

Self-supervised coarse-to-fine monocular depth estimation using a lightweight attention module

Yuanzhen Li¹, Fei Luo¹ (✉), and Chunxia Xiao¹ (✉)

© The Author(s) 2022.

Abstract Self-supervised monocular depth estimation has been widely investigated and applied in previous works. However, existing methods suffer from texture-copy, depth drift, and incomplete structure. It is difficult for normal CNN networks to completely understand the relationship between the object and its surrounding environment. Moreover, it is hard to design the depth smoothness loss to balance depth smoothness and sharpness. To address these issues, we propose a coarse-to-fine method with a normalized convolutional block attention module (NCBAM). In the coarse estimation stage, we incorporate the NCBAM into depth and pose networks to overcome the texture-copy and depth drift problems. Then, we use a new network to refine the coarse depth guided by the color image and produce a structure-preserving depth result in the refinement stage. Our method can produce results competitive with state-of-the-art methods. Comprehensive experiments prove the effectiveness of our two-stage method using the NCBAM.

Keywords monocular depth estimation; texture copy; depth drift; attention module

1 Introduction

Depth information in a 2D image has a wide range of applications, including 3D reconstruction [1–5], simultaneous localization and mapping (SLAM) [6], shadow removal [7], and so on. Range finding sensors, such as LiDAR, time of flight cameras (TOF), and stereo cameras, are often used to extract depth information. However, it is unrealistic to rely on

such expensive or complex sensors in many cases. This has advanced the development of learning-based methods using large datasets [8, 9]. Supervised monocular depth estimation methods have made great progress [10]. However, collecting extensive and high-quality ground truth depth is challenging due to sensor noise and unpredictable complex environmental conditions. Supervised monocular depth estimation thus has limited generalization ability.

Recently, self-supervised monocular depth estimation approaches have been introduced, trained with stereo image pairs [11] or monocular video sequences [12–14], and supervised with geometric information. Compared to stereo-based supervision, monocular video is more attractive, as more sequenced frames are available for use as supervision signals. To enhance the performance of depth estimation, many works focus on masking strategies [12, 13, 15], loss functions [12], and multi-task learning [16, 17]. However, existing self-supervised monocular depth methods still suffer from texture-copy, depth drift, and incomplete structure.

Texture-copy in depth map is a situation in which the details of the color image are transferred to the depth map. Monodepth2 [12] upsamples the generated multi-scale depths to the input image resolution and then computes all losses, to partially alleviate the texture-copy phenomenon. Depth drift occurs when object depth largely differs from its surrounding environment in the wrong way. It is caused by the depth network incompletely understanding the spatial correlation between the object and its surrounding environment. The incomplete structure indicates that the object depth is not completely predicted, especially for sharp objects in the scene, as the smoothness loss mistakenly eliminates the depth differences of the

¹ School of Computer Science, Wuhan University, Wuhan 430072, China. E-mail: Y. Li, yuanzhen@whu.edu.cn; F. Luo, luofei@whu.edu.cn (✉); C. Xiao, cxxiao@whu.edu.cn (✉).

Manuscript received: 2022-01-07; accepted: 2022-02-22

sharp object. In Fig. 1, we illustrate some typical examples of the above problems in the predicted depth; our predicted depth maps are better than those of the comparator methods.

We propose a coarse-to-fine method with a normalized convolutional block attention module (NCBAM). Our pipeline includes coarse depth estimation and depth refinement, as shown in Fig. 2. Specially, we improve the lightweight CBAM attention module [19], to provide a normalized convolutional block attention module (NCBAM), and then incorporate it into networks to tackle the problems of texture-copy and depth drift. Furthermore, we design a network that uses the corresponding color image as a guide to refine the coarse depth, which can deal with the incomplete structure problem. The coarse depth network and depth refinement network are trained.

To summarize, this paper presents the following two main contributions:

- We tackle the texture-copy and depth drift problems by improving the CBAM and incorporating it into depth and pose networks.
- We tackle the incomplete structure problem by designing a new network using the color image as a guide to refine the coarse depth.

2 Related work

2.1 Background

Inferring depth from a single image is an ill-posed problem. However, deep learning has shown its ability to provide acceptable estimation results based on large-scale datasets. In this section, we mainly review related work on self-supervised monocular depth estimation.

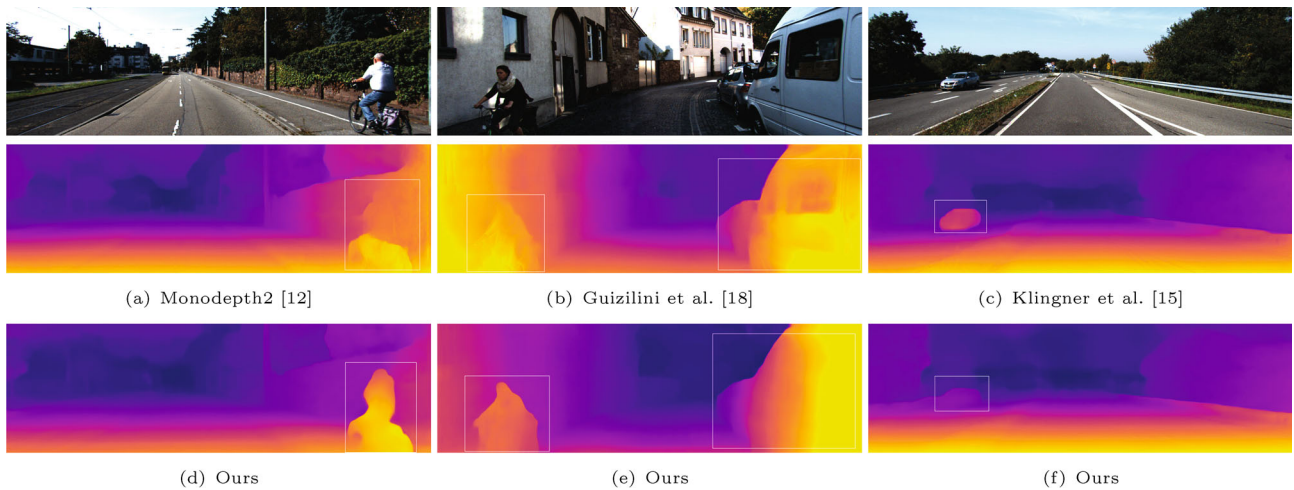


Fig. 1 Existing methods suffer from texture-copy, depth drift, and incomplete structure. (a) The depth of the person exhibits an incomplete structure problem in the Monodepth2 [12] output. (b) The depths of the person and the car suffer from a texture-copy problem in the result of Guizilini et al. [18]. (c) The car depth has a drift problem in the output of Klingner et al. [15].

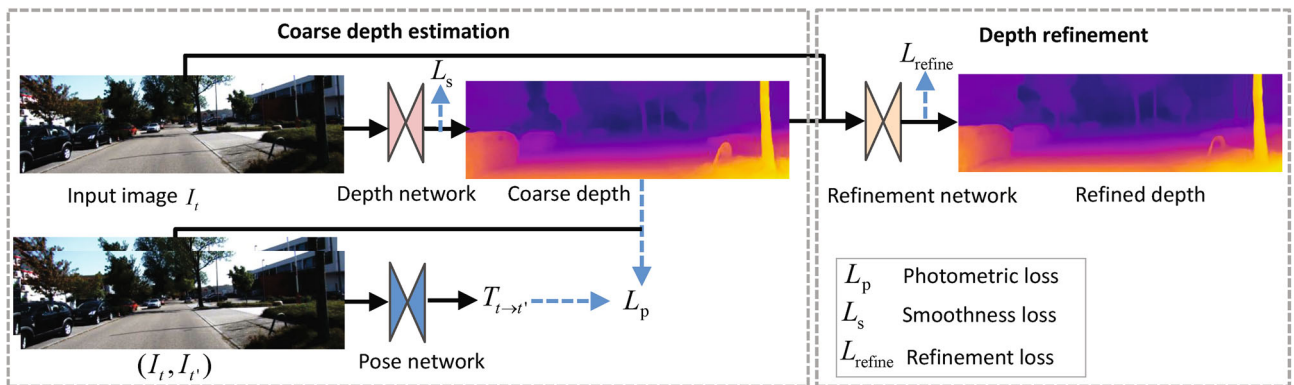


Fig. 2 Overview of our method. We use a coarse-to-fine method comprising coarse depth estimation and depth refinement. We train the depth and pose estimation models in the coarse depth estimation stage. In the depth refinement stage, we only train the depth refinement model.

Traditionally, structure-from-motion [20] and binocular stereo algorithms [21] have been used to estimate depth from a series of images or stereo image pairs, respectively. Recently, learning-based algorithms have made great progress in monocular depth estimation [22, 23]. Supervised methods train a network model on sparse depth labels provided using RGBD sensors. However, it is not easy to collect high-quality ground truth depths. As an alternative, self-supervised depth estimation has attracted attention, using stereo image pairs [11] or monocular video sequences [12, 24] as training datasets. Self-supervised depth estimation trains the depth estimation model by projecting one view to nearby views based on the predicted depth and minimizing the photometric re-projection loss between the projected image and the target image.

2.2 Self-supervised stereo training

Deep3D [25] uses a deep neural network to generate 3D stereo image pairs from 2D images or video frames and uses the photometric re-projection loss to train the depth network on the stereo image pair datasets. This network predicts a probabilistic disparity map for the input image, and the depth-based image rendering layer produces the right image in the context of binocular pairs. Garg et al. [26] proposed a deep neural network to directly estimate the depth and trained loss terms including a photometric re-projection loss and a depth smoothness loss. Monodepth [11] inputs a left image into a depth network and predicts left–right disparities to enforce mutual consistency. This method uses photometric re-projection loss and introduces a left–right disparity consistency loss. Both methods in Refs. [27] and [28] use generative adversarial networks to train the depth network.

2.3 Self-supervised monocular training

Monocular video is more attractive than stereo-based supervision, as more frame sequences are available for use as supervision signals. Self-supervised monocular training needs to estimate the parameters of the depth and pose estimation models. The pose estimation network takes a finite series of frames as input and outputs the relative camera pose. The source frame is warped to the target frame based on the predicted depth and relative camera pose, and then the photometric error between the warped frame and

the target frame is used to supervise the model during training [13].

