idea is to use nighttime image enhancement techniques to improve the quality of input images. But supervised nighttime image enhancers are intrinsically restricted by dataset bias as existing paired day/night datasets focus on indoor scenes and building these kinds of paired datasets for dynamic road scenes is nearly impossible. To this end, we propose the first learning framework that learns an enhancer and a single-view depth estimator both in a self-supervised manner, as shown in Fig. 2. Since these two modules can collaboratively work towards a common goal without any ground truth, our method significantly outperforms the previous SOTA method RNW [1] as shown in Fig. 1.

Delving into this framework, we identify an interesting but overlooked fact: nighttime images suffer not only from underexposed regions but also from overexposed regions (referred to as *unexpected regions* as a whole later). Both cause some detailed information loss and prevent the model from estimating accurate depth through local contextual cues. Moreover, the overexposed regions are often associated with the movement of cars (e.g., car light), which also violates photometric consistency. In the spirit of influential auto-masking techniques in the self-supervised depth estimation literature [23] [29], we strive to softly mask out unexpected regions.

We observe that the mid-product of image enhancement – enhancement ratio (or say the illumination component)

could provide hints to find the unexpected regions, *i.e.*, underexposed areas need a higher ratio and vice versa. This observation motivates us to design an uncertainty map built on the ratio to suppress unexpected areas. After that, we bridge depth estimation and image enhancement tightly, using a bridge-shaped model (Fig. 3) for soft masking. Apart from that, we also introduce a pre-trained denoising module to increase the image signal-to-noise ratio.

Furthermore, we find that the existing night driving datasets have only sparse ground truth due to the limit of LiDAR data which cannot cover all areas of interest during evaluation. Following the idea that transferring the knowledge in the simulation environment to the real world, we resort to CARLA [30], the simulator for autonomous driving research. However, huge domain gap between the rendered images and the real-world images makes it not straightforward to use the simulated data directly. Thus, we propose CARLA-EPE, a photo-realistically enhanced nighttime dataset based upon CARLA. We leverage the state-of-the-art photorealism enhancement network EPE [31] to transfer the style of the rendered images to a photorealistic one, resulting in a photorealistic nighttime dataset CARLA-EPE with dense ground truth of depth. From the experiment results, the task in our new dataset is more challenging than others, which brings meaningful new challenges to the field.

In brief, our contributions can be summed up in four-fold:

- We propose the first method that jointly learns nighttime image enhancement and depth estimation without using any paired ground truth.
- We identify that the illumination component in the self-supervised nighttime image enhancer can be used to identify unexpected regions and propose a bridge-shape model for soft auto-masking.
- We contribute a novel photorealistically enhanced nighttime dataset with dense depth ground truth.
- We achieve SOTA performance on public benchmarks and release our codes.

## II. RELATED WORKS

### A. Self-supervised Monocular Depth Estimation

Various methods have studied self-supervised depth estimation in computer vision owing to its less supervision demanded. SfMLearner [32] is the first work to address this problem by jointly learning the depth and relative pose between two adjacent frames to reconstruct the target frame with photometric consistency loss. Since the photometric constraint is easily affected by the blur, occlusion, and moving objects, several works proposed effective approaches, such as optical flow [33], instance segmentation [34], stationary pixel mask [23], point cloud consistency [35], and packing network [24]. However, these methods are designed for daytime depth estimation and fail to work in nighttime scenarios. The biggest challenge of nighttime depth estimation is that the assumption of temporal illumination consistency is invalid due to the low-light and non-uniform illuminations. Recently, a few works have explored it. ADDS

[36] proposed a domain-separated network to extract illumination and texture features for both day-time and nighttime depth estimation networks. Different from ADDS, RNW [1] trained on unpaired data by (1) leveraging a GAN-based method to adapt the daytime estimation network to apply for the nighttime data, (2) employing depth distribution from daytime to regularize the nighttime network, and (3) using a HE-based offline image enhancement module to deal with the low-light regions. Despite their better results, these models still suffer from underexposed and overexposed regions.

### B. Low-light Image Enhancement

Low-light image enhancement aims to restore details in low visibility regions. Histogram equalization (HE) [37], including its variants CLHE [37], is a classical image enhancement method that increases the global contrast of images. However, the HE-based approaches are based upon global distribution rather than local context, so the useable signal may be reduced while background noise contrast is increased. Retinex model-based methods [38]–[41] are another alternative ways for image enhancement. It assumes the low-light image can be decomposed into illumination and reflectance. Moreover, the reflectance is regarded as the result of enhancement. To better leverage the great power of deep learning, network-based Retinex method [42]–[44] combined CNNs and Retinex theory to pursue better accuracy and robustness. More recently, SCI [45] built a cascaded illumination learning process and developed a self-supervised framework. Nevertheless, their generalization performance in challenging driving scenarios is not satisfactory. In addition, a pre-trained image enhancement network is not necessarily suitable for nighttime depth estimation.

