GLOCALFUSE-DEPTH: FUSING TRANSFORMERS AND CNNS FOR ALL-DAY SELF-SUPERVISED MONOCULAR DEPTH ESTIMATION

Zezheng Zhang, Ryan K. Y. Chan The University of Hong Kong {zezhengz, cky166}@connect.hku.hk

Kenneth K. Y. Wong The University of Hong Kong, Advanced Biomedical Instrumentation Centre kywong@eee.hku.hk

ABSTRACT

In recent years, self-supervised monocular depth estimation has drawn much attention since it frees of depth annotations and achieved remarkable results on standard benchmarks. However, most of existing methods only focus on either daytime or nighttime images, thus their performance degrades on the other domain because of the large domain shift between daytime and nighttime images. To address this problem, in this paper we propose a two-branch network named GlocalFuse-Depth for self-supervised depth estimation of all-day images. The daytime and nighttime image in input image pair are fed into the two branches: CNN branch and Transformer branch, respectively, where both fine-grained details and global dependency can be efficiently captured. Besides, a novel fusion module is proposed to fuse multi-dimensional features from the two branches. Extensive experiments demonstrate that GlocalFuse-Depth achieves state-of-the-art results for all-day images on the Oxford RobotCar dataset, which proves the superiority of our method.

1 Introduction

Monocular depth estimation has been one of the most basic tasks in computer vision field, which is widely applied in various tasks, such as simultaneous localization and mapping (SLAM) [1, 2, 3, 4], augmented reality (AR) [5] and autonomous driving [6]. However, it is inherently an ill-posed problem: a single image could be produced from an infinite number of distinct 3D scenes, thus the depth map can be completely different. Benefiting from the advance of deep learning method and availability of large-scale training data, significant improvement has been achieved in monocular depth estimation in recent years. A variety of deep learning methods have manifested their effectiveness to recover the pixel-level dense depth map from an single image in an end-to-end manner [7, 8, 9].

However, collecting large-scale training datasets with accurate ground-truth depth maps for supervised learning is costly, due to the large size and high prize of depth sensors (like RGB-D cameras and LIDAR). In addition, the internal error and noise characteristics of depth sensors will also affect the learning of mapping between object appearance and their depth. Therefore, self-supervised depth estimation methods have been introduced. These methods use spatial consistency (stereo depth reconstruction) or temporal consistency (monocular video depth reconstruction) as supervisory signal during the training process. Some of them [10, 11, 12, 13, 14] have achieved performance comparable to the supervised methods on widely used benchmarks, such as KITTI [15] and Cityscapes [16].

Even so, most of current self-supervised depth estimation methods only handle daytime image depth estimation problem. They fail to generalize well on nighttime images owing to the large domain shift between daytime and nighttime images and not suitable for all-day tasks such as automatic driving. Depth estimation of nighttime images are challenging owing to varying illuminations and low visibility at nighttime. One solution is to apply image-to-image translation methods, such as CycleGAN [17], to translate nighttime images into daytime images, then use a pretrained daytime

model to estimate depth from these translated images. Unfortunately, due to the large domain shift between the two domain, it is difficult to obtain natural daytime images, then the performance is also limited. Fig. 1 (b) demonstrate the depth estimation results of Monodepth2 [18], a effective self-supervised daytime depth estimation approach, on nighttime images. Fig. 1 (c) demonstrate the depth estimation results of Monodepth2 [18], specially edges, are failed to be estimated due to the varying illuminations and low visibility of nighttime images.



Figure 1: Comparison with other methods on Oxford RobotCar dataset [19]. From left to right: (a) Nighttime Images,(b) Monodepth2 [18], (c) Monodepth2+CycleGAN [17], (d) ADDS-DepthNet [20], (e) GlocalFuse-Depth.

For images captured with a fixed viewpoint, their depth map remains the same, regardless of other terms changing with time such as illumination. Also, [21] proves that texture information plays more important roles in depth estimation than exact color information. To cater to above issues, inspired by [20] (ADDS-DepthNet, depth estimation results shown in Fig. 1 (d)), here we propose GlocalFuse-Depth, which adopts the day-night image pair (with same depth map, generated by CycleGAN [17]) as input to feed into two branches: CNN branch and Transformer branch. We choose this two-branch network architecture for the following reasons: although CNN is proved powerful at building hierarchical feature representation, its lack of efficiency in capturing global context information remains a challenge. CNN captures global information usually at the expense of efficiency, which requires stacked convolutional layers until the receptive field becomes large enough. While Transformer is good at modeling global context, it shows limitations in capturing fine-grained details due to the lack of spatial inductive-bias. In order to take full advantage of both structures, they are both adopted as two branches in GlocalFuse-Depth to capture both local and global feature (where "Glocal" comes from). Output features from the two branches with same dimension are fused to jointly make predictions. As shown in Fig. 1 (e), our method effectively relieves the problems of varying illuminations and low visibility, and achieves more attractive results for nighttime images.