The method in Ref. [13] was the first work that used monocular video to train end-to-end depth and camera pose estimation networks. Mahjourian et al. [29] used a 3D geometry consistency loss to train the model. Godard et al. [12] made the following three innovations. First, they proposed a minimum photometric re-projection loss to address the problem of occluded pixels. Then, they designed an auto-masking loss to ignore training pixels that violate relative camera motion assumptions. Finally, they upsampled the predicted depth maps to the input resolution and computed all losses to reduce texture-copy artifacts.

Multi-task training strategies are also available to improve the accuracy of depth estimation. Yang et al. [17] constrained the depth to be consistent with the surface normal and image edges. Ying and Shi [24] learned depth, optical flow, and pose together and used the predicted depth and optical flow to mask moving objects during training. Zhu et al. [30] used edge consistency between the semantic segmentation and depth map as a supervision signal. Klingner et al. [15] used the learned semantic information to eliminate the influence of moving objects when computing photometric re-projection loss.

Self-attention (Transformer) [31] has improved the performance of natural language processing systems by better handling of long-range dependencies between words. In addition, self-attention has been applied in computer vision tasks such as semantic segmentation [32] and depth estimation [33–35]. Johnston and Carneiro [36] used the ResNet-101 network to encode the input image and then passed it through a self-attention module [31] to explore contextual information, allowing the inference of similar depth values in discontinuous regions of the input image. However, as the self-attention module requires much memory, they only incorporated it into the encoder output layer.

Attention mechanisms have achieved great success in many visual tasks, such as image classification, object detection, and semantic segmentation [37]. We improve the lightweight attention module CBAM [19], to give a normalized convolutional block attention module (NCBAM). Then, we generalize the NCBAM model in multiple places, including the depth

estimation network, relative pose estimation network, and depth refinement network, to improve the accuracy of the depth and pose estimation models.

3 Method

In this section, we give a detailed description of our method (see Fig. 3). First, we introduce the improvements to CBAM to give NCBAM. Then, we describe the coarse depth and pose estimation methods. Finally, we present the depth refinement approach. We use the Monodepth2 [12] network as a baseline.

3.1 NCBAM attention module

The convolutional block attention module (CBAM) [19] is a lightweight module, which can aggregate deep features. It sequentially infers attention maps along with two separate sub-modules: channel and spatial, as shown in Fig. 4. The attention maps are multiplied by the input feature map for adaptive feature refinement. The CBAM attention module can learn correlations between the object itself and

the surrounding environment. When incorporating the CBAM model into depth estimation networks, it cannot completely solve the problems of texture-copy and depth drift (as shown in Fig. 12 later).

We improve upon the CBAM module in the following ways, in a normalized convolutional block attention module (NCBAM). To reduce the differences between global average pooling and global max-pooling in the channel and spatial attention modules, we convert the input feature to the range $(-1, 1)$ using the tanh function. We use the activation function $\text{softplus}(x) = \log(1 + e^x)$ to replace $\text{relu}(x) = \max(0, x)$ in the shared network of the channel submodule and the convolution layer of the spatial submodule. The activation function softplus can be seen as smoothing relu, avoiding excessive neuronal death during training. Experiment comparisons show that the NCBAM module can produce better results than the CBAM module.

3.2 Coarse depth estimation

We need to train depth and pose estimation networks simultaneously based on monocular video training.

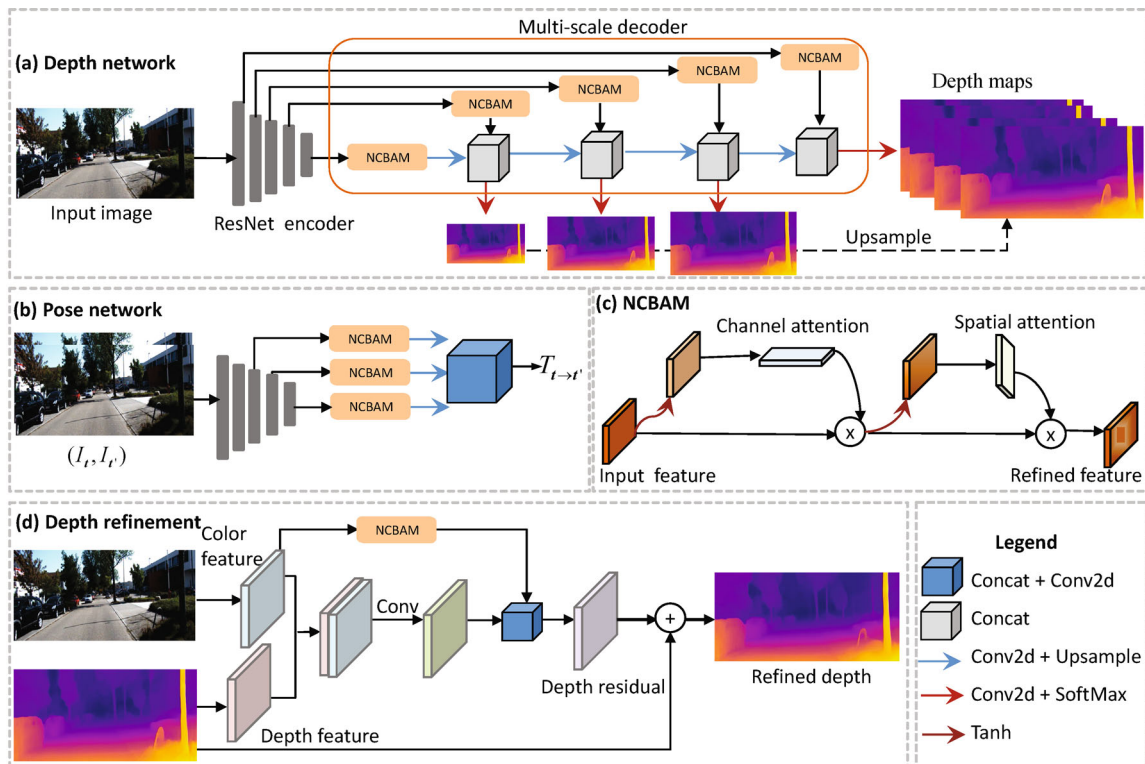


Fig. 3 Detailed architecture of our method. (a) Depth network. We incorporate the NCBAM module into the depth estimation network [12], which is a U-Net network, an encoder with residual blocks and a decoder with skip connections. (b) Pose network. We incorporate the NCBAM module into the standard pose network [12]. (c) NCBAM module. (d) Depth refinement network. This uses the color image as guidance to refine the coarse depth.

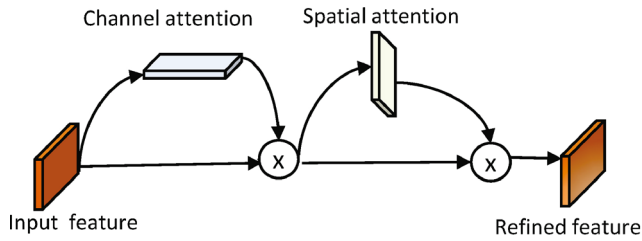


Fig. 4 The CBAM attention module [19] has two sequential sub-modules: channel and spatial.

The depth estimation network $f_D : I \rightarrow D$ predicts the depth for every pixel in the target image I_t . The pose estimation network $f_T : (I_t, I_{t'}) \rightarrow T_{t \rightarrow t'}$, predicts the camera transformation relating the target image I_t to the source image $I_{t'}$. Based on the learned depth and pose, we warp the frame into adjacent frames $\{t - 1, t + 1\}$ using the photometric re-projection loss as the optimization objective function. We follow the multi-scale depth output strategy proposed in Ref. [11]. First, we input encoder features into deep convolutions and output a low-resolution depth map. Then, we apply three additional stages of upsampling-convolution that receive skip connections from the ResNet encoder to generate corresponding resolution depth maps.

The NCBAM attention module can aggregate deep features and extract correlations between the object and the surrounding environment. Therefore, we incorporate the NCBAM model into the ResNet multi-scale features and the skip connections associated with decoder layers for the U-Net architecture (*SK* feature), as shown in Fig. 3(a). Skip connections between encoder and associated decoder layers can keep high-level information in the final depth output. When incorporating the NCBAM module, the depth estimation model can overcome with the texture-copy and depth drift problems.

The output of our depth network is a pixel-wise disparity probability for multiple disparity layers, to give a discrete disparity volume (DDV) [38]. We input the *SK* feature into a 2D convolutional layer with filters of size 3×3 , and output a K channel disparity probability volume $P = \{P_1, \dots, P_K\}$ with K disparity layers:

$$d_k = d_{\min} + \Delta_d(k - 1), \quad k = 1, \dots, K \quad (1)$$

where d_{\min} and Δ_d are the minimum disparity value and disparity interval, respectively. A depth-wise softmax operation processes P to produce an actual probability map for each disparity plane:

$$P^d = \text{softmax}(P) \quad (2)$$

We extract the final disparity as a weighted sum of the disparity probabilities P^d :

$$D = \sum_{k=1}^K d_k P_k^d \quad (3)$$

We use a widely used backbone [12], which takes two images at different time steps as input and learns the relative camera pose transformation $T_{t \rightarrow t'} \in SE(3)$ between the images recorded at time steps t and t' :

$$T_{t \rightarrow t'} = f_T(I_t, I_{t'}) \quad (4)$$

The special Euclidean group $SE(3)$ defines the set of all possible rotations and translations. Such transformations are usually represented by 4×4 matrices. Following Ref. [12], we predict the six degrees of freedom pose.

We also incorporate the NCBAM model into the pose network, as shown in Fig. 3(d). The NCBAM model can also learn other related features in the two input images, enhancing the accuracy of the pose model, furthermore increasing the accuracy of the depth estimation model. First, we input a pair of color images to the ResNet-18 network to extract corresponding deep features. Then, we input those deep features into the NCBAM module to learn their correlation (CF features). Finally, we concatenate the CF features and input them to a series of 2D convolution layers to output a single six degrees of freedom relative pose.