## III. METHOD

### A. Self-supervised Nighttime Depth Estimation

Given a single image  $I_t \in \mathbb{R}^{H \times W \times 3}$ , the goal of learning-based depth estimation is to predict a depth map  $D_t \in \mathbb{R}^{H \times W}$  by a trainable network  $\Phi_d : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^{H \times W}$ . To achieve self-supervision, the key idea is to reconstruct the target frame  $I_t$  from the source one  $I_s$  according to the geometry constraint. To be specific, given a known camera intrinsic matrix  $K$ , a predicted depth map  $D_t$  and relative pose  $P_{t \rightarrow s} \in \mathbb{R}^{4 \times 4}$  between source and target frames via a trainable network  $\Phi_p : \mathbb{R}^{H \times W \times 6} \rightarrow \mathbb{R}^{4 \times 4}$ , each point  $p_t$  in  $I_t$  can be projected onto the source view  $p_s$  in  $I_s$  by

$$p_s \sim K P_{t \rightarrow s} D_t(p_t) K^{-1} p_t, \quad (1)$$

where  $\sim$  represents the homogeneous equivalence. With Eq. 1, we can recover the target frame  $\hat{I}_t$  from  $I_s$  by:

$$\hat{I}_t = I_s \left\langle \text{proj}(D_t, P_{t \rightarrow s}, K) \right\rangle, \quad (2)$$

where  $\langle \cdot \rangle$  is a differentiable bilinear interpolation proposed by [46] and  $\text{proj}(\cdot, \cdot, \cdot)$  is the projection operation in Eq. 1.

The training signal is the photometric error between the target frame  $I_t$  and the reconstructed frame  $\hat{I}_t$ . Following

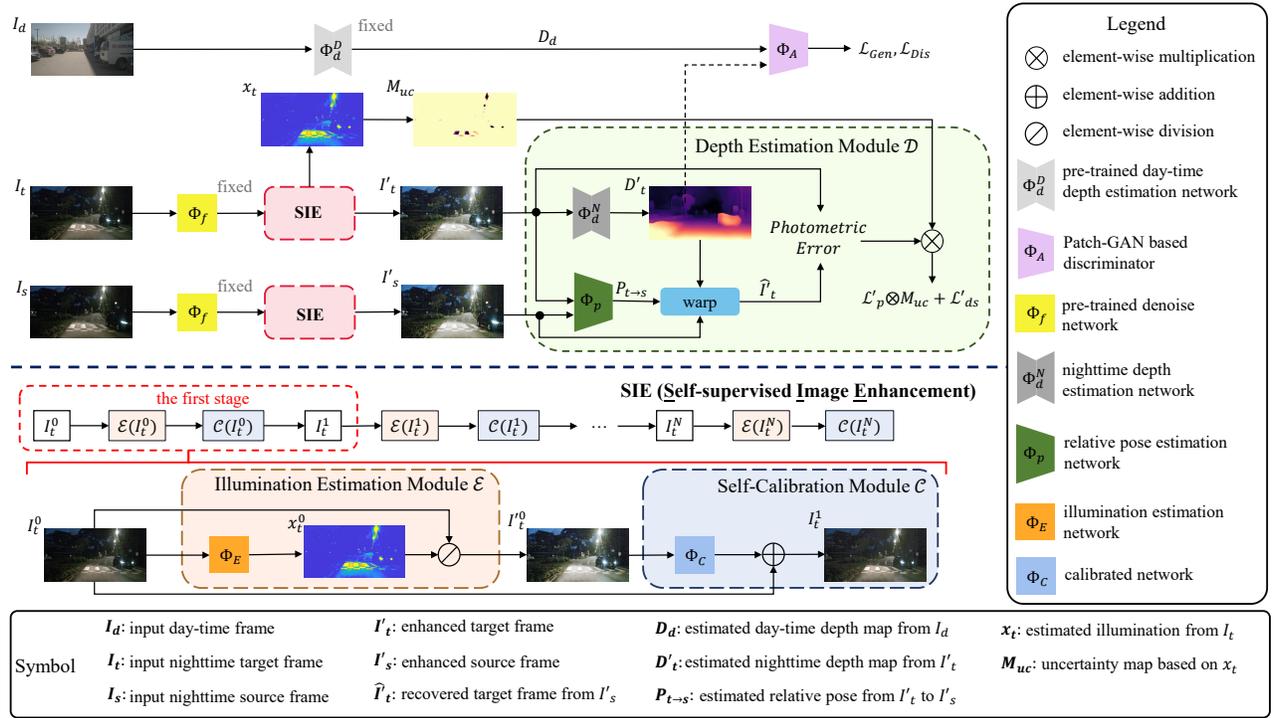


Fig. 2. The architecture of our proposed framework (see section III for details)

[23], we combine  $\ell_1$  and SSIM losses as the photometric loss  $L_p$ , which is defined as:

$$\mathcal{L}_p(I_t, \hat{I}_t) = \alpha \frac{(1 - \text{SSIM}(I_t, \hat{I}_t))}{2} + (1 - \alpha) |I_t - \hat{I}_t|_1, \quad (3)$$

where we set  $\alpha = 0.85$  for all experiments as in [23]. Additionally, as stated in [23], minimizing Eq. 3 can only enforce a necessary but not sufficient condition. Therefore, we follow [23] to avoid depth ambiguity by enforcing smoothness of the predicted depth map, *i.e.*

$$\mathcal{L}_{ds}(D_t, I_t) = |\partial_x D_t| e^{-|\partial_x I_t|} + |\partial_y D_t| e^{-|\partial_y I_t|}. \quad (4)$$

