

# Learning to Fuse Monocular and Multi-view Cues for Multi-frame Depth Estimation in Dynamic Scenes

Rui Li<sup>1</sup>, Dong Gong<sup>2\*</sup>, Wei Yin<sup>3\*</sup>, Hao Chen<sup>4</sup>, Yu Zhu<sup>1</sup>, Kaixuan Wang<sup>3</sup>,  
Xiaozhi Chen<sup>3</sup>, Jinqiu Sun<sup>1</sup>, Yanning Zhang<sup>1\*</sup>

<sup>1</sup>Northwestern Polytechnical University, <sup>2</sup>The University of New South Wales, <sup>3</sup>DJI, <sup>4</sup>Zhejiang University

<https://github.com/ruili3/dynamic-multiframe-depth>

## Abstract

Multi-frame depth estimation generally achieves high accuracy relying on the multi-view geometric consistency. When applied in dynamic scenes, e.g., autonomous driving, this consistency is usually violated in the dynamic areas, leading to corrupted estimations. Many multi-frame methods handle dynamic areas by identifying them with explicit masks and compensating the multi-view cues with monocular cues represented as local monocular depth or features. The improvements are limited due to the uncontrolled quality of the masks and the underutilized benefits of the fusion of the two types of cues. In this paper, we propose a novel method to learn to fuse the multi-view and monocular cues encoded as volumes without needing the heuristically crafted masks. As unveiled in our analyses, the multi-view cues capture more accurate geometric information in static areas, and the monocular cues capture more useful contexts in dynamic areas. To let the geometric perception learned from multi-view cues in static areas propagate to the monocular representation in dynamic areas and let monocular cues enhance the representation of multi-view cost volume, we propose a cross-cue fusion (CCF) module, which includes the cross-cue attention (CCA) to encode the spatially non-local relative intra-relations from each source to enhance the representation of the other. Experiments on real-world datasets prove the significant effectiveness and generalization ability of the proposed method.

## 1. Introduction

Depth estimation is a fundamental and challenging task for 3D scene understanding in various application scenarios, such as autonomous driving [10, 16, 27]. With the advent of convolutional neural networks (CNNs) [12, 18], depth estimation methods [2, 24–26, 40, 45, 46] are capable of predicting promising results given either single or mul-

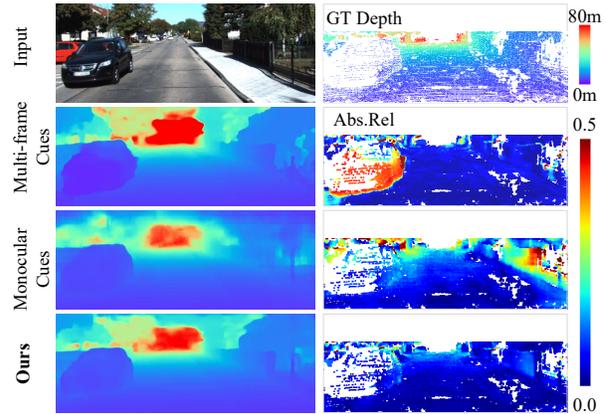


Figure 1. **Depth estimation in dynamic scenes.** Multi-frame predictions reserve high overall accuracy while degrading in dynamic areas. The monocular method better handles moving areas while suffering in static areas. Our method fuses both multi-frame and monocular depth cues for final prediction, yielding superior performance of the whole scene.

iple images. The single image-based methods learn the monocular cues, e.g., the texture or object-level features, to predict the depth [2, 42, 44], while multi-frame methods [36, 37, 40] can generally obtain higher overall accuracy relying on the multi-view geometric cues. Specifically, the 3D cost volume has been proven simple and effective for depth estimation, which encodes the multi-frame matching probabilities with a set of depth hypotheses [15, 37, 40].

Although multi-frame methods are widely used in scene reconstruction [15, 34, 40], they encounter non-negligible challenges in dynamic scenes with dynamic areas (e.g., moving cars and pedestrians). The dynamic areas cause corrupted values in the cost volume due to the violation of multi-view consistency [9, 33] and mislead the network predictions. However, depth estimation for the dynamic areas is usually crucial in most applications [10, 16, 19]. As shown in Fig. 1, multi-frame depth estimation for the dynamic cars is more challenging than the static backgrounds.

To handle the dynamic areas violating the multi-view consistency, a few multi-frame depth estimation methods [9, 36, 37] try to identify and exclude the dynamic areas

\*Corresponding author

through an *explicit* mask obtained relying on some assumptions or heuristics. Specifically, some method [37] excludes the multi-frame cost volume relying on a learned dynamic mask and compensates the excluded areas with monocular features; some methods directly adjust the dynamic object locations in input images [9] or supervise multi-frame depth [36] with predicted monocular depth. However, these methods are usually sensitive to the explicit mask’s quality, and the masks are obtained from additional networks or manually crafted criteria [9, 36, 37]. Despite better performances than pure multi-frame methods, these methods exhibit limited performance improvement compared with the additionally introduced *monocular* cues (as shown in Tab. 4), implying underutilized benefits from the *fusion* of the multi-view and monocular cues. Although some self-supervised monocular depth estimation methods [4, 5, 11, 14, 22, 23] also address the multi-view inconsistency issues, they mainly focus on handling the unfaithful self-supervision signals.