3.2.1 Training the coarse depth network

Following Ref. [12], training our coarse depth estimation model is mainly based on minimizing per-pixel photometric re-projection loss between the source image $I_{t'}$ and target image I_t , using the learned relative pose $T_{t \rightarrow t'}$ and depth D_t . The photometric re-projection loss is defined as

$$L_p = \mu \min_{t'} \text{pe}(I_t, I_{t' \rightarrow t}) \quad (5)$$

where $\text{pe}(\cdot)$ is the photometric reconstruction error; $t' \in \{t - 1, t + 1\}$: we use the two frames temporally adjacent to I_t as the source frames [12]. μ is a binary mask that filters out stationary points:

$$\mu = [\min_{t'} \text{pe}(I_t, I_{t' \rightarrow t}) < \min_{t'} \text{pe}(I_t, I_{t'})] \quad (6)$$

where $[\cdot]$ is the Iverson bracket. The binary mask μ includes pixels where the re-projection error of $I_{t' \rightarrow t}$ is lower than the un-warped image $I_{t'}$, indicating that the object is stationary relative to the camera. Re-projection loss minimization significantly reduces

artifacts along the object boundaries in the image, leading to better accuracy. The re-projected image is defined as

$$I_{t \rightarrow t'} = I_{t'} \langle \text{proj}(\sigma(D_t), T_{t \rightarrow t'}, K) \rangle \quad (7)$$

where $\langle \cdot \rangle$ is the sampling operator; $K \in \mathbb{R}^{3 \times 3}$ is the camera intrinsic parameter matrix, identical for all images. We also apply differentiable bilinear sampling [39] to sample the source images. $\text{proj}(\cdot)$ returns 2D coordinates of the projected depths D_t in $I_{t'}$ [40]:

$$\text{proj}(\sigma(D_t), T_{t \rightarrow t'}, K) = K T_{t \rightarrow t'} D_t(p_t) K^{-1} p_t \quad (8)$$

where p_t denotes a pixel. The photometric reconstruction error function pe [11] is

$$\text{pe}(I_t, I_{t'}) = \frac{\alpha}{2}(1 - \text{ssim}(I_t, I_{t'})) + (1 - \alpha) \|I_t - I_{t'}\|_1 \quad (9)$$

where $\alpha = 0.85$, and $\text{ssim}(\cdot)$ is the structural similarity measured as in Ref. [41] with a 3×3 block filter.

Monodepth2 [12] encourages neighbouring pixels to have similar depths, and uses an edge-aware depth smoothness loss L_s weighted by image gradients to improve predictions around object boundaries:

$$L_s = |\partial_x \hat{D}_t| \exp(-|\partial_x I_t|) + |\partial_y \hat{D}_t| \exp(-|\partial_y I_t|) \quad (10)$$

where ∂_x, ∂_y are gradient operators in x, y , respectively, and $\hat{D}_t = D_t / \overline{D}_t$ is the mean-normalized inverse depth from Ref. [42] to discourage shrinking of the estimated depth. The final loss is computed as the weighted sum of masked image photometric re-projection loss L_p and smoothness loss L_s :

$$L_{\text{coarse}} = L_p + \lambda L_s \quad (11)$$

where λ weights the smoothness term. In our experiments, we set it to 0.01.

3.3 Depth refinement

3.3.1 Background

In Section 3.2, we incorporate the NCBAM module into the depth estimation network to tackle the problems of texture-copy and depth drift, which can improve the accuracy of the initial depth estimates. However, these estimates are still imperfect. In particular, the depths may exhibit incomplete structure in sharp object regions, as shown in Fig. 7. We design a refinement network that uses the color image as guidance to refine the coarse depth, and deal with the incomplete structure problem.

First, we input the color and corresponding coarse depth images into a series of 2D convolution layers

to extract their features. Then, we concatenate the output features and pass them through a 2D convolution to generate the color–depth feature. Meanwhile, we input the color feature into the NCBAM module, allowing the network to learn more about the color feature. Finally, we concatenate the color–depth feature and the refined color feature and pass them through a series of 2D convolution layers to output a depth residual. The depth residual is added to the coarse depth to get the refined depth. Table 1 presents a detailed specification of the depth refinement network.

3.3.2 Training the depth refinement model

Depth and normal are two highly correlated entities. Inspired by Ref. [43], we design a normal consistency loss for the coarse depth \hat{D} and refined depth D :

$$L_n(u, \hat{u}) = \frac{1}{N} \sum_i \left(1 - \frac{\langle \hat{u}_i, u_i \rangle}{\|\hat{u}_i\| \|u_i\|} \right) \quad (12)$$

where N denotes the number of pixels, i indexes pixels, and $u_i = (\partial_x D_i, \partial_y D_i)$. Angle minimization is performed by maximizing the dot-product.

We also use the photometric re-projection loss in Eq. (5) with the camera pose model trained in the coarse depth estimation work. Here, we use multi-scale structural similarity, MS-SSIM [44]. The photometric reconstruction error function pe is

$$\text{pe}'(I_t, I_{t'}) = (1 - \alpha) \|I_t - I_{t'}\|_1 + \frac{\alpha}{2}(1 - \text{msssim}(I_t, I_{t'}, s)) \quad (13)$$

where $s = 4$ is the number of scales employed. Here, the photometric re-projection loss is

$$L_{p'} = \mu \min_{t'} \text{pe}'(I_t, I_{t \rightarrow t'}) \quad (14)$$

Table 1 Refinement network architecture. \mathbf{k} = kernel size, \mathbf{s} = stride, \mathbf{d} = kernel dilation, \mathbf{chns} = number of input and output channels for each layer, \mathbf{input} = input source of each layer, and + indicates concatenation

| Layer | \mathbf{k} | \mathbf{s} | \mathbf{d} | \mathbf{chns} | \mathbf{active} | \mathbf{input} |
|---------|--------------|--------------|--------------|-----------------|-------------------|-------------------|
| conv1.1 | 3 | 1 | 1 | 3/32 | PReLU | I_t |
| conv1.2 | 3 | 1 | 1 | 32/32 | | conv1.1 |
| conv1.3 | 3 | 1 | 2 | 32/32 | | conv1.2 |
| conv1.4 | | | | | | NCBAM (conv1.3) |
| conv2.1 | 3 | 1 | 1 | 1/32 | PReLU | D_t |
| conv2.2 | 3 | 1 | 1 | 32/32 | | conv2.1 |
| conv3.1 | 3 | 1 | 1 | 64/24 | PReLU | conv1.2 + conv1.2 |
| conv3.2 | 3 | 1 | 1 | 24/24 | | conv3.1 |
| conv4.1 | 3 | 1 | 2 | 56/24 | PReLU | conv3.2 + conv1.4 |
| conv4.2 | 3 | 1 | 2 | 24/24 | | conv4.1 |
| conv5.1 | 3 | 1 | 2 | 24/24 | PReLU | conv4.2 |
| conv5.2 | 3 | 1 | 1 | 24/1 | | conv5.1 |

The final learning objective function of our depth refinement network is

$$L_{\text{refine}} = L_{p'} + \gamma_n L_n + \gamma_s L_s \quad (15)$$

where γ_n and γ_s are hyperparameters to control the significance of the normal term L_n and patch smoothness term L_s , respectively. In our experiments, γ_n is set to 10^{-4} , and γ_s is set to 10^{-5} .

4 Results and discussion

This section presents experimental results to verify the effectiveness of our approach. Firstly, we describe the experimental datasets and implementation details. Secondly, we present evaluations of our method on various testing configurations. Thirdly, we perform an ablation study to demonstrate that the NCBAM module can improve the accuracy of the predicted depths. Finally, we apply our predicted depths to novel view synthesis and describe limitations of our depth estimation model.

4.1 Datasets and evaluation metrics

4.1.1 Datasets

We trained our overall network model using the standard KITTI benchmark [58]. The KITTI dataset collects rectified stereo pairs of 61 scenes (containing about 42,382 stereo frames) mainly concerned with driving scenarios. The image size is 1242×375 pixels. We follow the Eigen split test dataset [22]. It contains 39,810 monocular training sequences consisting of three frames, 4424 validation sequences, and 697 for evaluation. Following previous work [13], we remove static frames before training and only evaluate depths up to a fixed range of 80 m per standard practice [12]. We use the same intrinsic parameters for all images; we set the camera principal point to the image center and the focal length to the average of all focal lengths in KITTI.

4.1.2 Training dataset augmentation

For the training dataset, we resized all images to a standard resolution (640×192), or high resolution (1024×320). We augmented the training dataset with horizontal flips, and 50% were processed by random adjustments to contrast ± 0.2 , saturation ± 0.2 , hue ± 0.1 , and brightness ± 0.2 . The additional color images were only used as depth and pose network input, but the original color images were used to compute the training loss.

4.1.3 Depth evaluation metrics

To evaluate the depth estimation model, we used four error metrics and three accuracy metrics as in Ref. [23]. The four error metrics measure the difference between predicted depth and ground-truth depth, namely the absolute relative error (Abs Rel), the squared relative error (Sq Rel), the root mean square error (RMSE), and the logarithmic root mean square error (RMSE log). The three accuracy metrics give the fraction δ of predicted depths in an image whose ratio and inverse ratio to the ground truth are within the thresholds 1.25 , 1.25^2 , and 1.25^3 .

4.2 Implementation details

In our experiments, we set the number of disparity layers $K = 98$, the minimum disparity value $d_{\min} = 10^{-5}$, and disparity interval $\Delta_d = 0.01$. We used the PyTorch framework [59] to implement our work and trained on a single Nvidia 2080Ti. We used ResNet-18, ResNet-50, and ResNet-101 as the encoders for the depth estimation network. For coarse depth estimation, we used the Adam optimizer [60] with $\alpha = 10^{-4}$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$, training for 25 epochs with a batch size of 8. For depth refinement, $\alpha = 10^{-5}$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$, training for 13 epochs with a batch size of 6. Because the ResNet-101 network needs more memory, the batch size was set to 4 in coarse depth estimation, and we did not train it in the depth refinement work.

4.3 Depth evaluation on KITTI dataset

We evaluated our depth estimation model on the Eigen split test dataset [22]. Table 2 presents the results, which demonstrate that our method is better than existing methods in terms of the seven evaluation metrics. A qualitative evaluation of our coarse depth estimation model is provided in Fig. 5, showing a comparison to results generated by the methods in Refs. [12, 18, 36, 56]. Unlike those methods, our estimated depth maps have complete structures for objects, such as the human body. The estimated depth maps from the comparator methods exhibit texture-copy phenomena in areas such as the car, but our approach overcomes these problems. The depth evaluation results verify that the NCBAM module effectively estimates monocular depth.