Due to the poor quality of nighttime images, training the learning system with Eq. 3 can only provide noisy gradients. To alleviate that, we follow [1] to introduce a pre-trained day-time depth estimation model as prior to direct the nighttime model training via an adversarial manner. As shown in Fig. 2, we build a nighttime depth estimation network  $\Phi_d^N$  as a generator, aiming to make its prediction  $D_t$  indistinguishable with  $D_d$  which is the output of a pre-trained and fixed day-time depth estimation network  $\Phi_d^D$ . A Patch-GAN-based discriminator  $\Phi_A$  is a trainable network to distinguish  $D_t$  and  $D_d$ .  $\Phi_d^N$  and  $\Phi_A$  are trained by minimizing the GAN-based loss functions, which are formulated as:

$$\mathcal{L}_{\text{Dis}} = \frac{1}{2|I_d|} \sum_{D_d} (\Phi_A(D_d) - 1)^2 + \frac{1}{2|I_t|} \sum_{D_t} (\Phi_A(D_t))^2, \quad (5)$$

$$\mathcal{L}_{\text{Gen}} = \frac{1}{2|I_t|} \sum_{D_t} (\Phi_A(D_t) - 1)^2, \quad (6)$$

where  $|I_d|$  and  $|I_t|$  are the number of day-time and nighttime training images,  $D_d = \Phi_d^D(I_d)$  and  $D_t = \Phi_d^N(I_t)$ .

## B. Joint Training Framework

As discussed before, nighttime image enhancement could improve the quality of input images to help depth estimation. But supervised nighttime image enhancers are intrinsically restricted by dataset bias. Hence, we propose a framework to jointly train depth estimation and image enhancement (SIE) in a self-supervised manner, as shown in Fig. 2.

According to the Retinex theory [38], given a low-light image  $I_t$ , the enhanced image can be obtained by  $I'_t = I_t \odot x$ , where  $x$  is the illumination map which is the most vital part of image enhancement. An inaccurate illumination estimation may bring an over-enhancement result. In order to improve performance stability and decrease computational burden, we follow the stage-wise illumination estimation with a self-calibrated module structure from SCI [45], as shown in the bottom of Fig. 2. The enhancement process is formulated by

$$I_t^{n+1} = \mathcal{C}(\mathcal{E}(I_t^n)), \quad (7)$$

where  $n$  ( $0 < n < 3$ ) is the stage,  $\mathcal{E}$  and  $\mathcal{C}$  represent illumination estimation and calibration module respectively. For stage  $n$ ,  $\mathcal{E}$  and  $\mathcal{C}$  are implemented by

$$\mathcal{E}(I_t^n) = \begin{cases} x_t^n = \Phi_E(I_t^n) \\ I_t^n = I_t^n \oslash x_t^n, \end{cases} \quad (8)$$

$$\mathcal{C}(I_t^n) = \begin{cases} \text{res}_t^n = \Phi_C(I_t^n) \\ I_t^{n+1} = I_t^n \oplus \text{res}_t^n, \end{cases} \quad (9)$$

where  $\Phi_E$  and  $\Phi_C$  are the trainable networks to estimate illumination  $x_t^n$  and generate calibrated residual map  $\text{res}_t^n$  respectively.  $\Phi_E$  and  $\Phi_C$  share the same parameters in each stage. The calibration module re-generates a pseudo nighttime image so that SIE can be applied in several stages

and empirically calibration brings faster convergence and better enhancement. The enhancement loss contains fidelity and smoothness loss, formulated as

$$\mathcal{L}_f = \frac{1}{|I_t|} \sum_{I_t} \sum_n \|x_t^n - I_t^n\|_2, \quad (10)$$

$$\mathcal{L}_{es} = \frac{1}{|I_t|} \sum_{I_t} \sum_n \sum_i \sum_{j \in \mathcal{W}(i)} \kappa_{i,j} |x_t^n(i) - x_t^n(j)|, \quad (11)$$

where  $\kappa_{i,j}$  is weight of a gaussian kernel,  $\mathcal{W}(i)$  is a window centered at  $i$  with  $5 \times 5$  adjacent pixels, and  $x(i)$  means the pixel value of  $x$  at  $i$ . The insight behind  $L_f$  is that for nighttime images, the illumination component is largely similar to the input image. Meanwhile  $L_{es}$  is a consistency regularization loss.

In Fig. 2, for joint training, the enhanced result, *i.e.*, the first stage output of SIE, is the input of  $\Phi_d$  and  $\Phi_p$ . Besides, during the warping (Eq. 1) and loss computing (Eq. 3, 4), the target frame  $I_t$  and the recovered frame  $\hat{I}_t$  are replaced by the enhanced frame  $I'_t$  and  $\hat{I}'_t$  respectively. We denote the modified losses in Eq. 3, 4 as  $\mathcal{L}'_p(I'_t, \hat{I}'_t)$  and  $\mathcal{L}'_{ds}(D'_t, I'_t)$ .

### C. Statistics-Based Pixel-wise Uncertainty Mask

Nighttime images often contain overexposed or underexposed regions where important details would be lost, as shown in Fig. 1. It extremely breaks the model from predicting accurate depth via local contextual cues. Moreover, the overexposed regions are often associated with the movement of cars (e.g., car light), which also violates the illumination consistency in Eq. 3. Therefore, we need to design a certain mechanism to bypass such regions. The success of Monodepth2 [23] demonstrates that the masking strategy is a simple and effective way to filter out regions that do not meet the assumption of photometric consistency. Inspired by them, we strive to softly mask out unexpected regions.