The main contributions of this paper can be summarized as:

- We propose a two-branch network GlocalFuse-Depth for self-supervised all-day image depth estimation. It adopts complementary encoders (CNNs and Transformers) to capture local and global information in a day-night image pair to better grasp the texture correspondence between RGB image and depth map. To the best of our knowledge, GlocalFuse-Depth is the first model synthesizing CNN and Transformer for all-day image depth estimation.
- A new fusion module is proposed. It combines and aggregates the information from the two branches to obtain a comprehensive representation for the final depth estimation. Both channel-wise and spatial-wise attention are applied to boost representation power of our two-branch network.
- Experimental results of our method outperform other state-of-the-art methods on the Oxford RobotCar dataset for all-day image depth estimation, which confirms the superiority of our method.

2 Related Work

2.1 Daytime Depth Estimation

In view of the shortcomings of supervised depth estimation, self-supervised depth estimation has been extensively studied in recent years. There are two main types of self-supervised methods: stereo depth reconstruction and

monocular video depth reconstruction. [11] is the pioneering work of stereo depth reconstruction. They use the left–right consistency constraints as the supervisory signal to train the depth estimation network. On the basis, other researchers make a series of improvements on this work. For example, [22] introduces an adaptive regularization scheme for the network to better handle the co-visible and occluded regions in a stereo pair. [23] proposes a sparsity-invariant autoencoder to improve the traditional visual odometry. On the other hand, [12] is the first self-supervised depth estimation method using geometry cues between monocular video frames as supervisory signal, which trains the depth network along with a separate pose network. This work provided many useful references for subsequent works. [24] trains their network with a 3D-based loss under the supervision of geometry cues from adjacent frames of monocular video. [13] incorporates self-attention and discrete disparity prediction into their network. These methods have achieved promising results on daytime benchmarks, such as KITTI [15] and Cityscapes [16]. However, self-supervised depth estimation approaches for all-day images have not been well addressed due to the challenge of varying illuminations and low visibility of nighttime images, and the performance of most current methods degrades a lot on nighttime images due to large domain shift between daytime and nighttime images.

2.2 Nighttime Depth Estimation

Self-supervised depth estimation for nighttime images is still a relatively underexplored topic even with today's complete deep learning method. Different approaches have also been proposed to address this challenging task from both hardware and software side. [25] and [26] benefit from thermal images acquired by additional sensors. [27] is proposed to simultaneously learn cross-domain feature representation and depth estimation to acquire more robust supervision for nighttime images. [28] utilizes a translation network which renders nighttime stereo images from daytime stereo images. [29] adapts a network trained on daytime images to work for nighttime images, aiming to transfer knowledge from daytime images to nighttime images. [30] propose Priors-Based Regularization to learn distribution knowledge from unpaired depth maps and Mapping-Consistent Image Enhancement module to enhance image visibility and contrast. [20] propose a domain-separated framework which partition the information of day-night image pairs into two complementary sub-spaces and relieve the influence of disturbing terms for all-day depth estimation.

Though remarkable progress has been achieved, the large domain shift between daytime and nighttime images is always hard to be fixed. It is difficult to obtain satisfactory depth map of daytime and nighttime images in one single network. To solve this problem, we propose the two branch network GlocalFuse-Depth for for all-day self-supervised depth estimation, which can effectively issue the problem of inherent domain gap between daytime and nighttime images.