We also have qualitatively evaluated our coarse depth estimation model on the Cityscapes test dataset [57]: see Fig. 6. The Cityscapes and KITTI

Table 2 Comparisons to state-of-the-art methods on the KITTI test dataset [22]. In the Train column: S = self-supervised stereo pair supervision, M = self-supervised monocular video supervision. The best results in each category are in bold

| Method | Train | Memory (MB) | Error (lower is better) | | | | Accuracy (higher is better) | | |
|----------------------------|-------|-------------|-------------------------|--------------|--------------|--------------|-----------------------------|-------------------|-------------------|
| | | | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
| Zhou et al. [13] | M | 126 | 0.183 | 1.595 | 6.709 | 0.270 | 0.734 | 0.902 | 0.959 |
| Yang et al. [45] | M | — | 0.182 | 1.481 | 6.501 | 0.267 | 0.725 | 0.906 | 0.963 |
| Mahjourian et al. [29] | M | 126 | 0.163 | 1.240 | 6.220 | 0.250 | 0.762 | 0.916 | 0.968 |
| GeoNet [24] | M | 229 | 0.149 | 1.060 | 5.567 | 0.226 | 0.796 | 0.935 | 0.975 |
| DDVO [42] | M | 142 | 0.151 | 1.257 | 5.583 | 0.228 | 0.810 | 0.936 | 0.974 |
| DF-Net [46] | M | — | 0.150 | 1.124 | 5.507 | 0.223 | 0.806 | 0.933 | 0.973 |
| LEGO [17] | M | 211.35 | 0.162 | 1.352 | 6.276 | 0.252 | — | — | — |
| Ranjan et al. [16] | M | 213.471 | 0.148 | 1.149 | 5.464 | 0.226 | 0.815 | 0.935 | 0.973 |
| EPC++ [47] | M | 146.1 | 0.141 | 1.029 | 5.350 | 0.216 | 0.816 | 0.941 | 0.976 |
| Struct2depth [48] | M | 66.82 | 0.141 | 1.026 | 5.291 | 0.215 | 0.816 | 0.945 | 0.979 |
| Monodepth2 [12] | M | 66.72 | 0.115 | 0.903 | 4.863 | 0.193 | 0.877 | 0.959 | 0.981 |
| Klingner et al. [15] | M | 110.15 | 0.113 | 0.835 | 4.693 | 0.191 | 0.879 | 0.961 | 0.981 |
| Guizilini et al. [18] | M | 106.38 | 0.120 | 1.018 | 5.136 | 0.198 | 0.865 | 0.955 | 0.980 |
| Johnston and Carneiro [36] | M | 250.74 | 0.110 | 0.872 | 4.714 | 0.189 | 0.878 | 0.958 | 0.980 |
| Ours | M | 69.13 | 0.101 | 0.811 | 4.674 | 0.179 | 0.890 | 0.962 | 0.983 |
| Ours-refine | M | 70.58 | 0.098 | 0.810 | 4.672 | 0.177 | 0.890 | 0.964 | 0.983 |
| Garg et al. [49] | S | 23 | 0.152 | 1.226 | 5.849 | 0.246 | 0.784 | 0.921 | 0.967 |
| Monodepth R50 [11] | S | 668 | 0.133 | 1.142 | 5.533 | 0.230 | 0.830 | 0.936 | 0.970 |
| StrAT [50] | S | — | 0.128 | 1.019 | 5.403 | 0.227 | 0.827 | 0.935 | 0.971 |
| Poggi et al. (VGG) [51] | S | 902 | 0.119 | 1.201 | 5.888 | 0.208 | 0.844 | 0.941 | 0.978 |
| SuperDepth [52] | S | — | 0.112 | 0.875 | 4.958 | 0.207 | 0.852 | 0.947 | 0.977 |
| Monodepth2 [12] | S | 61.7 | 0.109 | 0.873 | 4.960 | 0.209 | 0.864 | 0.948 | 0.975 |
| Watson et al. [53] | S | 132.14 | 0.111 | 0.912 | 4.977 | 0.205 | 0.862 | 0.950 | 0.977 |
| MonoResMatch [54] | S | 487 | 0.111 | 0.867 | 4.714 | 0.199 | 0.864 | 0.954 | 0.979 |
| UnDeepVO [55] | MS | — | 0.183 | 1.730 | 6.571 | 0.268 | — | — | — |
| EPC++ [47] | MS | 146.1 | 0.128 | 0.936 | 5.011 | 0.209 | 0.831 | 0.945 | 0.979 |
| Monodepth2 [12] | MS | 66.72 | 0.106 | 0.818 | 4.750 | 0.196 | 0.874 | 0.957 | 0.979 |
| WaveletMonodepth [56] | MS | — | 0.109 | 0.814 | 4.808 | 0.198 | 0.868 | 0.955 | 0.980 |

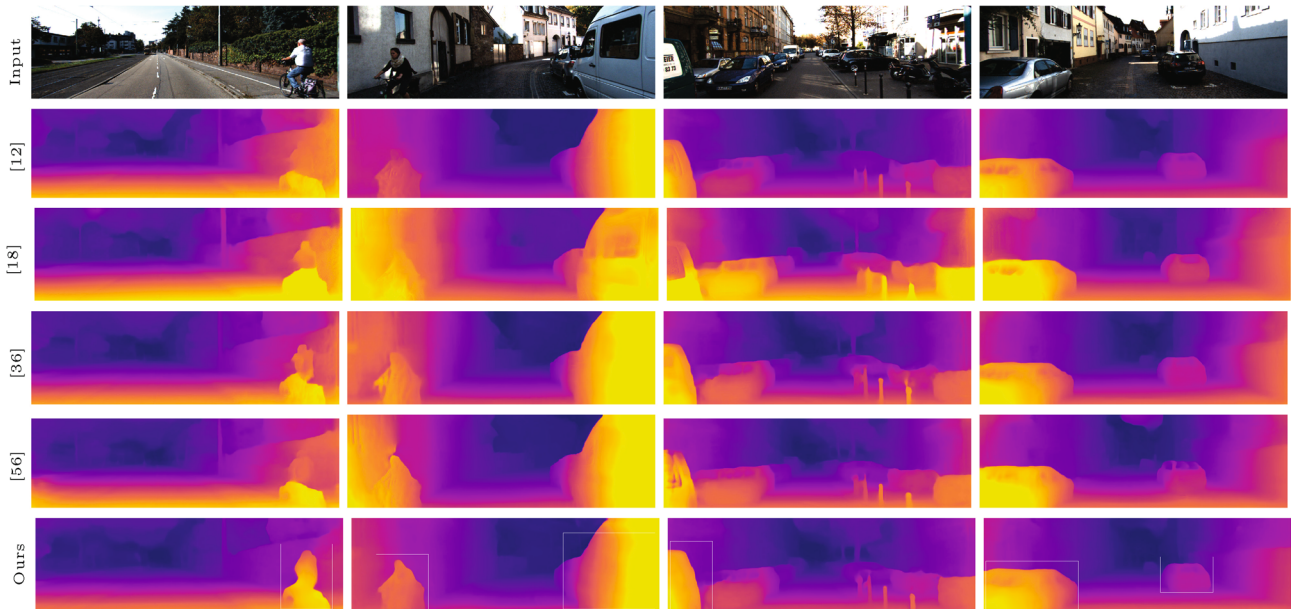


Fig. 5 Comparisons to state-of-the-art self-supervised monocular depth estimation methods: Monodepth2 [12], Guizilini et al. [18], Johnston and Carneiro [36], and WaveletMonodepth [56] using examples from the Eigen split test dataset [22].