To tackle the above issues, we propose a novel method that fuses the respective benefits from the monocular and multi-view cues, leading to significant improvement upon each individual source in dynamic scenes. We first analyze the behaviors of monocular and multi-frame cues in dynamic scenes, that the pure monocular method can generally learn good structural relations around the dynamic areas, and the pure multi-view cue preserves more accurate geometric properties in the static areas. We then unveil the effectiveness of leveraging the benefits of both depth cues by directly fusing depth volumes (Sec. 3). Inspired by the above observations, beyond treating monocular cues as a local supplement of multi-frame methods [9, 36, 37], we propose a *cross-cue fusion* (CCF) module to enhance the representations of multi-view and monocular cues with the other, and fuse them together for dynamic depth estimation. We use the spatially non-local relative intra-relations encoded in *cross-observation attention* (CCA) weights from each source to guide the representation of the other, as shown in Fig. 4. Specifically, the intra-relations of monocular cues can help to address multi-view inconsistency in dynamic areas, while the intra-relations from multi-view cues help to enhance the geometric property of the monocular representation, as visualized in Fig. 5. Unlike [1, 9, 36, 37], the proposed method unifies the input format of both cues as volumes and conducts fusion on them, which achieves better performances (as shown in Fig. 6). The proposed fusion module is learnable and does not require any heuristic masks, leading to better generalization and flexibility.

Our main contributions are summarized as follows:

- We analyze multi-frame and monocular depth estimations in dynamic scenes and unveil their respective advantages in static and dynamic areas. Inspired by this, we propose a novel method that fuses depth volumes from each cue to achieve significant improve-

ment upon individual estimations in dynamic scenes.

- We propose a *cross-cue fusion* (CCF) module that utilizes the *cross-cue attention* to encode non-local intra-relations from one depth cue to guide the representation of the other. Different from methods using local masks, the attention weights learn mask-free global geometric information according to the geometric properties of each depth cue (as shown in Fig. 5).
- The proposed method outperforms the state-of-the-art method in dynamic areas with a significant error reduction of 21.3% while retaining its superiority in overall performance on KITTI. It also achieves the best generalization ability on the DDAD dataset in dynamic areas than the competing methods.

## 2. Related Work

**Learning-based multi-frame depth estimation.** Learning depth estimation from multiple images has attracted much attention these years. Current methods [7, 15, 32, 39–41] typically construct a cost volume using homography warping [40] between multi-view images. Under the premise that the scene is static, the cost volume encodes the probabilities of different depth hypotheses for each pixel, which can be regularized by 3D CNNs to yield final depth prediction. Aiming at recovering the accurate structure of static scenes, many endeavors have been made in improving the quality of the cost volume [15, 28], network efficiency [15, 39, 41] as well as the temporal consistency [7, 32], *etc.* The effectiveness of static scene reconstruction has sparked several attempts [17, 36, 37] that extend the multi-frame depth estimation to the large-scale outdoor scenes. The cost-volume between temporal consistent images provides extra depth information than image contexts [36]. However, for the dynamic areas that violate the static scene assumption, the cost volume provides wrong depth probabilities, which mislead the network to produce wrong depth even under supervised learning [37]. In this regard, processing the dynamic areas has become one of the main challenges for multi-frame depth estimation in outdoor scenes.

**Dynamic depth estimation in outdoor scenes.** Dynamic areas are ubiquitous in real-world scenarios and are important in applications such as autonomous driving. Some literature [4, 5, 14, 22, 23] seeks to handle dynamic objects in self-supervised monocular depth estimation, where the network inputs a single image and the dynamic areas mainly affect the supervisory signals. They typically identify possible moving objects by semantic [21] or instance segmentation [3–5, 22], then conduct robust learning by masking out dynamic areas during loss computation [21] or directly model the object motion [3–5, 14]. However, these methods differ from the topic discussed in this paper in that: 1) we focus on addressing multi-frame dynamic depth estimation issue, where the main challenge lies in the corrupted cost