3 Approach

3.1 Input Image Pair and Self-supervised Depth Estimation Objective

For the daytime and nighttime images of the same scene, the depth information should be consistent, although their light conditions are quite different. However, it is almost impossible to guarantee a real-world daytime and nighttime image pair contain identical objects even they were captured at the same position with a fixed viewpoint, since there are always moving objects in outdoor scenes. This will mislead the network when mapping the RGB information with the depth information. Therefore, a pretrained CycleGAN [12] is used to translate daytime images to nighttime images before training, then a daytime and the corresponding nighttime image compose a day-night image pair. This pre-processing method ensures that the two images in an image pair contain identical depth information, and they can be fed into the two branches of our network respectively with the same supervisory signal when training. Note that other image translation methods can also be used here. During the inference process, CycleGAN [12] is also used to translate the input daytime or nighttime image to the other domain and form the day-night image pair for the network to estimate depth.

As for the optimization objective, similar to prior work, we formulate the self-supervised depth estimation task as a monocular video depth reconstruction problem. The goal is to minimize the photometric reprojection error [18] at training stage. For two frames I_t and $I_{t'}$ in an image sequence, the reprojection process is:

$$I_{t} = KT_{t' \to t} D(t') K^{-1} I_{t'}, \tag{1}$$

where D(t') is the depth map of the source frame $I_{t'}$, $T_{t' \to t}$ represents the spatial transformation from $I_{t'}$ to the target frame I_t (estimated by a pose network) and K is the camera intrinsic matrix. The photometric loss can be formulated as:

$$L_p = L_p^d + L_p^n,\tag{2}$$

$$L_p^{d/n} = \frac{\alpha}{2} (1 - SSIM(I_t^{d/n}, \hat{I}_t^{d/n})) + (1 - \alpha) ||(I_t^{d/n} - \hat{I}_t^{d/n})||_1,$$
(3)

where I represents the reconstructed target image obtained by Eq. (1) and I represents the original image. The superscript d and n represent daytime and nighttime images, respectively and α is empirically set as 0.85. Here we use the two frames temporally adjacent to I_t as our source frames, *i.e.* I_{t-1} and I_{t+1} .

3.2 Network Architecture



Figure 2: Overview of GlocalFuse-Depth architecture: two parallel branches - CNN branch (bottom right) and Transformer branch (top) fused by our proposed fusion module.

As shown in Fig. 2, GlocalFuse-Depth consists of two parallel branches: CNN branch and Transformer branch, which takes day-night image pairs as input (daytime images fed into the CNN branch and nighttime images fed into the Transformer branch). The CNN branch gradually increases the receptive field and encodes features from local to global, and the Transformer branch starts with global self-attention and recovers the local details by upsampling. Extracted features with same spatial resolutions from two branches are fed into our proposed fusion module, where branch channel selection and spatial attention are applied to selectively fuse the information from both branches. Then the multi-level fused feature maps are combined to generate the final prediction using attention-gated skip-connection [31]. There are two main benefits of our proposed architecture: firstly, with the characteristic of grasping long-range information of transformer, GlocalFuse-Depth can capture global depth information without stacking a large number of convolutional layers. At the same time, the use of CNN helps to preserve sensitivity of the network on local context. Secondly, our proposed fusion module comprises channel-selection and spatial attention for features from the two branches, which exploits different characteristics of the output feature from two branches, thus making the fused representation comprehensive and compact.

The forward propagation process of the Transformer branch is as follows. Firstly, the input image $x \in \mathbb{R}^{H*W*3}$ is divided into $N = \frac{H}{P} * \frac{W}{P}$ patches, where the patch size P is typically set as 16. These patches are then flattened and fed into a linear embedding layer with latent size D_0 and obtain the embedding sequence $E \in \mathbb{R}^{N*D_0}$. In order to utilize the spatial prior, a learnable positional embedding with the same dimension as E is added. The result vector $z^0 \in \mathbb{R}^{N*D_0}$ is fed into L consecutive transformer encoder block. In the transformer encoder, self-attention (SA) mechanism is applied as

$$SA(z^{i}) = softmax(\frac{qk^{T}}{\sqrt{D_{h}}})v, \tag{4}$$

where $[q, k, v] = z^i W_{qkv}$, $W_{qkv} \in \mathbb{R}^{D_0 * 3D_h}$ is the linear projection matrix and D_h is the dimension of q, k and v. Multi-head self-attention (MSA) is an extension of SA in which runs k self-attention operations in parallel and project the concatenated outputs back to \mathbb{R}^{D_0} . MSA is the basis of transformer for learning long-range dependencies. Layer normalization and MLP is followed to obtain the encoded sequence $z^L \in \mathbb{R}^{N*D_0}$ (refer to [32] for other details). In the decoder part, we first reshape z^L to $t^0 \in \mathbb{R}^{\frac{H}{16}*\frac{W}{16}*C_0}$, then use two consecutive upsampling layer to recover its spatial resolution, obtaining $t^1 \in \mathbb{R}^{\frac{H}{8}*\frac{W}{8}*C_1}$ and $t^2 \in \mathbb{R}^{\frac{H}{4}*\frac{W}{4}*C_2}$. These three feature maps with different dimension are saved for late fusion with the outputs of CNN branch.