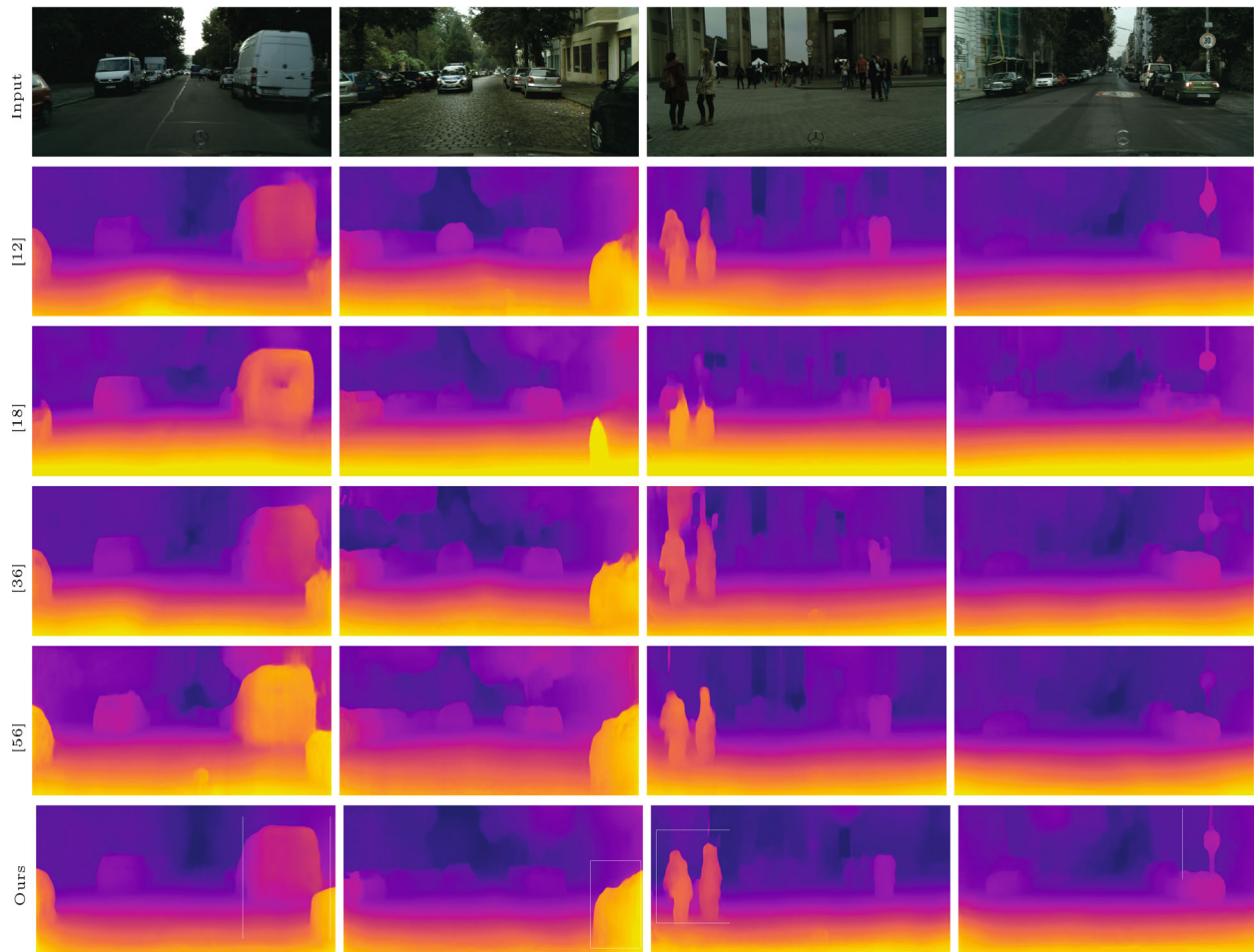


Fig. 6 Comparisons to existing self-supervised monocular depth estimation methods: Monodepth2 [12], Guizilini et al. [18], Johnston and Carneiro [36], and WaveletMonodepth [56], using the Cityscapes test dataset [57].

datasets have some differences. Our model only trained on the KITTI dataset was used to predict depths on the Cityscapes dataset. Our predicted depth maps are better than those from the methods in Refs. [12, 18, 36, 56], whose predicted depths exhibit texture-copy, depth drift, and incomplete structure problems, in areas such as cars, persons, and landmarks.

The benefits of the depth refinement model are shown in Table 2 and Fig. 7. Compared to the coarse depth estimation model, the results of the depth refinement model are further improved. In Fig. 7, we show qualitative results of the coarse and refined depth estimation models. Refined depths provide better results on thin structures such as poles. Table 2 and Fig. 7 show that our depth refinement network is effective, and can refine the coarse depth.

We have compared our method to Monodepth2 [12] and Johnston and Carneiro [36] using a variety of

encoder networks, including ResNet-18, ResNet-50, and ResNet-101. Table 3 shows the quantitative results. Our results are quantitatively better than those of these two methods. Figure 8 shows a qualitative comparison between our method and Monodepth2 [12] with the ResNet-18 encoder and input image size of 1024×320 . Unlike Monodepth2 [12], our depth estimation model can deal with the texture-copy problem and produce a clear depth for a sharp object in the image. Figure 9 shows a comparison between our method and the results of Johnston and Carneiro [36] using the ResNet-101 encoder network. Our approach can predict depths of delicate structures.

4.4 Depth evaluation on Make3D dataset

In Table 4, we provide quantitative evaluation results on the Make3D dataset [61] using our models trained on the KITTI dataset. We used the same testing protocol as Monodepth2 [12] and the evaluation

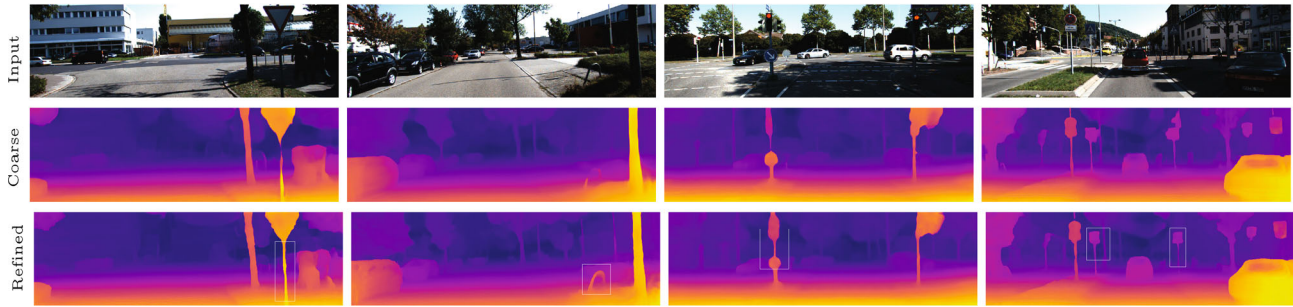


Fig. 7 Comparisons between the coarse depth and refined depth. Some depth maps of sharp objects were refined using the depth refinement network.

Table 3 Comparisons to Monodepth2 [12] and Johnston and Carneiro [36], with encoder network ResNet-18, ResNet-50, and ResNet-101. High denotes the input image resolution is 1024×320

| Encoder | Method | Error (lower is better) | | | | Accuracy (higher is better) | | |
|------------|----------------------------|-------------------------|--------------|--------------|--------------|-----------------------------|-------------------|-------------------|
| | | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
| ResNet-18 | Monodepth2 [12] | 0.115 | 0.903 | 4.863 | 0.193 | 0.877 | 0.959 | 0.981 |
| | Johnston and Carneiro [36] | 0.110 | 0.872 | 4.714 | 0.189 | 0.878 | 0.958 | 0.980 |
| | Ours | 0.101 | 0.811 | 4.674 | 0.179 | 0.890 | 0.962 | 0.983 |
| | Monodepth2 [12] (High) | 0.115 | 0.882 | 4.701 | 0.190 | 0.879 | 0.961 | 0.982 |
| | Ours (High) | 0.101 | 0.807 | 4.635 | 0.173 | 0.891 | 0.962 | 0.983 |
| ResNet-50 | Monodepth2 [12] | 0.110 | 0.831 | 4.642 | 0.187 | 0.883 | 0.962 | 0.982 |
| | Ours | 0.098 | 0.795 | 4.631 | 0.171 | 0.890 | 0.963 | 0.983 |
| ResNet-101 | Johnston and Carneiro [36] | 0.110 | 0.872 | 4.714 | 0.189 | 0.878 | 0.958 | 0.980 |
| | Ours | 0.097 | 0.788 | 4.623 | 0.170 | 0.892 | 0.963 | 0.984 |

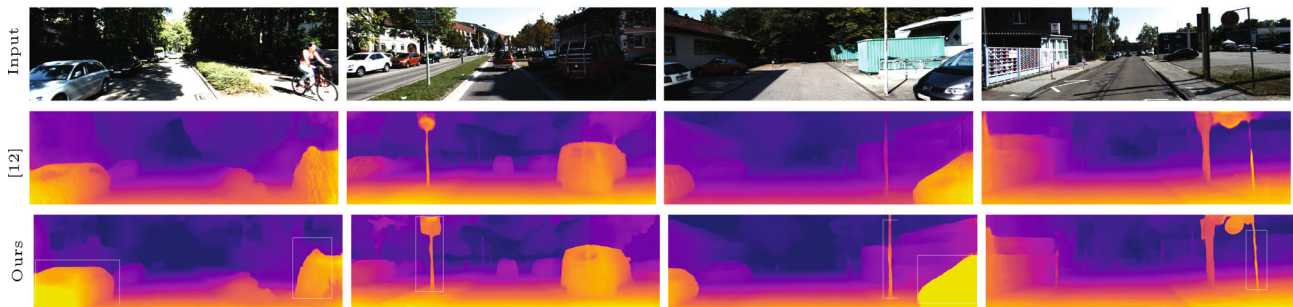


Fig. 8 Comparisons to Monodepth2 [12] on an image with 1024×320 resolution. Our method can overcome the texture-copy problem and produce clear depths for sharp objects in the image.

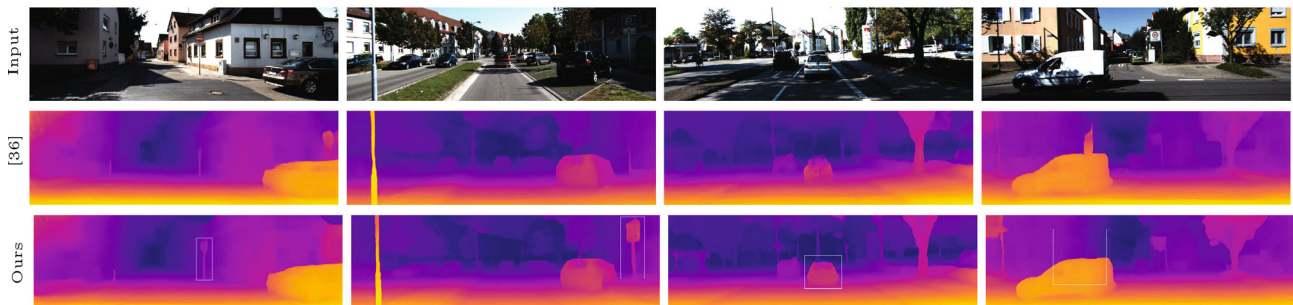


Fig. 9 Comparison to results from Johnston and Carneiro [36] using the ResNet-101 encoder network. Our method can estimate accurate depths of thin structures.

criteria from Monodepth [11]. For the Make3D dataset, we evaluated on a center crop of 2×1 ratio and used median scaling [12] because the ground truth depth and corresponding input image were not well aligned. Table 4 shows that our method can produce better results than previous self-supervision methods. Figure 10 shows our depth prediction results on the Make3D test dataset. These results demonstrate that the estimated depth is credible even though depth estimation model training did not use the Make3D dataset.

4.5 Odometry evaluation

We evaluated the pose model to validate the effectiveness of the NCBAM model on pose. Following Monodepth2 [12], the training set of our pose network model was sequences 0–8 of the KITTI odometry dataset [58], and the test dataset was sequences 9 and 10. Generally, the pose network takes five frames as input [13, 16] and predicts transformations. Our pose network baseline is Monodepth2 [12]: the input

to the pose network comprises two frame images, and the output is a relative pose transformation between that pair of frames. To evaluate the two-frame model on the five-frame test sequences, Monodepth2 [12] makes separate predictions for each of the four pairs of frame transformations for each set of five frames and combines them to form local trajectories. We follow the Monodepth2 [12] testing protocol to evaluate our pose network. Here, our pose model was trained for 11 epochs. Evaluation results are shown in Table 5: the accuracy of our pose model is better than that of Monodepth2 [12], and show that the NCBAM can enhance the results of the pose model.