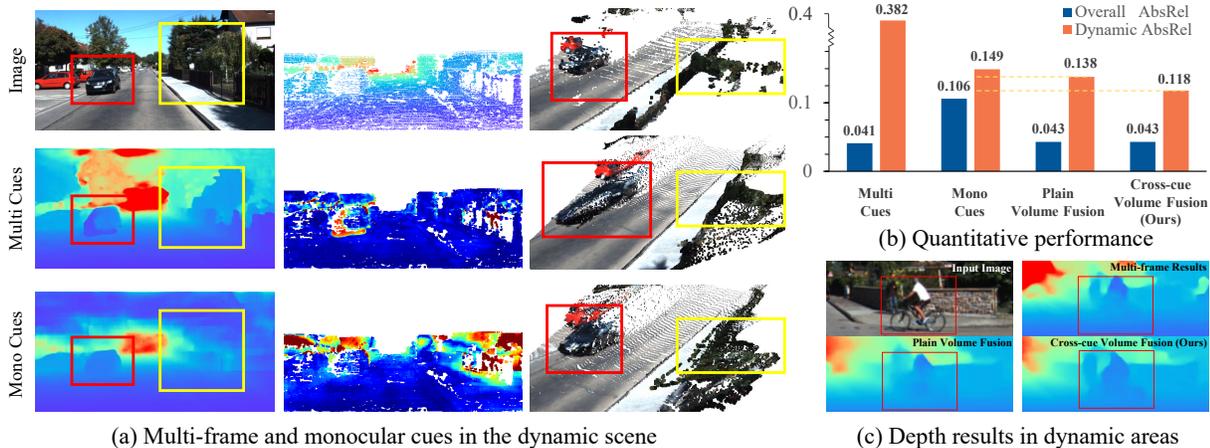


Figure 2. **Multi-frame and monocular cues in the dynamic scene.** a) Multi-frame cues preserve accurate geometric properties in the static area (yellow box), while the monocular cues learn good structural relations in dynamic areas (red box). b) Multi-frame and monocular cues show respective benefits in different areas. While our plain volume fusion scheme shows obvious performance improvement, the proposed *cross-cue fusion* demonstrates better capabilities to handle dynamic depth. c) Depth predictions that show progressive improvements from the plain volume fusion and cross-cue volume fusion.

volumes than the self-supervised loss; 2) our method is supervised so that the robustness of loss function is beyond the main concern of this task.

In the context of multi-frame depth estimation, current methods leverage the depth information from the single image to improve depth in dynamic areas. Manydepth [36] proposes a self-discovered mask and supervises the potential dynamic areas with monocular depth. MonoRec [37] proposes a motion segmentation network to mask out the dynamic areas in the cost volume and use only monocular image features to infer depth. Feng *et al.* [9] proposes to correct the dynamic object locations (with instance mask) in the image plane using monocular depth before computing the cost volume. Despite the higher dynamic results than pure multi-frame estimations, their performances are quite comparable [1,9,36], if not worse than, their proposed monocular branch (as shown in Tab. 4). Moreover, they require pre-computed instance masks [9,37] that bring extra computation burden for network training or inference. There also exists another multi-frame method [1] which guides multi-frame cost reconstruction with a large monocular network. However, due to the reliance on monocular network predictions, it demonstrates weaker generalization capabilities than multi-frame methods (Tab. 3). In contrast, our methods do not require any pre-computed object masks and achieve obvious improvement upon both monocular and multi-frame predictions in dynamic areas while retaining good generalization abilities across datasets.

### 3. Analyses on Multi-view & Monocular Cues

Focusing on the dynamic scenes, we first analyze the behaviors of the depth estimation methods relying on multi-view and monocular cues, as in Sec. 3.1. Beyond the overall performance on the whole scene, we specifically study their behaviors in the dynamic area. Then we analyze how

the proper fusion of multi-view and monocular cues may benefit depth estimation in dynamic areas in Sec. 3.2.

#### 3.1. Behaviours of Multi-view and Monocular Cues

We implement two depth estimation methods with the two types of cues, *i.e.*, Multi Cues and Mono Cues in Fig. 2, and train them on KITTI [10], where the U-Net [37] is used. Specifically, the model with multi-view cues (Multi Cues) takes the cost volume as the input. To analyze the performance in the dynamic areas, we use the dynamic mask from [37] computed by thresholding photometric error and depth inconsistency.

As shown in Fig. 2 (b), multi-frame depth estimation generally achieves high overall accuracy relying on the multi-view clues. But the performance in the dynamic areas is degraded due to the violation of the multi-view consistency in the input cost volume. Fig. 2 (a) shows that the multi-frame method generates fine 3D structures in static areas (the yellow box), whereas its estimation of the moving car (the red box) is corrupted. Compared with the multi-frame method, the monocular depth estimation method achieves worse overall performance as shown in Fig. 2 (b), without the geometric cues. On the other hand, it does not suffer from multi-view inconsistency and thus performs better in dynamic areas. In Fig. 2 (a), we can observe that the two depth cues show respective benefits in different areas and different aspects of the dynamic scene.

#### 3.2. Potential Benefits of Fusion

To handle the dynamic areas violating the multi-view consistency, previous multi-frame depth estimation methods [9,36,37] try to use the monocular information as the input for the dynamic areas. However, the improvement is limited due to the underutilized benefits of both cues, as discussed in Sec. 1. Especially the performance in the dynamic

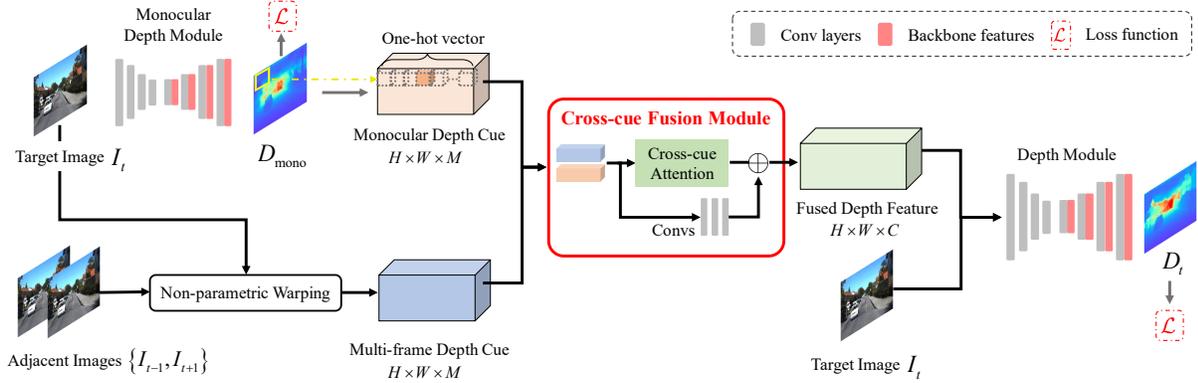


Figure 3. **Overview of the proposed method.** We first extract multi-frame depth cues with cost volume and monocular depth cues using one-hot depth volume. Then, we fuse the two volumes with the proposed cross-cue fusion module (CCF) to yield an improved fused depth feature. The fused depth feature is sent to the depth module for final depth estimation.

areas is easily bounded by the monocular cues.