Recall that we have used an image enhancement module, SIE. It can predict an illumination map  $x_t$  to determine the enhancement ratio of the color of each pixel. As shown in Fig. 2, the ratio tends to be large in the underexposed regions and small in the overexposed regions. If we use this to weigh the importance of each pixel in photometric loss, the unexpected regions, as mentioned above, are more likely to be filtered. Specifically, we define an uncertainty map  $M_{uc} \in \mathbb{R}^{H \times W}$ , which gives low confidence in the unexpected regions and high confidence in the reasonable regions.  $M_{uc}$  is built upon  $x_t$  and formulated as

$$M_{uc} = \begin{cases} \frac{1}{1+p^2(x-a)^2}, & x_{min} \leq x \leq a \\ 1, & a \leq x \leq b \\ \frac{1}{1+q^2(x-a)^2}, & b \leq x \leq x_{max} \end{cases} \quad (12)$$

where  $a$  and  $b$  are the statistics-based bound of illumination to filter reasonable areas,  $p$  and  $q$  are the attenuation coefficients. Intuitively, this function looks like a bridge as shown in Fig.3. Suggested values of  $a$ ,  $b$ ,  $p$  and  $q$  can be found in the code release.

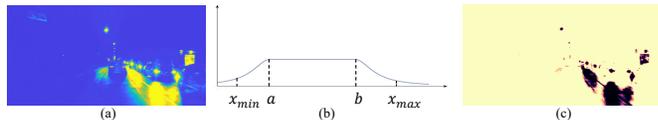


Fig. 3. This figure shows the principle of  $M_{uc}$ . (a) is the illumination map  $x_t$ . (b) shows the  $M_{uc}$  function, which could softly mask out unexpected regions. (c) is the visualization of uncertainty mask.

### D. Image Denoising

Image denoising is another useful component in the nighttime depth estimation since the low-light images captured by the sensor in night scenes contain much more noise. Especially after the low-light image is enhanced, the noise is inevitably amplified. It could affect the performance of the training signal as it further breaks the illumination consistency between adjacent frames. We use a network AP-BSN [47] with the model pre-trained on the SIDD [48] dataset. In order to decrease the training burden, the weights of the denoising network AP-BSN are fixed during training.

### E. Full Pipeline

In summary, for SIE module, the total loss  $\mathcal{L}_{SIE}$  is formulated by

$$\mathcal{L}_{SIE} = \beta \mathcal{L}_f + \gamma \mathcal{L}_{es}, \quad (13)$$

and for depth estimation module, the total loss  $\mathcal{L}_{DE}$  is formulated by

$$\mathcal{L}_{DE} = \lambda \mathcal{L}'_p \otimes M_{uc} + \mu \mathcal{L}'_{ds}. \quad (14)$$

The total loss of the whole pipeline is defined as

$$\mathcal{L}_{total} = \eta \mathcal{L}_{SIE} + \zeta \mathcal{L}_{DE} + \xi \mathcal{L}_{Gen} + \rho \mathcal{L}_{Dis}. \quad (15)$$

The  $\beta, \gamma, \lambda, \mu, \eta, \zeta, \xi$ , and  $\rho$  are hyper-parameters.

## IV. EXPERIMENTS

### A. Dataset

**CARLA-EPE.** Existing datasets for road scene depth estimation exploit LiDAR to acquire the groundtruth depth, which can only generate sparse depth maps and require a higher cost. The sparse ground truth may not reveal the overall performance of depth estimation methods. Although the RGB image and the corresponding dense depth map can be easily collected in the simulator (e.g., CARLA [30]), the large domain gap between the simulated and the real image dramatically affects the application of the trained model in real scenes. Hence, we propose a nighttime depth estimation dataset based on CARLA and the enhancing photorealism enhancement (EPE) network [31], which can provide dense depth ground truth and photo-realistic images. The pipeline of the dataset generation is as follows. Inspired from [49] and [50], we first capture rendered images as well as intermediate rendering buffers (G-buffers), which contains geometry, materials and motion information shown in Fig. 4 (a) from CARLA simulator. Then we find matches between rendered dataset and the DarkZurich dataset [51] taken from the real-world with semantic label predicted by DANNet [52]. After that, we train the EPE [31] network to transfer the render