In former works, to obtain global information, CNN encoders are always designed very deep to decrease the dimension of an image, which takes up massive computational resources and increases training difficulty. Considering that global information can be grasped by the Transformer branch, here we only adopt the first 4 blocks of ResNet34 as our CNN branch. This also benefits the CNN branch from retaining richer local information. We take the outputs from the 4th $(g^0 \in \mathbb{R}^{\frac{H}{16}*\frac{W}{16}*C_0})$, 3rd $(g^1 \in \mathbb{R}^{\frac{H}{8}*\frac{W}{8}*C_1})$ and 2nd $(g^2 \in \mathbb{R}^{\frac{H}{4}*\frac{W}{4}*C_2})$ blocks to fuse with the outputs from the Transformer branch.

To generate final segmentation, fused feature with different scales are combined using the attention-gated (AG) skipconnection [31]. Features with smaller spatial dimension are upsampled and used as the gating signal to allow the model to focus more on prominent regions while suppressing feature activations in irrelevant regions at late stage.

3.3 Fusion Block



Figure 3: The proposed fusion module.

To effectively combine the encoded features from the CNN branch and the Transformer branch, we propose a new fusion module (as in Fig. 3). For each scale of output of the CNN and the Transformer branch, inspired by [33], we first select representative channel of the two outputs for the final estimation. Outputs from the two branches is aggregated via an element-wise summation:

$$\hat{u}^i = t^i + g^i \tag{5}$$

A gate vector z^i used to control information flows from the two branches is computed by consecutive global average pooling (GAP) layer and fully connected (FC) layer.

$$z^{i} = FC(GAP(\hat{u}^{i})) \tag{6}$$

GAP layer is applied to embed the global information of the summation feature \hat{u}_i and FC layer is used for better efficiency by reducing the dimension. After that, a soft attention vector, guided by the gate vector z^i , is computed to select channels of the features from the two branches:

$$s_t^i = \frac{e^{W_t^i * z^i}}{e^{W_t^i * z^i} + e^{W_g^i * z^i}} \qquad s_c^i = \frac{e^{W_g^i * z^i}}{e^{W_t^i * z^i} + e^{W_g^i * z^i}}$$
(7)

Besides, the spatial attention block adopted from CBAM [34], which is complementary to the soft channel attention, is used as spatial filters to highlight informative regions and suppress irrelevant regions for both branches: outputs of max pooling and average pooling along the channel axis are concatenated and fed into a convolutional layer to produce the spatial attention map:

$$\tilde{t}^{i} = Conv([Max Avg]Pool(t^{i})) \qquad \tilde{g}^{i} = Conv([Max Avg]Pool(g^{i})) \tag{8}$$

Then we multiply soft channel-selection vector $\tilde{t}_i(\tilde{g}_i)$, spatial attention vector $s_t(s_g)$ and original outputs $t_i(g_i)$ from each branch. Finally, we perform element-wise summation to get the final fuse result for each scale:

$$t^{i} = \tilde{t}^{i} * s^{i}_{t} * t^{i} \qquad g^{i} = \tilde{g}^{i} * s^{i}_{q} * g^{i}$$

$$\tag{9}$$

$$u^i = t^i + g^i \tag{10}$$

4 Experiments

In this section, we compare the performance of our method with state-of-the-art methods on Oxford RobotCar dataset [19], for both daytime and nighttime images.

4.1 Dataset

The KITTI [15] and Cityscapes [16] datasets are widely used in depth estimation task. However, these two datasets only consist of daytime images, which can not meet the requirements for all-day depth estimation. Therefore, we choose the Oxford RobotCar dataset [19] for training and testing in this work. Oxford RobotCar dataset is a large-scale dataset captured by cameras and sensors with fixed position and viewpoint during a long-term autonomous driving, which includes both daytime and nighttime images. Following [20], we use the left images collected by the front stereo-camera (Bumblebee XB3) with the resolution of 960 * 1280 for self-supervised depth estimation. Sequence "2014-12-09-13-21-02" and "2014-12-16-18-44-24" are used for training of daytime and nighttime, respectively. Both training data are selected from the first 5 splits. The testing images are collected from the other splits of the Oxford RobotCar dataset, which contains 451 daytime images and 411 nighttime images. We use the depth data captured by the front LMS-151 depth sensors as the ground truth in the testing phase. The images are first center-cropped to 640 * 1280, then resized to 256 * 512 as the inputs of the network.