We seek to explore better fusion schemes to use the complementary benefits of these two cues, as discussed in Sec. 3.1. Apart from using the monocular cues in the dynamic area, we also hope the geometric information learned with the multi-view cues in the static areas can be propagated to boost the monocular cues in dynamic areas, which requires a more comprehensive fusion process.

To unify the information from both multi-frame and monocular cues, we encode both depth cues in the form of a cost volume [40] and a one-hot depth volume transferred from the estimated monocular depth map (in Sec. 4.2). In the analyses, we first consider the *plain volume fusion* using concatenation and convolution operations. This learnable fusion scheme brings obvious improvement in dynamic areas, implying proper fusion methods help to utilize the potential benefits in dynamic scenes. Based on the feasibility of learnable fusion, we further investigate the *cross-cue volume fusion* and it exhibits further improvement (in Fig. 2 (b), (c)), which indicates a better way to utilize depth cues as introduced in the next section.

## 4. The Proposed Method

We aim to learn depth  $D_t$  of the target image  $I_t$  from a short image sequence, with known or estimated camera parameters  $K, T$ . In this paper, we define the image sequence as  $\{I_{t-1}, I_t, I_{t+1}\}$ , where  $\{I_{t-1}, I_{t+1}\}$  are adjacent images to the target frame  $I_t$ ,  $K, T$  are camera intrinsic and extrinsic provided as known values. The goal is to conduct accurate estimations of the dynamic scenes that contain various challenging dynamic objects.

### 4.1. Overview

As shown in Fig. 3, the proposed method consists of three major parts - the multi-view and monocular cues are first represented as volumes through cost volume construction and depth one-hot vector transformation, the cross-cue fusion (CCF) module then fuses both multi-frame

and monocular cues by leveraging attention mechanisms [13, 31, 38], *i.e.*, extracting relative intra-relations of each cue to guide the other to yield improved geometric representations of dynamic scene structure. The depth module takes the fused representation to estimate the final depth.

### 4.2. Representing Monocular and Multi-view Cues

**Multi-view cues as cost volume.** We represent the multi-view cues by computing the cost volume following the pipeline of multi-view stereo (MVS) [37, 40]. We warp the adjacent images  $\{I_{t-1}, I_{t+1}\}$  to the target view using  $K, T$  and a set of depth hypotheses  $d \in \{d_k\}_{k=1}^M$  uniformly sampled in the inverse depth space  $[\frac{1}{d_{\min}}, \frac{1}{d_{\max}}]$ , where  $M$  denotes the number of depth hypotheses and is set to 32 in our paper. We construct the multi-frame cost volume  $C_{\text{multi}} \in \mathbb{R}^{H \times W \times M}$  by measuring the pixel-wise similarity between the warped images and the target image, using SSIM [35] as introduced in [37]. For each pixel  $(i, j)$  of  $C_{\text{multi}} \in [0, 1]^{H \times W \times M}$ , channel positions  $k \in \{1, \dots, M\}$  with large matching scores indicates a higher possibility to include real scene depth.

**Monocular cues as depth volume.** We construct the monocular cues by estimating the single-view depth map and then transform the depth into one-hot depth volume. We leverage a U-Net architecture [37] without any complex designs for single-view depth estimation, yielding monocular prediction  $D_{\text{mono}} = f_{\theta}^{\text{mono}}(I_t)$ , where  $D_{\text{mono}} \in \mathbb{R}^{H \times W}$  and  $f_{\theta}^{\text{mono}}$  is the monocular network. To facilitate smooth fusion between the multi-frame and monocular cues, different from methods using single-view features [34, 37] or depth values [9, 36], we transform the whole depth map  $D_{\text{mono}}$  into the depth volume  $C_{\text{mono}} \in \{0, 1\}^{H \times W \times M}$  by converting each absolute depth value to a one-hot vector with

$$C_{\text{mono},(i,j)}[k] = \{1 \mid d_{\text{mono}} \in (d_{k-1}, d_k]\}_{k=1}^M, \quad (1)$$

where  $C_{\text{mono},(i,j)} \in \{0, 1\}^M$  represents the one-hot vector in  $C_{\text{mono}}$  corresponding to the pixel at  $(i, j)$ ,  $d_{\text{mono}}$  is the