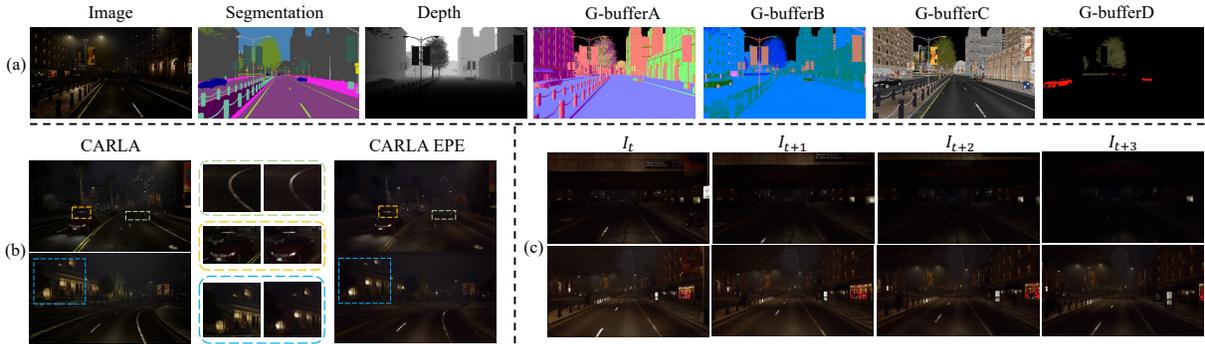


Fig. 4. G-buffers and CARLA-EPE Dataset. (a) shows G-buffers of the rendered image from CARLA, representing geometry (Depth, G-bufferA), materials (G-bufferB, G-bufferC) and motion (G-bufferD). (b) is the comparison of the original CARLA image and the corresponding CARLA-EPE image. The dashed box clearly demonstrates the enhancement of details. The orange and green dotted boxes in the middle show the details of the specular highlights from the car roof and lane lines under the street lights, and the blue dotted boxes show the more realistic lighting in the night scene. (c) contains two continuous sequences in CARLA-EPE, which exhibits drastic lighting changes between adjacent frames.

TABLE I

COMPARISON WITH OTHER METHODS. LOWER IS BETTER FOR ERROR AND HIGHER IS BETTER FOR ACCURACY.

Method	Max Depth	Error				Accuracy		
		Abs Rel	Sq Rel	RMSE	RMSE log	a1	a2	a3
<b>nuScenes</b>								
MonoDepth2 [23]	60 m	1.185	42.306	21.613	1.567	0.184	0.360	0.504
RNW [1]	60 m	0.326	3.999	9.932	0.417	0.492	0.765	0.870
Ours	60 m	<b>0.292</b>	<b>3.363</b>	<b>9.120</b>	<b>0.390</b>	<b>0.572</b>	<b>0.805</b>	<b>0.908</b>
<b>RobotCar</b>								
MonoDepth2 [23]	60 m	0.580	21.446	12.771	0.521	0.552	0.840	0.920
DeFeat-Net [55]	60 m	0.335	4.339	9.111	0.389	0.603	0.828	0.914
ADFA [56]	60 m	0.233	3.783	10.089	0.319	0.668	0.884	0.924
ADDS [36]	60 m	0.231	2.674	8.800	0.268	0.620	0.892	0.956
RNW [1]	60 m	0.185	1.894	7.319	0.246	0.735	0.910	0.965
Ours	60 m	<b>0.170</b>	<b>1.686</b>	<b>6.797</b>	<b>0.234</b>	<b>0.758</b>	<b>0.923</b>	<b>0.968</b>

images to realistic styles of DarkZurich. As for depth, we use the depth channel, one of the G-buffers we extract. We finally generate 12,000 pairs of nighttime images and dense depth maps. Fig. 4 (b) shows some comparisons of original and enhanced images. It shows that specular highlights on cars and lane lines to cars and lane lines under street lights are visible and building lighting is more realistic at night. This significantly increases the realism of the rendered image. Fig. 4 (c) shows two image sequences whose brightness varies greatly between adjacent frames due to changes in lighting, which is similar to real night scenes.

**nuScenes [53].** nuScenes is a large autonomous driving dataset comprising 1000 video clips collected in diverse road scenes and weather conditions. These scenes are pretty challenging, with a fair amount of unexpected regions.

**RobotCar [54].** RobotCar is a large-scale autonomous driving dataset including driving videos captured on a consistent route during various weather conditions, traffic conditions, and times of day and night.

### B. Implementation Details

During test, we only evaluate predictions where the groundtruth depth is within 60 meters (m) in RobotCar and nuScenes, and 40 m in CARLA-EPE. For hyper-parameters in training process, please check the code release.

### C. Quantitative and Qualitative Results

Here, we compare our approach with strong baselines on nighttime depth estimation, as shown in Table I and demonstrate more qualitative results in Fig. 5. Table I shows the comparison results with the daytime self-supervised depth

TABLE II

ABLATION STUDY OF EACH PROPOSED COMPONENT.

Method	Max Depth	Error		Accuracy	
		Abs Rel	RMSE	a1	a2
<b>nuScenes</b>					
Separate Training	60 m	0.325	10.196	0.543	0.782
Joint Training (J.T.)	60 m	0.317	9.779	0.546	0.786
J.T. + $M_h$	60 m	0.314	9.503	0.551	0.791
J.T. + $M_{uc}$	60 m	0.302	9.201	0.556	0.793
J.T. + $M_{uc}$ + Denoise (Full)	60 m	<b>0.292</b>	<b>9.126</b>	<b>0.572</b>	<b>0.805</b>

estimation method MonoDepth2 [23] and other competitive nighttime methods such as DeFeat-Net [55], ADFA [56], ADDS [36], and RNW [1] on the nuScenes and RobotCar dataset. The metrics in Table I are standard ones used in prior works like [32]. Overall, our method achieves a significant improvement over the other baselines in all metrics on both datasets and sets SOTA. In the following, we choose Abs Rel and a1 as the representative metrics to indicate error and accuracy, respectively. On the nuScenes and RobotCar datasets, our method improves the accuracy of RNW by 16.2% and 3.5% and reduces the error by 10.4% and 7.6%. The improvement on the nuScenes dataset is more significant because the overexposed and underexposed regions are prevalent in its scenes, and the images are noisier. This is well aligned with the theoretical expectations of our method, as introduced in the method section.