4.6 Ablation study

To better validate the effectiveness of NCBAM in coarse depth estimation, we have performed an ablation study. Table 6 shows the results. We start from the baseline Monodepth2 [12] + DDV with ResNet-18 encoder network (1st row). Then, we incorporate the NCBAM module into the depth estimation network (2nd row), pose estimation

Table 4 Evaluation results on the Make3D test dataset [61]

| Method | Train | Error (lower is better) | | | |
|----------------------------|-------|-------------------------|--------------|--------------|-------------------|
| | | Abs Rel | Sq Rel | RMSE | log ₁₀ |
| Karsch et al. [62] | D | 0.425 | 4.948 | 8.290 | 0.151 |
| Liu et al. [63] | D | 0.476 | 6.611 | 10.030 | 0.167 |
| Laina et al. [64] | D | 0.205 | 1.724 | 5.578 | 0.086 |
| Monodepth [11] | S | 0.544 | 10.94 | 11.760 | 0.193 |
| Zhou et al. [13] | M | 0.383 | 5.321 | 10.470 | 0.478 |
| DDVO [42] | M | 0.387 | 4.720 | 8.090 | 0.204 |
| Monodepth2 [12] | M | 0.322 | 3.589 | 7.417 | 0.163 |
| Johnston and Carneiro [36] | M | 0.306 | 3.100 | 7.126 | 0.160 |
| Ours | M | 0.285 | 2.798 | 6.950 | 0.147 |
| Ours-refine | M | 0.283 | 2.974 | 6.949 | 0.146 |

Table 5 Odometry evaluation results on testing sequences 9 and 10 of the KITTI odometry dataset. Results are average absolute trajectory error and standard deviation in meters

| Method | Seq. 09 | Seq. 10 | Frames |
|------------------------|----------------------|----------------------|--------|
| ORB-Slam [65] | 0.014 ± 0.008 | 0.012 ± 0.001 | — |
| DDVO [42] | 0.045 ± 0.108 | 0.033 ± 0.074 | 3 |
| Zhou et al. [13] | 0.050 ± 0.039 | 0.034 ± 0.028 | 5 → 2 |
| Zhou et al. [13] | 0.021 ± 0.017 | 0.020 ± 0.015 | 5 |
| Mahjourian et al. [29] | 0.013 ± 0.010 | 0.012 ± 0.011 | 3 |
| GeoNet [24] | 0.012 ± 0.007 | 0.012 ± 0.009 | 5 |
| Ranjian et al. [16] | 0.012 ± 0.007 | 0.012 ± 0.008 | 5 |
| Monodepth2 [12] | 0.017 ± 0.008 | 0.015 ± 0.010 | 2 |
| Ours | 0.015 ± 0.001 | 0.012 ± 0.004 | 2 |

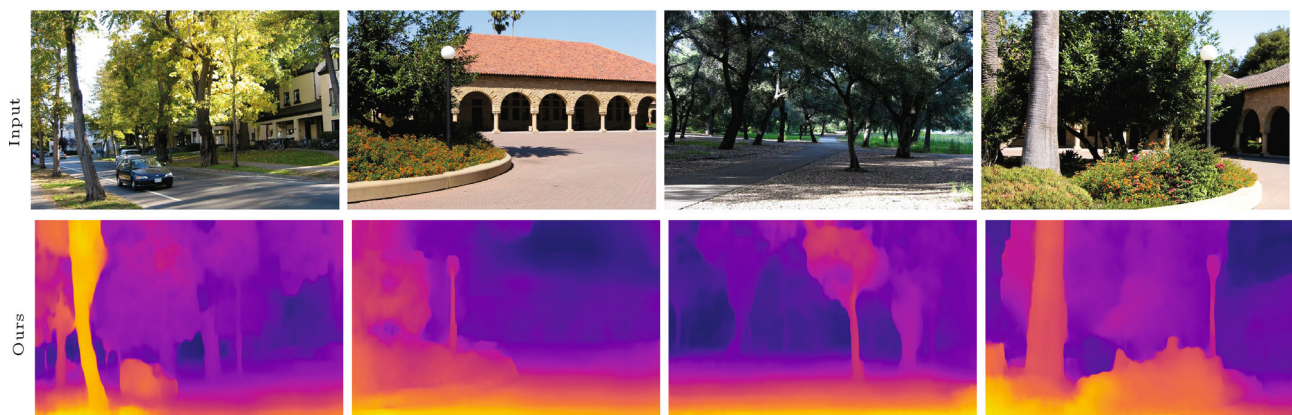


Fig. 10 Examples of our approach on the Make3D dataset [61].

Table 6 Ablation study. Evaluation of the depth estimation model with CBAM and NCBAM on the Eigen split test dataset [22]. First, we evaluate the performance of the NCBAM module used in the depth network (NCBAM-D) and pose network (NCBAM-P). The baseline is Monodepth2 [12] with ResNet-18 + DDV. Then, we compare the effectiveness of the CBAM and NCBAM attention modules. The best results are marked in bold

| Method | Error (lower is better) | | | | Accuracy (higher is better) | | |
|-------------------------|-------------------------|--------------|--------------|--------------|-----------------------------|-------------------|-------------------|
| | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
| Baseline | 0.112 | 0.893 | 4.843 | 0.190 | 0.879 | 0.961 | 0.982 |
| Baseline + NCBAM-D | 0.104 | 0.835 | 4.705 | 0.182 | 0.884 | 0.960 | 0.982 |
| Baseline + NCBAM-P | 0.106 | 0.838 | 4.709 | 0.185 | 0.883 | 0.961 | 0.981 |
| Baseline + (full) | 0.101 | 0.811 | 4.674 | 0.179 | 0.890 | 0.962 | 0.983 |
| Baseline + (full), CBAM | 0.103 | 0.812 | 4.678 | 0.181 | 0.890 | 0.960 | 0.981 |

network (3rd row), and depth and pose estimation networks (4th row). In each part of our method NCBAM improves evaluation measures. The results are significantly improved when adding the NCBAM module to both depth and pose networks. For comparison, we incorporate the CBAM module into the depth + pose network (5th row). The depth estimation model with NCBAM in row 4 is better than the one with CBAM on the seven evaluation metrics. NCBAM can better aggregate deep features and extract correlations between the

object and surrounding environment to improve depth estimation accuracy.

Figure 11 shows qualitative results of the ablation study. The best results are those in which we incorporate NCBAM into depth and pose networks. In Fig. 12, we compare depth estimation models with NCBAM and CBAM modules. Better results are provided by the version with NCBAM. Table 6, Fig. 11, and Fig. 12 show that the proposed NCBAM model is effective.

In Table 7, we illustrate the effectiveness of the

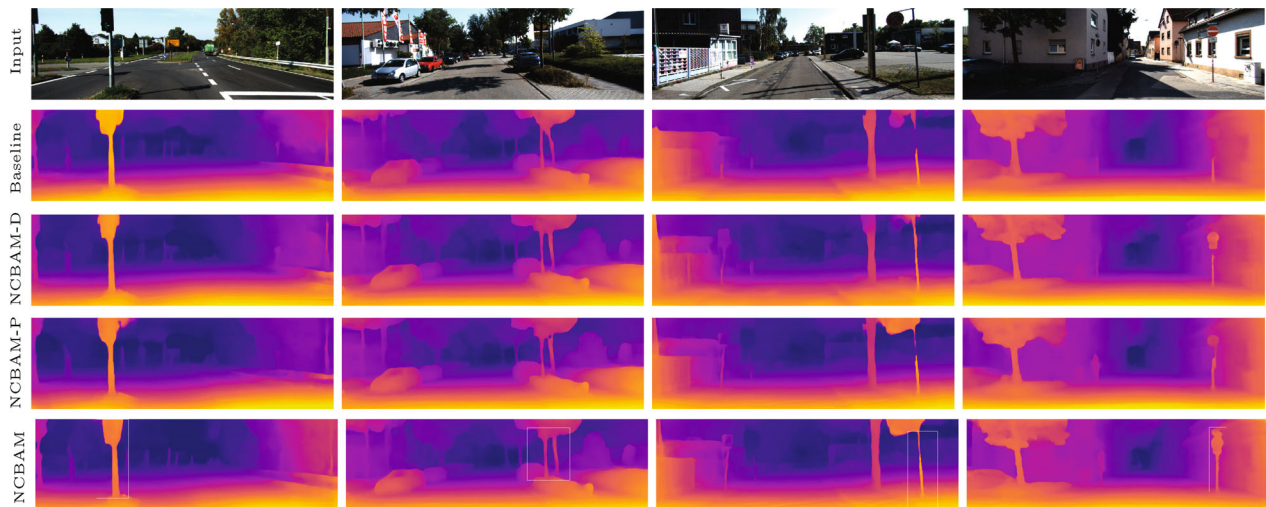


Fig. 11 Ablation study. Incorporating NCBAM into both the depth and pose networks produces the best results.

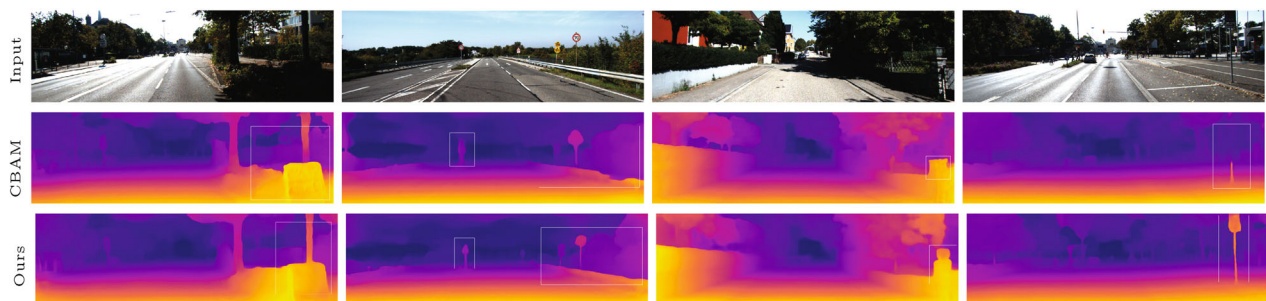


Fig. 12 Comparison of CBAM and NCBAM. Incorporating NCBAM into the depth and pose networks produces better results than CBAM.