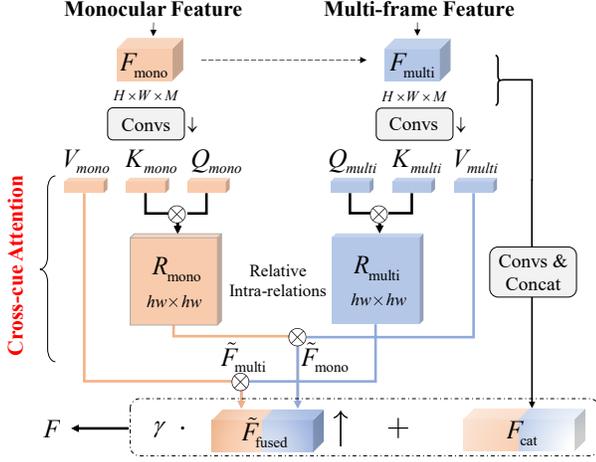


Figure 4. **The cross-cue fusion (CCF) module.** Taking both multi-frame and monocular depth volume as input, the CCF module enhances multi-frame depth features with the relative intra-relations of monocular depth volume ( $R_{mono}$ ). Meanwhile, the intra-relations of the multi-frame depth volume ( $R_{multi}$ ) also enhance the monocular module, yielding enhanced depth features for the final depth estimation.

pixel-wise depth value. Note that we also try some soft representations [1, 2] but observe no obvious improvement.

### 4.3. The Cross-cue Fusion Module

We propose a cross-cue fusion (CCF) module to fuse the multi-frame and monocular depth cues. Different from methods [9, 36, 37] that utilize the masked local depth cues for dynamic depth estimation, the CCF module fuses the whole depth cues using depth volumes, retaining the potential to leverage both benefits for further improvement.

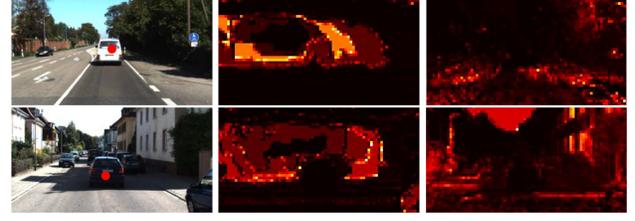
**The Cross-cue fusion pipeline.** Given volumes  $C_{multi}, C_{mono} \in \mathbb{R}^{H \times W \times M}$ , we first process them via shallow convolution layers, yielding down-sampled monocular and multi-frame depth features  $F_{multi}, F_{mono}$  in shape  $\mathbb{R}^{h \times w \times M}$ . We then feed  $F_{multi}$  and  $F_{mono}$  to the proposed *cross-cue attention* (CCA) to enhance each depth cue by extracting the relative intra-relations from the other

$$\begin{aligned} \tilde{F}_{multi} &= \text{CCA}_{multi}(F_{mono}, F_{multi}), \\ \tilde{F}_{mono} &= \text{CCA}_{mono}(F_{multi}, F_{mono}), \end{aligned} \quad (2)$$

the enhanced features are concatenated to yield the fused feature  $\tilde{F}_{fused}$ . To retain detailed information from initial depth cues, we process the input depth cues via  $F_{cat} = \text{Cat}(\text{Conv}(C_{multi}), \text{Conv}(C_{mono}))$  and add the residual connection. The final cross-cue features can be written as

$$F = \gamma \tilde{F}_{fused} \uparrow + F_{cat}, \quad (3)$$

where  $\gamma$  is a learned weighting factor and ‘ $\uparrow$ ’ denotes the up-sampling operation. The fused feature  $F$  is then sent to the depth network along with the image context features to yield the final depth prediction  $D_t \in \mathbb{R}^{H \times W}$ .



(a) Input image (b)  $R_{mono}$  atten. map (c)  $R_{multi}$  atten. map

Figure 5. **Attention maps of the dynamic position (red dots) examples.**  $R_{mono}$  boosts multi-frame features by attending monocular features around the dynamic areas for better use of the monocular cues in the corresponding areas.  $R_{multi}$  enhances monocular features by attending multi-frame features in the static areas which have more reliable geometric information.

**Cross-cue attention.** The cross-cue attention (CCA) targets utilizing the relative intra-relation of one depth cue to improve the geometric information of another. Since the CCA modules are deployed in a parallel manner as shown in Fig. 4, we introduce  $\tilde{F}_{multi} = \text{CCA}_{multi}(F_{mono}, F_{multi})$  in detail as an example.

Given depth features  $F_{mono}, F_{multi} \in \mathbb{R}^{h \times w \times M}$ , we transform  $F_{mono}$  into query feature  $Q_{mono}$  and key feature  $K_{mono}$ , then transform  $F_{multi}$  into value feature  $V_{multi}$  using convolution operation  $f(\cdot, \theta)$

$$\begin{aligned} Q_{mono} &= f(F_{mono}, \theta_{mono}^Q), \\ K_{mono} &= f(F_{mono}, \theta_{mono}^K), \\ V_{multi} &= f(F_{multi}, \theta_{multi}^V), \end{aligned} \quad (4)$$

after reshaping all features into size  $(hw, M)$ , we compute the non-local relative intra-relations of monocular features by matrix multiplication  $\otimes$  followed by softmax operation

$$R_{mono} = \text{Softmax}(Q_{mono} \otimes K_{mono}^T), \quad (5)$$

where  $R_{mono} \in \mathbb{R}^{hw \times hw}$  stands for the non-local intra-relations of the monocular depth cue. We then utilize  $R_{mono}$  to improve the geometric representations of the multi-frame feature  $V_{multi}$  by

$$\tilde{F}_{multi} = R_{mono} \otimes V_{multi}, \quad (6)$$

where  $\tilde{F}_{multi}$  denotes the improved multi-frame representations benefited from monocular depth cues. Similar operations are done for  $\tilde{F}_{mono} = \text{CCA}_{mono}(F_{multi}, F_{mono})$ , where  $R_{multi}$  stands for the intra-relations of multi-frame cues that can be used to improve monocular depth feature  $V_{mono}$ . Please refer to Fig. 4 for details.