We also provide the qualitative visualization of the predicted depth maps in Fig. 5. The blue boxes show that the RNW suffers severely from overexposure, predicting obviously wrong depth. The red boxes also show that the RNW cannot estimate the object depth in underexposed regions. Owing to the new framework and soft mask strategy we proposed, our method could predict more reasonable depths in unexpected regions and capture the correct boundary of the small or moving objects in nighttime images (such as lampposts and cars). In addition, our model can reach 25 FPS on a single 2080Ti during inference.

### D. Ablation Study

Here, we provide ablations on nuScenes dataset for the designs of each component we proposed.

**Joint Training Framework:** The result of our joint

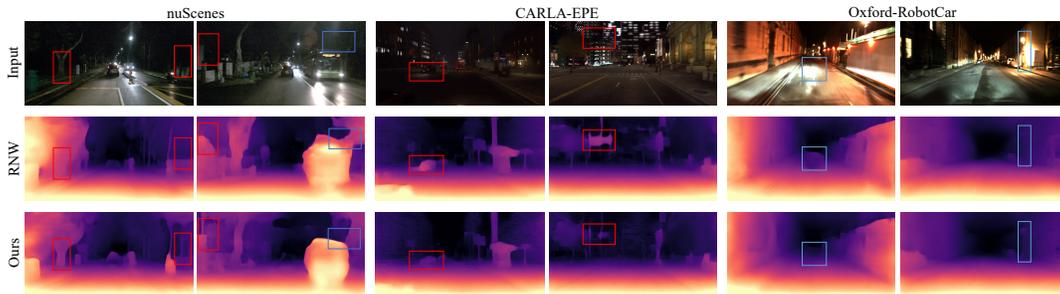


Fig. 5. Qualitative results on nuScenes, CARLA-EPE, and RobotCar datasets.

TABLE III

ABLATION STUDY OF THE EVALUATION ON SPARSE AND DENSE DEPTH GROUND TRUTH.

Method	Max Depth	Error		Accuracy	
		Abs Rel	RMSE	a1	a2
RNW w. dense	40 m	1.164	9.184	0.173	0.330
Ours w. dense	40 m	1.121	8.992	0.174	0.336
RNW w. sparse	40 m	0.975	8.210	0.283	0.582
Ours w. sparse	40 m	<b>0.941</b>	<b>7.987</b>	<b>0.310</b>	<b>0.592</b>

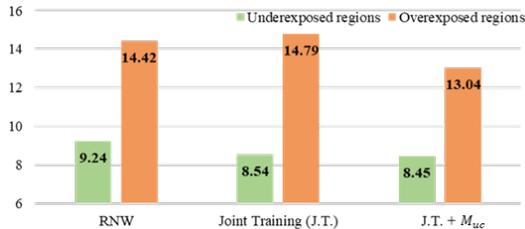


Fig. 6. We manually select 63 and 65 images in the test set of nuScenes, containing underexposed and overexposed regions respectively, and then count RMSE of these regions under three settings (RNW, Joint Training (J.T.), and J.T. + Mask  $M_{uc}$ ).

training framework without masking and denoising is shown in the first row of Table II. We find that it still significantly outperforms the RNW [1] in Table I and credit it to the new formulation. Since our method makes it possible to learn these two tasks collaboratively towards the final goal, the performance is largely promoted. The second row demonstrates our method performs better than a baseline that combines a pre-trained SIE with the self-supervised depth estimator, showing the importance of joint self-supervised training. As expected, Fig. 6 shows that the image enhancement could improve the performance in underexposed regions but not in overexposed regions.

**Soft Masking:** The fourth row of Table II shows the advantage of the proposed uncertainty masking strategy combined with the joint training framework. Given that the self-supervised enhancement loss  $L_f$  (Eq. 10) is built upon the insight that input nighttime images are largely proportional to the illumination component. Thus, a natural question is whether it is effective to use the original input image for illumination distribution modelling to suppress unexpected regions as  $M_{uc}$  does. So, we build a new baseline called  $M_h$  and the result of it shows that (1) it also works, verifying the idea of softly masking out unexpected regions is important and the statistical model well serves its purpose. (2)  $M_{uc}$

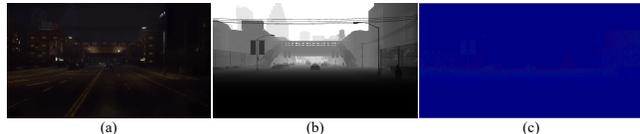


Fig. 7. Comparison between sparse and dense ground truth. (a) is the RGB image, (b) is the dense ground truth of the simulator, and (c) is the sparse ground truth from LiDAR. Please zoom in electronically to inspect the sparse depth map.

works better, showing that using the mid-product illumination map is a better choice and the necessity of combining these two tasks because the illumination map is provided by the self-supervised enhancement task. Fig. 6 shows the effectiveness of the masking strategy in overexposed regions.

**Denoising:** The last row of Table II indicates the positive impacts of introducing the denoising module, which we credit improved photometric consistency after denoising.

**CARLA-EPE Dataset:** In the first two rows of Table III, our method still outperforms RNW in CARLA-EPE dataset with the dense ground truth, and the dense protocol is more challenging since the results are worse than those on nuScenes dataset with sparse ground truth. Then we sample the dense depth to generate a sparse depth map according to the distribution of the LiDAR data, as shown in Fig. 7. As expected, the last two rows of Table III show that the performance is improved in this setting, which indicates the limitation of the evaluation with sparse depth maps.