4.2 Implementation Details

At the training phase, the pretrained CycleGAN[17] is used to translate daytime images to nighttime images. The generated 'fake' nighttime images and the corresponding daytime images compose the image pairs and fed into GlocalFuse-Depth.

The network is trained with an NVIDIA GeForce RTX3090 GPU for 30 epochs. The parameters are optimized with Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.99$) and the batch size is set as 16. We also used learning rate warmup scheme to increase the training stability: the learning rate increases to 1e-5 linearly in first 5 epoches, then decreases as the cosine function in remaining epoches.

4.3 Quantitative Results

Table. 1 demonstrates the quantitative comparison results between our method and some state-of-the-art methods. The maximum depth is set to be 60m. In Table. 1, Monodepth2 [18] (day) and Monodepth2 [18] (night) mean training with

Table 1: Quantitative comparison with state-of-the-art methods. The maximum depth is set to be 60m. Lower value is better for the first four metrics, higher value is better for the other three. The best results are presented in **bold** for each metric. Monodepth2 [18] (day) and Monodepth2 [18] (night) mean training with daytime and nighttime images of the Oxford dataset, respectively. Monodepth2+CycleGAN [17] means training the Monodepth2 model with daytime image, then using 'fake' daytime images, which is translated from the nighttime images, to test its performance.

Method (test at night)	AbsRel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	δ <1.25 ²	δ <1.25 ³
Monodepth2[18] (day)	0.440	6.039	14.148	0.481	0.365	0.645	0.840
Monodepth2[18] (night)	0.513	8.266	16.668	0.528	0.575	0.860	0.917
Monodepth2+CycleGAN[17]	0.240	2.341	9.082	0.281	0.636	0.880	0.955
ADFA[29]	0.233	3.783	10.089	0.319	0.668	0.844	0.924
RNW-Net [30]	0.242	3.496	11.962	0.323	0.637	0.858	0.939
ADDS-DepthNet [20]	0.231	2.674	8.800	0.286	0.620	0.892	0.956
GlocalFuse-Depth	0.217	2.299	8.520	0.261	0.672	0.901	0.961
Method(test at day)	AbsRel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	δ <1.25 ²	δ <1.25 ³
Monodepth2[18] (day)	0.125	1.248	6.634	0.190	0.822	0.948	0.985
Monodepth2[18] (night)	0.305	3.648	11.592	0.396	0.459	0.745	0.885
ADDS-DepthNet [20]	0.115	0.794	4.855	0.168	0.863	0.967	0.989
GlocalFuse-Depth	0.113	0.785	4.787	0.165	0.871	0.969	0.989

daytime and nighttime images of the Oxford dataset, respectively. Monodepth2+CycleGAN [17] means training the Monodepth2 with daytime images, then using 'fake' daytime images, which are translated from the nighttime images, to test its performance. Monodepth2 [18] is proved to be an effective self-supervised depth estimation method for daytime images as shown in Table. 1. However, its performances are limited for nighttime images because of the large domain shift between daytime and nighttime images. Besides, because of varying illumination and the low visibility, different degrees of information is lost due to too bright or too dark regions on nighttime images, which causes that training directly on nighttime images also cannot get good enough results.

Meanwhile, although combining Monodepth2 [18] and CycleGAN [17] could transform the nighttime images to 'fake' daytime images and reduce the domain shift between train and test images, the performances are also limited due to the loss of CycleGAN itself. ADFA [29] and RNW-Net [30], which is specialized to estimate depth of nighttime images, could reduce the domain shift between day and night images at the feature level, but the performance is limited by daytime results. ADDS-DepthNet [20] also use day-night image pair as input and could improve the depth estimation results of both the daytime and nighttime images to a certain extent, its performance can still be improved. As shown in Table. 1, the use of complementary encoders and effective fuse module of our GlocalFuse-Depth could grasp inherent texture features from both daytime and nighttime images, then improve the depth estimation performance for all-day images. Almost all the evaluation metrics for daytime and nighttime images are largely improved by our approach, which proves the superiority of our method.