**Effectiveness of the CCA.** As shown in Fig. 5, despite using the same CCA operation, the intra-relations  $R_{mono}$  and  $R_{multi}$  attend different areas for improving the dynamic depth, showing their ability to capture respective benefits from monocular (better dynamic depth) and multi-frame cues (better static depth). This learned property enables

Eval	Method	Back.	Reso.	Sup.	Abs Rel	Sq Rel	RMSE	RMSE <sub>log</sub>	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Overall	Manydepth [36]	Res-18	MR	M	0.071	0.343	3.184	0.108	0.945	0.991	0.998
	DynamicDepth [9]	Res-18	MR	M	0.068	0.296	3.067	0.106	0.945	0.991	0.998
	MonoRec [37]	Res-18	MR	D*	0.050	0.290	2.266	0.082	0.972	0.991	0.996
	<b>Ours</b>	Res-18	MR	D	<b>0.043</b>	<b>0.151</b>	<b>2.113</b>	<b>0.073</b>	<b>0.975</b>	<b>0.996</b>	<b>0.999</b>
	MaGNet [1]	Effi-B5	MR	D	0.057	0.215	2.597	0.088	0.967	<b>0.996</b>	<b>0.999</b>
	<b>Ours</b>	Effi-B5	MR	D	0.046	0.155	2.112	0.076	0.973	<b>0.996</b>	<b>0.999</b>
	MaGNet [1]	Effi-B5	HR	D	0.043	0.135	2.047	0.082	0.981	<b>0.997</b>	<b>0.999</b>
	<b>Ours</b>	Effi-B5	HR	D	<b>0.039</b>	<b>0.103</b>	<b>1.718</b>	<b>0.067</b>	<b>0.981</b>	<b>0.997</b>	<b>0.999</b>
Dynamic	Manydepth [36]	Res-18	MR	M	0.222	3.390	7.921	0.237	0.676	0.902	0.964
	DynamicDepth [9]	Res-18	MR	M	0.208	2.757	7.362	0.227	0.682	0.911	0.971
	MonoRec [37]	Res-18	MR	D*	0.360	9.083	10.963	0.346	0.590	0.882	0.780
	<b>Ours</b>	Res-18	MR	D	0.118	0.835	4.297	0.146	0.871	0.975	0.990
	MaGNet [1]	Effi-B5	MR	D	0.141	1.219	4.877	0.168	0.830	0.955	0.986
	<b>Ours</b>	Effi-B5	MR	D	<b>0.111</b>	<b>0.768</b>	<b>4.117</b>	<b>0.135</b>	<b>0.881</b>	<b>0.980</b>	<b>0.994</b>
	MaGNet [1]	Effi-B5	HR	D	0.140	1.060	4.581	0.202	0.834	0.954	0.982
	<b>Ours</b>	Effi-B5	HR	D	<b>0.112</b>	<b>0.830</b>	<b>4.101</b>	<b>0.137</b>	<b>0.885</b>	<b>0.978</b>	<b>0.992</b>

Table 1. **Quantitative comparisons on KITTI [10] Odometry dataset.** ‘Back.’ denotes the network backbone. ‘Reso.’ denotes the image resolutions, where ‘MR’ refers to the resolution of  $256 \times 512$  and ‘HR’ is  $352 \times 1216$ . In the ‘Sup.’ column, ‘M’ are self-supervised methods, ‘D\*’ refers to semi-supervised methods trained with pseudo GT depth, while ‘D’ denotes fully-supervised methods. Color blue denotes ‘lower is better’, while red means ‘higher is better’. The best results are in bold.

unbounded dynamic depth performance upon both predictions, *i.e.*, the monocular depth can be improved by attending static depth (using  $R_{\text{multi}}$ ) from multi-frame cues, and the improved monocular depth in dynamic areas will further be propagated (with  $R_{\text{mono}}$ ) to multi-frame predictions.

#### 4.4. Loss Function

Given the predicted depth  $D$  and the ground truth depth  $\hat{D}$ , the loss can be described as

$$\mathcal{L}(D, \hat{D}) = \beta \mathcal{L}_{\text{SI}}(D, \hat{D}) + \mathcal{L}_{\text{VNL}}(D, \hat{D}), \quad (7)$$

where  $\mathcal{L}_{\text{SI}}$  denotes the scale-invariant loss [2, 8],  $\mathcal{L}_{\text{VNL}}$  is the virtual normal loss [42, 43],  $\beta$  is the weighing factor which is set to 4. Since we have both monocular depth  $D_{\text{mono}}$  and the final depth prediction  $D_t$ , the final loss  $\mathcal{L}_{\text{final}}$  is

$$\mathcal{L}_{\text{final}} = \mathcal{L}(D_{\text{mono}}, \hat{D}) + \mathcal{L}(D_t, \hat{D}). \quad (8)$$

### 5. Experiments

In this section, we compare our method with state-of-the-art multi-frame depth estimation methods [1, 9, 36, 37] for the dynamic scenes, including self-supervised methods [9, 36], semi-supervised method [37] as well as fully supervised method [6] (Sec. 5.3). We then conduct ablation studies (Sec. 5.4) to evaluate different variants of our method. We also validate the generalization ability (Sec. 5.5) and evaluate different methods’ improvements upon monocular networks for depth estimation in dynamic areas (Sec. 5.6).