Table 7 Improvements due to two improvements in CBAM. The best results are marked in bold

| Method | Error (lower is better) | | | | Accuracy (higher is better) | | |
|-----------------------|-------------------------|--------------|--------------|--------------|-----------------------------|-------------------|-------------------|
| | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
| CBAM | 0.103 | 0.812 | 4.678 | 0.181 | 0.890 | 0.960 | 0.981 |
| CBAM_Tanh | 0.102 | 0.811 | 4.676 | 0.180 | 0.890 | 0.960 | 0.982 |
| CBAM_softplus | 0.103 | 0.811 | 4.675 | 0.179 | 0.889 | 0.961 | 0.982 |
| NCBAM | 0.101 | 0.811 | 4.674 | 0.179 | 0.890 | 0.962 | 0.983 |
| CBAM_sigmoid | 0.104 | 0.812 | 4.677 | 0.183 | 0.890 | 0.959 | 0.981 |
| CBAM_sigmoid_softplus | 0.103 | 0.812 | 4.676 | 0.181 | 0.891 | 0.961 | 0.982 |

improvements made to CBAM. The first converts the input feature to the range $(-1, 1)$ using the tanh function in the channel and spatial attention modules (CBAM_Tanh). The second is to use the activation function softplus instead of relu in the shared network of the channel sub-module and the convolution layer of the spatial sub-module (CBAM_softplus). We can see that both changes to CBAM improve depth estimation, and the best results are those when we utilize both improvements simultaneously, i.e., NCBAM.

The normalization sigmoid activation $\text{sigmoid}(x) = 1/(1 + e^{-x})$ maps the input to the range $(0, 1)$. Table 7 compares use of tanh and sigmoid functions for normalization, and shows that tanh function produces better results.

4.7 Limitations

Although our method can overcome the above three targeted problems, our approach also has some limitations in common with other methods. One is that it cannot effectively predict the depths of moving objects. Figure 13 presents an example generated by our approach and other state-of-the-art self-supervised monocular depth estimation methods. Unfortunately, all methods fail to predict the person's depth since the training set KITTI dataset is collected in scenes nearly completely lacking in humans.

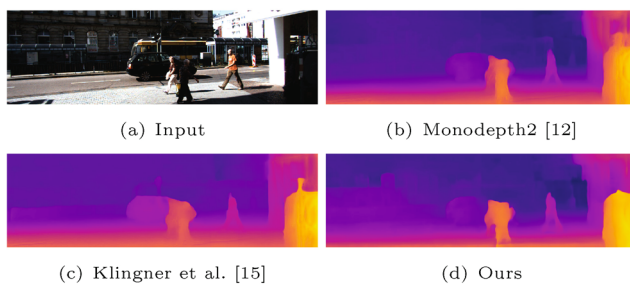


Fig. 13 Limitation. Our approach and other state-of-the-art methods fail to predict the depths of moving persons.

5 Conclusions

In this paper, we have presented a new self-supervised monocular depth estimation method. Previous methods typically produce predicted depth maps with incomplete structures, texture-copy issues, and depth drift problems. We improved the attention model CBAM, to provide NCBAM, incorporated it into networks, and proposed a coarse-to-fine approach to address the above problems. We have performed extensive experiments to compare our method to state-of-the-art methods which validate its effectiveness. In future, we will further investigate depth estimation in complex scenes containing motion objects.

Author contributions

Yuanzhen Li conceived and designed the study, and collected the data. All authors analyzed the data and were involved in writing the manuscript.

Acknowledgements

This work is partially supported by the Key Technological Innovation Projects of Hubei Province (2018AAA062), National Natural Science Foundation of China (61972298), Wuhan University–Huawei GeoInformatics Innovation Lab.

Declaration of competing interest

The authors have no competing interests to declare that are relevant to the content of this article.

References

- [1] Cao, Y. P.; Kobbelt, L.; Hu, S. M. Real-time high-accuracy three-dimensional reconstruction with consumer RGB-D cameras. *ACM Transactions on Graphics* Vol. 37, No. 5, Article No. 171, 2018.
- [2] Fu, Y. P.; Yan, Q. G.; Liao, J.; Xiao, C. X. Joint texture and geometry optimization for RGB-D reconstruction.

- In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 5949–5958, 2020.
- [3] Yang, L.; Yan, Q. G.; Fu, Y. P.; Xiao, C. X. Surface reconstruction via fusing sparse-sequence of depth images. *IEEE Transactions on Visualization and Computer Graphics* Vol. 24, No. 2, 1190–1203, 2018.
 - [4] Fu, Y. P.; Yan, Q. G.; Yang, L.; Liao, J.; Xiao, C. X. Texture mapping for 3D reconstruction with RGB-D sensor. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 4645–4653, 2018.
 - [5] Fu, Y. P.; Yan, Q. G.; Liao, J.; Zhou, H. J.; Tang, J.; Xiao, C. X. Seamless texture optimization for RGB-D reconstruction. *IEEE Transactions on Visualization and Computer Graphics* doi: 10.1109/TVCG.2021.3134105, 2021.
 - [6] Luo, H. C.; Gao, Y.; Wu, Y. H.; Liao, C. Y.; Yang, X.; Cheng, K. T. Real-time dense monocular SLAM with online adapted depth prediction network. *IEEE Transactions on Multimedia* Vol. 21, No. 2, 470–483, 2019.
 - [7] Fan, X. Y.; Wu, W. J.; Zhang, L.; Yan, Q. G.; Fu, G.; Chen, Z. P.; Long, C.; Xiao, C. Shading-aware shadow detection and removal from a single image. *The Visual Computer* Vol. 36, Nos. 10–12, 2175–2188, 2020.
 - [8] Karsch, K.; Liu, C.; Kang, S. B. Depth transfer: Depth extraction from video using non-parametric sampling. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 36, No. 11, 2144–2158, 2014.
 - [9] Watson, J.; Aodha, O. M.; Turmukhambetov, D.; Brostow, G. J.; Firman, M. Learning stereo from single images. In: *Computer Vision – ECCV 2020. Lecture Notes in Computer Science, Vol. 12346*. Vedaldi, A.; Bischof, H.; Brox, T.; Frahm, J. M. Eds. Springer Cham, 722–740, 2020.
 - [10] Guo, X.; Li, H.; Yi, S.; Ren, J.; Wang, X. Learning monocular depth by distilling cross-domain stereo networks. In: *Computer Vision – ECCV 2018. Lecture Notes in Computer Science, Vol. 11215*. Ferrari, V.; Hebert, M.; Sminchisescu, C.; Weiss, Y. Eds. Springer Cham, 506–523, 2018.
 - [11] Godard, C.; Aodha, O. M.; Brostow, G. J. Unsupervised monocular depth estimation with left-right consistency. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 6602–6611, 2017.
 - [12] Godard, C.; Aodha, O. M.; Firman, M.; Brostow, G. Digging into self-supervised monocular depth estimation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 3827–3837, 2019.
 - [13] Zhou, T. H.; Brown, M.; Snavely, N.; Lowe, D. G. Unsupervised learning of depth and ego-motion from video. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 6612–6619, 2017.
 - [14] Zhao, W.; Liu, S. H.; Shu, Y. Z.; Liu, Y. J. Towards better generalization: Joint depth-pose learning without PoseNet. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 9148–9158, 2020.
 - [15] Klingner, M.; Termöhlen, J. A.; Mikolajczyk, J.; Fingscheidt, T. Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance. In: *Computer Vision – ECCV 2020. Lecture Notes in Computer Science, Vol. 12365*. Vedaldi, A.; Bischof, H.; Brox, T.; Frahm, J. M. Eds. Springer Cham, 582–600, 2020.
 - [16] Ranjan, A.; Jampani, V.; Balles, L.; Kim, K.; Sun, D. Q.; Wulff, J.; Black, M. J. Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 12232–12241, 2019.
 - [17] Yang, Z. H.; Wang, P.; Wang, Y.; Xu, W.; Nevatia, R. LEGO: Learning edge with geometry all at once by watching videos. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 225–234, 2018.
 - [18] Guizilini, V.; Ambrus, R.; Pillai, S.; Raventos, A.; Gaidon, A. 3D packing for self-supervised monocular depth estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2482–2491, 2020.
 - [19] Woo, S.; Park, J.; Lee, J. Y.; Kweon, I. S. CBAM: Convolutional block attention module. In: *Computer Vision – ECCV 2018. Lecture Notes in Computer Science, Vol. 11211*. Ferrari, V.; Hebert, M.; Sminchisescu, C.; Weiss, Y. Eds. Springer Cham, 3–19, 2018.
 - [20] Schonberger, J. L.; Frahm, J. M. Structure-from-motion revisited. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 4104–4113, 2016.
 - [21] Hirschmuller, H. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 30, No. 2, 328–341, 2008.
 - [22] Eigen, D.; Fergus, R. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In: Proceedings of the IEEE International Conference on Computer Vision, 2650–2658, 2015.