4.4 Qualitative Results

The qualitative nighttime image depth estimation results comparison of our method with some state-of-the-art methods are shown in Fig. 4, where (b) shows the CycleGAN [17] generated 'fake' daytime images from nighttime images, (c) shows the depth estimation results of Monodepth2 [18] trained with daytime images and tested with nighttime images, (d) shows the results of Monodepth2 [18] trained with daytime images and tested with 'fake' daytime images. Compared with (c), it is obvious that (d) obtains better depth estimation results, which proves that the use of CycleGAN [17] could improve the performance of nighttime image depth estimation substantially. (e) shows the result of RNW [30] and (f) shows the result of ADDS-DepthNet [20]. (g) shows our result. Compared with (e) and (f), our result could obtain more accurate segmentation on some objects, such as the car on the left of the first and the third nighttime image. Besides, our method could be unaffected from some distortion of the nighttime image. In the second row, a defect line caused by image processing on nighttime images appears on ADDS-DepthNet [20] estimation result. However, it disappears on our result, which proves the robustness of our method. Different from the first three row, which shows the scene of the car driving on a straightway, the nighttime image of the fourth row shows a rarely seen scene in training set: car turning. Our method could recover more details on the cars, while RNW [30] and ADDS-DepthNet [20] fails to do so, which proves the generalizability of our method.



Figure 4: Qualitative depth estimation result comparison with other state-of-the-art methods of nighttime images. From left to right: (a) Nighttime images, (b) 'Fake' daytime images generated by CycleGAN [17], (c) Monodepth2 [18], (d) Monodepth2+CycleGAN[17], (e) RNW-Net [30], (f) ADDS-DepthNet [20], (g) GlocalFuse-Depth.

Fig. 5 demonstrates the qualitative daytime result comparison with other methods. We can see that more depth details can be recovered by our method, such as the tree on the left of the first daytime image and the car on the middle right of the second image, which clarifies the depth estimation of our method for daytime images.



Figure 5: Qualitative depth estimation result comparison with other state-of-the-art methods of daytime images. From left to right: (a) Daytime images, (b) Monodepth2 [18], (c) ADDS-DepthNet [20], (d) GlocalFuse-Depth.

4.5 Ablation Study

Here, we conduct a series of experiments to demonstrate the effectiveness of GlocalFuse-Depth and report the results in Table 2. For the sake of simplicity, we quantitatively compare the performance using AbsRel.

Firstly, we compare the depth estimation performance with different encoder design. In Table. 2, the quantitative depth estimation result of our proposed method (CNN-Transformer as encoders), CNN only as encoder in both branches and transformer only as encoder in both branches is shown. As can be seen, 5.2% and 0.4% improvement of AbsRel is achieved by using both CNN and transformer as encoders compared to using transformer only for nighttime and daytime images, respectively. Similarly, 12.9% and 15.7% improvement of AbsRel is achieved compared to using CNN only. These result proves the effectiveness of use complementary encoders in extracting local and global information from the images.

Secondly, we show the performance improvement on using day-night image pair as input. In Table. 2, the quantitative result of our proposed method (day-night image pair), using day-day image pair and using night-night image pair as input is shown. Compared with using day-day image pair, 6.5% and 0.3% performance improvement of AbsRel

Study subject	Method (night)	AbsRel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	δ <1.25 ²	δ <1.25 ³
	GlocalFuse-Depth	0.217	2.299	8.520	0.261	0.672	0.901	0.961
Encoder design	Transformer only	0.229	3.438	11.724	0.326	0.641	0.863	0.933
	CNN only	0.249	3.783	12.635	0.351	0.575	0.846	0.925
Image pair	Day-day	0.232	3.626	11.853	0.329	0.649	0.862	0.931
	Night-night	0.475	5.672	12.069	0.494	0.348	0.620	0.812
Fusion method	Bifusion module[35]	0.229	3.559	11.863	0.331	0.650	0.863	0.930
	Concatenation	0.248	3.658	12.372	0.340	0.584	0.851	0.930
	Dot product	0.286	3.573	12.149	0.339	0.502	0.809	0.950
Study subject	Method (day)	AbsRel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
	GlocalFuse-Depth	0.113	0.785	4.787	0.165	0.871	0.969	0.989
Encoder design	Transformer only	0.118	1.139	6.418	0.181	0.847	0.957	0.985
	CNN only	0.134	1.404	7.181	0.206	0.805	0.940	0.980
Image pair	Day-day	0.117	1.225	6.637	0.186	0.839	0.951	0.983
	Night-night	0.266	4.344	13.014	0.370	0.570	0.832	0.917
Fusion method	Bifusion module[35]	0.123	1.256	6.678	0.190	0.834	0.952	0.983
	Concatenation	0.137	1.454	7.388	0.211	0.796	0.936	0.978
	Dot product	0.218	2.285	8.933	0.276	0.649	0.879	0.955