#### 5.1. Datasets

**KITTI.** We follow [37] to evaluate dynamic sequential data on KITTI Odometry dataset, which contains 13666

training and 8634 testing samples. We use the dynamic masks provided by [37] for evaluation. More than 1300 samples in the test set contain dynamic objects. We evaluate both medium-resolution (MR:  $256 \times 512$ ) and high-resolution (HR:  $352 \times 1216$ ) depth results, and the metrics are computed within the 0-80m range.

**DDAD.** We use the forward-facing cameras with 3848 testing samples to evaluate the methods’ generalization ability. Since there is no pre-defined dynamic mask for DDAD, we construct it by warping adjacent images with GT depth and selecting instance masks with high photometric error. More than 70% of the test set contains dynamic objects. All methods are evaluated within the 0-80m depth range.

#### 5.2. Implementation Details

We implement our method with Pytorch [29] and train it using NVIDIA TITAN RTX GPUs. Unless specified, both the monocular network and the depth network are consistent with the depth module of [37], with ResNet-18 [18] as backbone using ImageNet [6] pre-trained parameters. We train our model for 80 epochs using the Adam [20] optimizer and learning rate of  $10^{-4}$ , which further drops to  $10^{-5}$  after 65 epochs. The batch size is set to 8.

#### 5.3. Results on KITTI

We show both overall and dynamic performances of different methods in Table 1. Our method achieves the best performance in dynamic depth estimation. Specifically, it outperforms the SOTA fully-supervised MaGNet [1] with 21.28% reduction of Abs.Rel ( $0.141 \rightarrow 0.111$ ), using the same Efficient-B5 [30] backbone. Significant improvements are also observed compared to self/semi-supervised

#	Category	Variant	Dynamic					Overall
			Abs Rel	Abs Sq	RMSE	RMSE <sub>log</sub>	$\delta \leq 1.25$	AbsRel
1	Depth with single cues	Pure multi-frame cues	0.382	7.167	10.292	0.35	0.509	0.041
2		Pure monocular cues	0.149	1.369	5.282	0.178	0.810	0.106
3	Volume fusion with masks	Self-discovered mask [36]	0.130	0.990	4.692	0.160	0.837	0.043
4		MaskNetwork [37]	0.220	2.896	6.299	0.223	0.735	0.040
5	Volume fusion without masks	Stack & 3D Convs	0.154	1.479	5.866	0.189	0.777	0.046
6		Stack & 3D U-Net [15]	0.155	1.444	5.762	0.191	0.772	<b>0.040</b>
7		Concat & 2D Convs	0.138	1.124	5.110	0.174	0.815	0.043
8		<b>Ours</b> CCF w./o. $R_{\text{multi}}$	0.124	0.939	4.610	0.154	0.855	0.043
9		<b>Ours</b> CCF w./o. $R_{\text{mono}}$	0.123	0.926	4.545	0.153	0.861	0.043
10		<b>Ours</b> CCF w./ only intra-cue self-attention	0.122	0.896	4.544	0.152	0.860	0.042
11		<b>Ours</b> CCF w./o. residual connection	0.130	0.961	4.616	0.157	0.840	0.048
12		<b>Ours</b> depth module w./o. $I_t$	0.126	0.954	4.636	0.155	0.844	0.042
13		<b>Ours</b> CCF - full	<b>0.118</b>	<b>0.835</b>	<b>4.297</b>	<b>0.146</b>	<b>0.871</b>	0.043

Table 2. **Ablation experiments on KITTI.** We show the results of different fusion types of multi-frame and monocular cues. ‘CCF’ denotes the proposed cross-cue fusion module. ‘Dynamic’ denotes dynamic depth errors while ‘Overall’ refers to the overall depth error.

methods [9, 36, 37]. Besides the obvious improvement in dynamic areas, our method also outperforms other methods in most metrics for overall performance. Qualitative results are shown in Fig. 6. While the dynamic areas generally lead to performance decline for other multi-frame methods, our method conducts obviously better estimations in the moving objects, while retaining the overall accuracy.

#### 5.4. Ablation Study

As shown in Sec. 3, the way to fuse the two volumes in our method influences the final dynamic depth performance. As shown in Tab. 2, we fuse the multi-frame and monocular volumes in ways that the explicit mask is needed (‘Volume fusion with masks’) and the explicit mask is not needed (‘Volume fusion without masks’). We also evaluate the variants of our method (‘**Ours** CCF’) as individual mask-free methods. Row #1~2 shows individual depth results leveraging pure multi-frame and monocular cues.

**Volume fusion with explicit masks.** We leverage the computed [36] or learned masks [37] to fuse our volumes from two depth cues. As shown in row #3~4, The dynamic mask generated by [37] does not surpass the baseline fusion method. While the self-discovered mask [36] shows certain improvements in dynamic areas, it is computed using heuristics and thus has uncontrolled dynamic mask quality, which leads to more restricted performances than ours.