- [23] Eigen, D.; Puhrsch, C.; Fergus, R. Depth map prediction from a single image using a multi-scale deep network. In: Proceedings of the 27th International Conference on Neural Information Processing Systems, Vol. 2, 2366–2374, 2014.
- [24] Yin, Z. C.; Shi, J. P. GeoNet: Unsupervised learning of dense depth, optical flow and camera pose. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 1983–1992, 2018.
- [25] Xie, J. Y.; Girshick, R.; Farhadi, A. Deep3D: Fully automatic 2D-to-3D video conversion with deep convolutional neural networks. In: *Computer Vision – ECCV 2016. Lecture Notes in Computer Science, Vol. 9908*. Leibe, B.; Matas, J.; Sebe, N.; Welling, M. Eds. Springer Cham, 842–857, 2016.
- [26] Garg, R.; Vijay Kumar, B. G.; Carneiro, G.; Reid, I. Unsupervised CNN for single view depth estimation: Geometry to the rescue. In: *Computer Vision – ECCV 2016. Lecture Notes in Computer Science, Vol. 9912*. Leibe, B.; Matas, J.; Sebe, N.; Welling, M. Eds. Springer Cham, 740–756, 2016.
- [27] Pilzer, A.; Xu, D.; Puscas, M.; Ricci, E.; Sebe, N. Unsupervised adversarial depth estimation using cycled generative networks. In: Proceedings of the International Conference on 3D Vision, 587–595, 2018.
- [28] Aleotti, F.; Tosi, F.; Poggi, M.; Mattoccia, S. Generative adversarial networks for unsupervised monocular depth prediction. In: *Computer Vision – ECCV 2018 Workshops. Lecture Notes in Computer Science, Vol. 11129*. Leal-Taixé, L.; Roth, S. Eds. Springer Cham, 337–354, 2019.
- [29] Mahjourian, R.; Wicke, M.; Angelova, A. Unsupervised learning of depth and ego-motion from monocular video using 3D geometric constraints. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 5667–5675, 2018.
- [30] Zhu, S. J.; Brazil, G.; Liu, X. M. The edge of depth: Explicit constraints between segmentation and depth. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 13113–13122, 2020.
- [31] Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; Polosukhin, I. Attention is all you need. In: Proceedings of the 31st International Conference on Neural Information Processing Systems, 6000–6010, 2017.
- [32] Yuan, Y. H.; Huang, L.; Guo, J. Y.; Zhang, C.; Chen, X. L.; Wang, J. D. OCNNet: Object context for semantic segmentation. *International Journal of Computer Vision* Vol. 129, No. 8, 2375–2398, 2021.
- [33] Ranftl, R.; Bochkovskiy, A.; Koltun, V. Vision transformers for dense prediction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 12159–12168, 2021.
- [34] Li, Z. S.; Liu, X. T.; Drenkow, N.; Ding, A.; Creighton, F. X.; Taylor, R. H.; Unberath, M. Revisiting stereo depth estimation from a sequence-to-sequence perspective with transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 6177–6186, 2021.
- [35] Yang, G. L.; Tang, H.; Ding, M. L.; Sebe, N.; Ricci, E. Transformer-based attention networks for continuous pixel-wise prediction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 16249–16259, 2021.
- [36] Johnston, A.; Carneiro, G. Self-supervised monocular trained depth estimation using self-attention and discrete disparity volume. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 4755–4764, 2020.
- [37] Guo, M.; Xu, T.; Liu, J.; Liu, Z.; Jiang, P.; Mu, T.; Zhang, S.; Martin, R. R.; Cheng, M.; Hu, S. Attention mechanisms in computer vision: A survey. *Computational Visual Media* Vol. 8, No. 3, 331–368, 2022.
- [38] Kendall, A.; Martirosyan, H.; Dasgupta, S.; Henry, P.; Kennedy, R.; Bachrach, A.; Bry, A. End-to-end learning of geometry and context for deep stereo regression. In: Proceedings of the IEEE International Conference on Computer Vision, 66–75, 2017.
- [39] Jaderberg, M.; Simonyan, K.; Zisserman, A.; Kavukcuoglu, K. Spatial transformer networks. In: Proceedings of the 28th International Conference on Neural Information Processing Systems, Vol. 2, 2017–2025, 2015.
- [40] Chen, Y. H.; Schmid, C.; Sminchisescu, C. Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 7062–7071, 2019.
- [41] Wang, Z.; Bovik, A. C.; Sheikh, H. R.; Simoncelli, E. P. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing* Vol. 13, No. 4, 600–612, 2004.
- [42] Wang, C. Y.; Buenaposada, J. M.; Zhu, R.; Lucey, S. Learning depth from monocular videos using direct methods. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022–2030, 2018.

- [43] Ramamonjisoa, M.; Lepetit, V. SharpNet: Fast and accurate recovery of occluding contours in monocular depth estimation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshop, 2109–2118, 2019.
- [44] Zhao, H.; Gallo, O.; Frosio, I.; Kautz, J. Loss functions for image restoration with neural networks. *IEEE Transactions on Computational Imaging* Vol. 3, No. 1, 47–57, 2017.
- [45] Yang, Z. H.; Wang, P.; Xu, W.; Zhao, L.; Nevatia, R. Unsupervised learning of geometry from videos with edge-aware depth-normal consistency. *Proceedings of the AAAI Conference on Artificial Intelligence* Vol. 32, No. 1, 7493–7500, 2018.
- [46] Zou, Y. L.; Luo, Z. L.; Huang, J. B. DF-Net: Unsupervised joint learning of depth and flow using cross-task consistency. In: *Computer Vision – ECCV 2018. Lecture Notes in Computer Science, Vol. 11209*. Ferrari, V.; Hebert, M.; Sminchisescu, C.; Weiss, Y. Eds. Springer Cham, 38–55, 2018.
- [47] Luo, C. X.; Yang, Z. H.; Wang, P.; Wang, Y.; Xu, W.; Nevatia, R.; Yuille, A. Every pixel counts++: Joint learning of geometry and motion with 3D holistic understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 42, No. 10, 2624–2641, 2020.
- [48] Casser, V.; Pirk, S.; Mahjourian, R.; Angelova, A. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. *Proceedings of the AAAI Conference on Artificial Intelligence* Vol. 33, 8001–8008, 2019.
- [49] Garg, R.; Vijay Kumar, B. G.; Carneiro, G.; Reid, I. Unsupervised CNN for single view depth estimation: Geometry to the rescue. In: *Computer Vision – ECCV 2016. Lecture Notes in Computer Science, Vol. 9912*. Leibe, B.; Matas, J.; Sebe, N.; Welling, M. Eds. Springer Cham, 740–756, 2016.
- [50] Mehta, I.; Sakurikar, P.; Narayanan, P. J. Structured adversarial training for unsupervised monocular depth estimation. In: Proceedings of the International Conference on 3D Vision, 314–323, 2018.
- [51] Poggi, M.; Tosi, F.; Mattoccia, S. Learning monocular depth estimation with unsupervised trinocular assumptions. In: Proceedings of the International Conference on 3D Vision, 324–333, 2018.
- [52] Pillai, S.; Ambruş R.; Gaidon, A. SuperDepth: Self-supervised, super-resolved monocular depth estimation. In: Proceedings of the International Conference on Robotics and Automation, 9250–9256, 2019.
- [53] Watson, J.; Firman, M.; Brostow, G.; Turmukhambetov, D. Self-supervised monocular depth hints. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2162–2171, 2019.
- [54] Tosi, F.; Aleotti, F.; Poggi, M.; Mattoccia, S. Learning monocular depth estimation infusing traditional stereo knowledge. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 9791–9801, 2019.
- [55] Li, R. H.; Wang, S.; Long, Z. Q.; Gu, D. B. UnDeepVO: Monocular visual odometry through unsupervised deep learning. In: Proceedings of the IEEE International Conference on Robotics and Automation, 7286–7291, 2018.
- [56] Ramamonjisoa, M.; Firman, M.; Watson, J.; Lepetit, V.; Turmukhambetov, D. Single image depth prediction with wavelet decomposition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 11084–11093, 2021.
- [57] Cordts, M.; Omran, M.; Ramos, S.; Rehfeld, T.; Enzweiler, M.; Benenson, R.; Franke, U.; Roth, S.; Schiele, B. The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 3213–3223, 2016.
- [58] Geiger, A.; Lenz, P.; Urtasun, R. Are we ready for autonomous driving? The KITTI vision benchmark suite. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 3354–3361, 2012.
- [59] Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; Lerer, A. Automatic differentiation in pytorch. In: Proceedings of the 31st Conference on Neural Information Processing Systems, 2017.
- [60] Kingma, D. P.; Ba, J. Adam: A method for stochastic optimization. In: Proceedings of the International Conference on Learning Representations, 2015.
- [61] Saxena, A.; Sun, M.; Ng, A. Y. Make3D: Learning 3D scene structure from a single still image. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 31, No. 5, 824–840, 2009.
- [62] Karsch, K.; Liu, C.; Kang, S. B. Depth transfer: Depth extraction from video using non-parametric sampling. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 36, No. 11, 2144–2158, 2014.
- [63] Liu, M. M.; Salzmann, M.; He, X. M. Discrete-continuous depth estimation from a single image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 716–723, 2014.

- [64] Laina, I.; Rupprecht, C.; Belagiannis, V.; Tombari, F.; Navab, N. Deeper depth prediction with fully convolutional residual networks. In: Proceedings of the 4th International Conference on 3D Vision, 239–248, 2016.
- [65] Mur-Artal, R.; Montiel, J. M. M.; Tardós, J. D. ORB-SLAM: A versatile and accurate monocular SLAM system. *IEEE Transactions on Robotics* Vol. 31, No. 5, 1147–1163, 2015.



Yuanzhen Li is working towards a Ph.D. degree in the School of Computer Science, Wuhan University. Her research interests include image editing and computer vision.



Fei Luo received his B.Sc. degree from the School of Computer Science of Hubei University of Technology in 2003. He received his M.Sc. and Ph.D. degrees from the School of Computer Science of Wuhan University in 2005 and 2008, respectively. He is now an assistant professor at the School of Computer Science, Wuhan University, Wuhan, China. In 2009, he worked as a research assistant at the School of Computer Engineering of Nanyang Technological University, Singapore. From December 2012 to December 2014, he worked as a postdoc at the Human Polymorphism Study Center, Paris, France. His research interests include data mining and computer vision.



Chunxia Xiao received his B.Sc. and M.Sc. degrees from the Mathematics Department of Hunan Normal University in 1999 and 2002, respectively, and his Ph.D. degree from the State Key Lab of CAD&CG of Zhejiang University in 2006. Currently, he is a professor at the School of Computer Science, Wuhan University.

From October 2006 to April 2007, he worked as a postdoc in the Department of Computer Science and Engineering, Hong Kong University of Science and Technology, and from February 2012 to February 2013, he visited the University of California Davis for 1 year. His main interests include computer graphics, computer vision, virtual reality, and augmented reality.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made.

The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Other papers from this open access journal are available free of charge from <http://www.springer.com/journal/41095>. To submit a manuscript, please go to <https://www.editorialmanager.com/cvmj>.