Table 2: Ablation study on encoder design, composition of input image pair and fusion method.

is achieved by our method for nighttime and daytime images. Besides, 54.3% and 57.5% improvement of AbsRel for nighttime and daytime images is achieved compared with using night-night image pair. The results above proves that the use of day-night image pair could help the network extract domain invariant information of both daytime and nighttime images, thus the estimation performance is improved.

Finally, we compare the performance our proposed fusion module with other fusion methods. In Table. 2, compared with the Bifusion module in [35], simply do the concatenation and dot product, our proposed fusion module achieves improvement of AbsRel for nighttime and daytime images by 5.2% and 8.1%, 12.5% and 17.5% and 24.1% and 48.2%, respectively. These results is owing to the use of effective channel-selection and spatial attention in our proposed module.

5 Conclusion

In this paper, we propose a two-branch network GlocalFuse-Depth which combines CNN and Transformer with late fusion for self-supervised monocular depth estimation of all-day images. The resulting architecture leverages the inherent characteristics of CNNs on modeling fine-grained feature and the capability of Transformers on modelling global context. A novel fusion module is proposed to selectively fuse multi-dimensional features from the two branches with channel-selection and spatial attention. Experiments results on the Oxford RobotCar dataset demonstrate that GlocalFuse-Depth achieves state-of-the-art results for all-day images.

Acknowledgments

Research Grants Council of the Hong Kong Special Administrative Region of China (HKU 17210522, HKU C7074-21G, HKU 17205321, HKU 17200219, HKU 17209018, CityU T42-103/16-N) and Health@InnoHK program of the Innovation and Technology Commission of the Hong Kong SAR Government.

References

 Gibson Hu, Shoudong Huang, Liang Zhao, Alen Alempijevic, and Gamini Dissanayake. A robust rgb-d slam algorithm. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 1714–1719, 2012.

- [2] Raúl Mur-Artal and Juan D. Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, 2017.
- [3] Keisuke Tateno, Federico Tombari, Iro Laina, and Nassir Navab. Cnn-slam: Real-time dense monocular slam with learned depth prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), pages 6243–6252, July 2017.
- [4] Lina Liu, Xibin Song, Xiaoyang Lyu, Junwei Diao, Mengmeng Wang, Yong Liu, and Liangjun Zhang. Fcfr-net: Feature fusion based coarse-to-fine residual learning for depth completion. *Proceedings of the AAAI Conference* on Artificial Intelligence, 35(3):2136–2144, May 2021.
- [5] Richard A. Newcombe, Steven J. Lovegrove, and Andrew J. Davison. Dtam: Dense tracking and mapping in real-time. In 2011 International Conference on Computer Vision, pages 2320–2327, 2011.
- [6] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3061–3070, June 2015.
- [7] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2002–2011, June 2018.
- [8] Wei Yin, Yifan Liu, Chunhua Shen, and Youliang Yan. Enforcing geometric constraints of virtual normal for depth prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5684–5693, October 2019.
- [9] Zhenyu Zhang, Zhen Cui, Chunyan Xu, Yan Yan, Nicu Sebe, and Jian Yang. Pattern-affinitive propagation across depth, surface normal and semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4106–4115, June 2019.
- [10] Ravi Garg, Vijay Kumar B.G., Gustavo Carneiro, and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 740–756, Cham, 2016. Springer International Publishing.
- [11] Clement Godard, Oisin Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (*CVPR*), pages 270–279, July 2017.
- [12] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and egomotion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1851–1858, July 2017.
- [13] Adrian Johnston and Gustavo Carneiro. 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 (CVPR)*, pages 4756–4765, June 2020.
- [14] Yunhan Zhao, Shu Kong, Daeyun Shin, and Charless Fowlkes. Domain decluttering: Simplifying images to mitigate synthetic-real domain shift and improve depth estimation. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition (CVPR), pages 3330–3340, June 2020.
- [15] A Geiger, P Lenz, C Stiller, and R Urtasun. Vision meets robotics: The kitti dataset. The International Journal of Robotics Research, 32(11):1231–1237, 2013.
- [16] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 3213–3223, June 2016.
- [17] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision* (ICCV), pages 2223–2232, Oct 2017.
- [18] Clement Godard, Oisin Mac Aodha, Michael Firman, and Gabriel J. Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (ICCV), pages 3828–3838, October 2019.
- [19] Will Maddern, Geoffrey Pascoe, Chris Linegar, and Paul Newman. 1 year, 1000 km: The oxford robotcar dataset. *The International Journal of Robotics Research*, 36(1):3–15, 2017.
- [20] Lina Liu, Xibin Song, Mengmeng Wang, Yong Liu, and Liangjun Zhang. Self-supervised monocular depth estimation for all day images using domain separation. In *Proceedings of the IEEE/CVF International Conference* on Computer Vision (ICCV), pages 12737–12746, October 2021.