## E. Conclusion

In this paper, we propose a self-supervised framework to jointly learn image enhancement and depth estimation. Delving into this framework, we find an effective mid-product of the image enhancement to generate a pixel-wise mask to suppress the overexposed and underexposed regions. Benefiting from these improvements, we set SOTA on public nighttime depth estimation benchmarks. Additionally, we contribute a novel photo-realistically enhanced nighttime dataset with dense depth ground truth.

## F. Acknowledgments

This work was supported by the National Natural Science Foundation of China (NSFC) under Grants No. 62173324 and Tsinghua-Toyota Joint Research Fund (20223930097).

## REFERENCES

- [1] K. Wang, Z. Zhang, Z. Yan, X. Li, B. Xu, J. Li, and J. Yang, "Regularizing nighttime weirdness: Efficient self-supervised monocular depth estimation in the dark," in *ICCV*, 2021.
- [2] F. Bu, T. Le, X. Du, R. Vasudevan, and M. Johnson-Roberson, "Pedestrian planar lidar pose (pplp) network for oriented pedestrian detection based on planar lidar and monocular images," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1626–1633, 2019.
- [3] K. Ryu, K.-i. Lee, J. Cho, and K.-J. Yoon, "Scanline resolution-invariant depth completion using a single image and sparse lidar point cloud," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 6961–6968, 2021.
- [4] N. Vödisch, O. Unal, K. Li, L. Van Gool, and D. Dai, "End-to-end optimization of lidar beam configuration for 3d object detection and localization," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 2242–2249, 2022.
- [5] C. Bai, T. Xiao, Y. Chen, H. Wang, F. Zhang, and X. Gao, "Fasterlio: Lightweight tightly coupled lidar-inertial odometry using parallel sparse incremental voxels," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4861–4868, 2022.
- [6] R. McCraith, E. Insafutdinov, L. Neumann, and A. Vedaldi, "Lifting 2d object locations to 3d by discounting lidar outliers across objects and views," in *ICRA*, 2022.
- [7] S. Huang, Y. Xie, S.-C. Zhu, and Y. Zhu, "Spatio-temporal self-supervised representation learning for 3d point clouds," in *ICCV*, 2021.
- [8] Y. Chen, H. Li, R. Gao, and D. Zhao, "Boost 3-d object detection via point clouds segmentation and fused 3-d giou-l1 loss," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 2, pp. 762–773, 2020.
- [9] X. Chen, H. Zhao, G. Zhou, and Y.-Q. Zhang, "Pq-transformer: Jointly parsing 3d objects and layouts from point clouds," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 2519–2526, 2022.
- [10] Y. Han, Y. Chen, H. Li, M. Ma, and D. Zhao, "Fast depth estimation of object via neural network perspective projection," in *2022 IEEE 11th Data Driven Control and Learning Systems Conference (DDCLS)*, 2022.
- [11] B. Tian, L. Luo, H. Zhao, and G. Zhou, "Vibus: Data-efficient 3d scene parsing with viewpoint bottleneck and uncertainty-spectrum modeling," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 194, pp. 302–318, 2022.
- [12] J. Zhang, H. Zhao, A. Yao, Y. Chen, L. Zhang, and H. Liao, "Efficient semantic scene completion network with spatial group convolution," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 733–749.
- [13] B. Jin, B. Tian, H. Zhao, and G. Zhou, "Language-guided semantic style transfer of 3d indoor scenes," in *Proceedings of the 1st Workshop on Photorealistic Image and Environment Synthesis for Multimedia Experiments*, 2022, pp. 11–17.
- [14] L. Jing, R. Yu, H. Kretschmar, K. Li, C. R. Qi, H. Zhao, A. Ayvaci, X. Chen, D. Cower, Y. Li *et al.*, "Depth estimation matters most: Improving per-object depth estimation for monocular 3d detection and tracking," *arXiv preprint arXiv:2206.03666*, 2022.
- [15] V. Licăret, V. Robu, A. Marcu, D. Costea, E. Slușanșchi, R. Sukthankar, and M. Leordeanu, "Ufo depth: Unsupervised learning with flow-based odometry optimization for metric depth estimation," in *ICRA*, 2022.
- [16] J. Choi, D. Jung, Y. Lee, D. Kim, D. Manocha, and D. Lee, "Selftune: Metrically scaled monocular depth estimation through self-supervised learning," *arXiv preprint arXiv:2203.05332*, 2022.
- [17] V. Patil, A. Liniger, D. Dai, and L. Van Gool, "Improving depth estimation using map-based depth priors," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3640–3647, 2022.
- [18] H. E. Boulahbal, A. Voicila, and A. I. Comport, "Instance-aware multi-object self-supervision for monocular depth prediction," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 10962–10968, 2022.
- [19] Y. Chen, D. Zhao, L. Lv, and Q. Zhang, "Multi-task learning for dangerous object detection in autonomous driving," *Information Sciences*, vol. 432, pp. 559–571, 2018.
- [20] H. Fu, M. M. Gong, C. H. Wang, K. Batmanghelich, and D. C. Tao, "Deep ordinal regression network for monocular depth estimation," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2002–2011, 2018.
- [21] Y. Kuznetsov, J. Stuckler, and B. Leibe, "Semi-supervised deep learning for monocular depth map prediction," in *CVPR*, 2017.
- [22] S. Li, X. Wu, Y. Cao, and H. Zha, "Generalizing to the open world: Deep visual odometry with online adaptation," in *CVPR*, 2021.