- [21] Tom van Dijk and Guido de Croon. How do neural networks see depth in single images? In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2183–2191, October 2019.
- [22] Alex Wong and Stefano Soatto. Bilateral cyclic constraint and adaptive regularization for unsupervised monocular depth prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), pages 5644–5653, June 2019.
- [23] Lorenzo Andraghetti, Panteleimon Myriokefalitakis, Pier Luigi Dovesi, Belen Luque, Matteo Poggi, Alessandro Pieropan, and Stefano Mattoccia. Enhancing self-supervised monocular depth estimation with traditional visual odometry. In 2019 International Conference on 3D Vision (3DV), pages 424–433, 2019.
- [24] Reza Mahjourian, Martin Wicke, and Anelia Angelova. Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5667–5675, June 2018.
- [25] Namil Kim, Yukyung Choi, Soonmin Hwang, and In So Kweon. Multispectral transfer network: Unsupervised depth estimation for all-day vision. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1):6983– 6991, Apr. 2018.
- [26] Yawen Lu and Guoyu Lu. An alternative of lidar in nighttime: Unsupervised depth estimation based on single thermal image. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 3833–3843, January 2021.
- [27] Jaime Spencer, Richard Bowden, and Simon Hadfield. Defeat-net: General monocular depth via simultaneous unsupervised representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14402–14413, June 2020.
- [28] Aashish Sharma, Loong-Fah Cheong, Lionel Heng, and Robby T. Tan. Nighttime stereo depth estimation using joint translation-stereo learning: Light effects and uninformative regions. In 2020 International Conference on 3D Vision (3DV), pages 23–31, 2020.
- [29] Madhu Vankadari, Sourav Garg, Anima Majumder, Swagat Kumar, and Ardhendu Behera. Unsupervised monocular depth estimation for night-time images using adversarial domain feature adaptation. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 443–459, Cham, 2020. Springer International Publishing.
- [30] Kun Wang, Zhenyu Zhang, Zhiqiang Yan, Xiang Li, Baobei Xu, Jun Li, and Jian Yang. Regularizing nighttime weirdness: Efficient self-supervised monocular depth estimation in the dark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16055–16064, October 2021.
- [31] Jo Schlemper, Ozan Oktay, Michiel Schaap, Mattias Heinrich, Bernhard Kainz, Ben Glocker, and Daniel Rueckert. Attention gated networks: Learning to leverage salient regions in medical images. *Medical Image Analysis*, 53:197–207, 2019.
- [32] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2020.
- [33] Xiang Li, Wenhai Wang, Xiaolin Hu, and Jian Yang. Selective kernel networks. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition (CVPR), pages 510–519, June 2019.
- [34] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, September 2018.
- [35] Yundong Zhang, Huiye Liu, and Qiang Hu. Transfuse: Fusing transformers and cnns for medical image segmentation. In Marleen de Bruijne, Philippe C. Cattin, Stéphane Cotin, Nicolas Padoy, Stefanie Speidel, Yefeng Zheng, and Caroline Essert, editors, *Medical Image Computing and Computer Assisted Intervention – MICCAI* 2021, pages 14–24, Cham, 2021. Springer International Publishing.