- [23] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, "Digging into self-supervised monocular depth estimation," in *ICCV*, 2019.
- [24] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3d packing for self-supervised monocular depth estimation," in *CVPR*, 2020.
- [25] Y. Almalioglu, M. R. U. Saputra, P. P. B. d. Gusmão, A. Markham, and N. Trigoni, "Ganvo: Unsupervised deep monocular visual odometry and depth estimation with generative adversarial networks," in *ICRA*, 2019.
- [26] S. Pillai, R. Ambrus, and A. Gaidon, "Superdepth: Self-supervised, super-resolved monocular depth estimation," in *ICRA*, 2019.
- [27] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [28] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *CVPR*, 2016.
- [29] A. Ranjan, V. Jampani, L. Balles, K. Kim, D. Sun, J. Wulff, and M. J. Black, "Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation," in *CVPR*, 2019.
- [30] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "Carla: An open urban driving simulator," in *CoRL*, 2017.
- [31] S. R. Richter, H. A. Al Hajja, and V. Koltun, "Enhancing photorealism enhancement," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [32] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, "Unsupervised learning of depth and ego-motion from video," in *CVPR*, 2017.
- [33] Z. Yin and J. Shi, "Geonet: Unsupervised learning of dense depth, optical flow and camera pose," in *CVPR*, 2018.
- [34] V. Casser, S. Pirk, R. Mahjourian, and A. Angelova, "Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos," in *AAAI*, 2019.
- [35] R. Mahjourian, M. Wicke, and A. Angelova, "Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5667–5675, 2018.
- [36] L. Liu, X. Song, M. Wang, Y. Liu, and L. Zhang, "Self-supervised monocular depth estimation for all day images using domain separation," in *ICCV*, 2021.
- [37] S. M. Pizer, E. P. Amburn, J. D. Austin, R. Cromartie, A. Geselowitz, T. Greer, B. ter Haar Romeny, J. B. Zimmerman, and K. Zuiderveld, "Adaptive histogram equalization and its variations," *Computer vision, graphics, and image processing*, vol. 39, no. 3, pp. 355–368, 1987.
- [38] E. H. Land, "The retinex theory of color vision," *Scientific american*, vol. 237, no. 6, pp. 108–129, 1977.
- [39] D. J. Jobson, Z.-u. Rahman, and G. A. Woodell, "A multiscale retinex for bridging the gap between color images and the human observation of scenes," *IEEE Transactions on Image processing*, vol. 6, no. 7, pp. 965–976, 1997.
- [40] ———, "Properties and performance of a center/surround retinex," *IEEE transactions on image processing*, vol. 6, no. 3, pp. 451–462, 1997.
- [41] X. Guo, Y. Li, and H. Ling, "Lime: Low-light image enhancement via illumination map estimation," *IEEE Transactions on image processing*, vol. 26, no. 2, pp. 982–993, 2016.
- [42] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep retinex decomposition for low-light enhancement," *arXiv preprint arXiv:1808.04560*, 2018.
- [43] Y. Zhang, X. Guo, J. Ma, W. Liu, and J. Zhang, "Beyond brightening low-light images," *International Journal of Computer Vision*, vol. 129, no. 4, pp. 1013–1037, 2021.
- [44] R. Liu, L. Ma, J. Zhang, X. Fan, and Z. Luo, "Retinex-inspired unrolling with cooperative prior architecture search for low-light image enhancement," in *CVPR*, 2021.
- [45] L. Ma, T. Ma, R. Liu, X. Fan, and Z. Luo, "Toward fast, flexible, and robust low-light image enhancement," in *CVPR*, 2022.
- [46] M. Jaderberg, K. Simonyan, A. Zisserman *et al.*, "Spatial transformer networks," in *Advances in neural information processing systems*, 2015, pp. 2017–2025.
- [47] W. Lee, S. Son, and K. M. Lee, "Ap-bsn: Self-supervised denoising for real-world images via asymmetric pd and blind-spot network," in *CVPR*, 2022.
- [48] A. Abdelhamed, S. Lin, and M. S. Brown, "A high-quality denoising dataset for smartphone cameras," in *CVPR*, 2018.

- [49] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *ECCV*, 2016.
- [50] S. R. Richter, Z. Hayder, and V. Koltun, "Playing for benchmarks," in *ICCV*, 2017.
- [51] C. Sakaridis, D. Dai, and L. Van Gool, "Map-guided curriculum domain adaptation and uncertainty-aware evaluation for semantic nighttime image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [52] X. Wu, Z. Wu, H. Guo, L. Ju, and S. Wang, "Dannet: A one-stage domain adaptation network for unsupervised nighttime semantic segmentation," in *CVPR*, 2021.
- [53] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nusenes: A multimodal dataset for autonomous driving," in *CVPR*, 2020.
- [54] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 year, 1000 km: The oxford robotcar dataset," *The International Journal of Robotics Research*, vol. 36, no. 1, pp. 3–15, 2017.
- [55] J. Spencer, R. Bowden, and S. Hadfield, "Defeat-net: General monocular depth via simultaneous unsupervised representation learning," in *CVPR*, 2020.
- [56] M. Vankadari, S. Garg, A. Majumder, S. Kumar, and A. Behera, "Unsupervised monocular depth estimation for night-time images using adversarial domain feature adaptation," in *ECCV*, 2020.