**Volume fusion without masks.** Besides our method that uses a mask-free fusion scheme, we also fuse the volumes without any mask using convolutions. We first stack the depth volumes in a new dimension and process the fused features with layers of 3D convolutions as well as the 3D U-Net, which is commonly used in the MVS methods [15, 40]. As shown in row # 5~6 of Table 2, the 3D convolutions have no obvious improvement in dynamic areas. However, as we concatenate the two depth volumes and process them with 2D convolution layers (as in Fig. 2), we observe further im-

provements (row #7) upon the pure monocular results, but this mask-free fusion scheme still lags behind our method with obvious margins.

**Design choices of the network modules.** We evaluate different variants of our method as an individual type of mask-free fusion method. We respectively remove the multi-frame (row #8) and monocular (row #9) intra-relations and directly feed the other cues to the fused feature. Then, we implement another variant of our method with intra-cue self-attention in the proposed CCF (row #10). We also disable the residual connection of  $F_{\text{cat}}$  (row #11) and the input target image  $I_t$  of the depth module (row #12). Results show the effectiveness and necessity of the proposed technical designs, which achieves the best performance.

#### 5.5. Generalization on DDAD

To validate the generalization ability of our method, we test the KITTI models of all supervised methods [1, 37] on the challenging DDAD dataset. As shown in Table 3, previous methods either degrade on the dynamic areas due to the violation of multi-view consistency (e.g., AbsRel 0.544 for [37]) or exhibit worse overall performance owing to the cross-domain issue of the monocular network (e.g., AbsRel 0.208 for [1]). Our method shows more promising results that it not only outperforms other methods in dynamic depth estimation but also retains the advantage of multi-frame estimation with competitive overall performance.

#### 5.6. Improvement upon Monocular Estimation

Though previous methods handle dynamic areas with monocular cues, the improvement is usually constrained by monocular results. In Tab. 4, we show the final depth results in dynamic areas and compare them against the monocular results of each method. Our method achieves the largest Abs.Rel reduction of 20.8% compared to other methods which explicitly use the monocular network. Meanwhile,

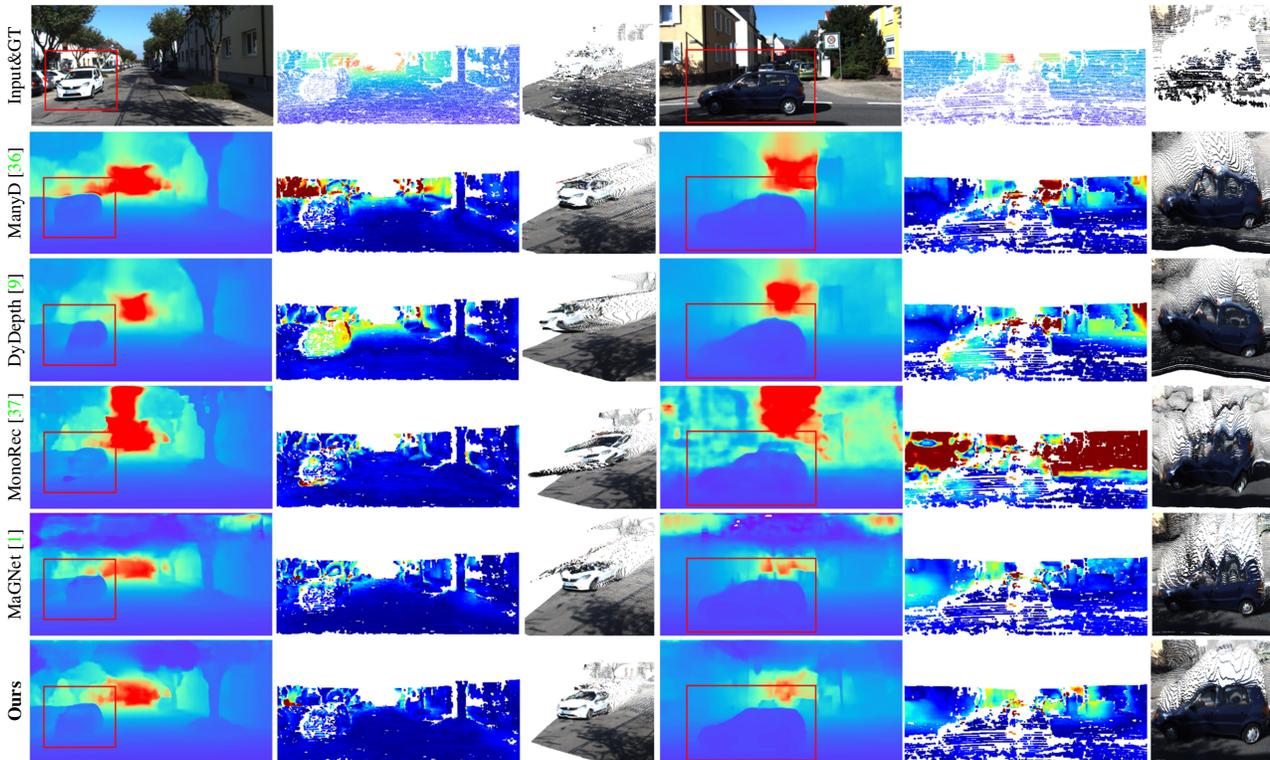


Figure 6. **Qualitative results on KITTI dataset [10]**. From left to right: depth predictions (dynamic objects are highlighted with red boxes), error maps, and the reconstructed point clouds of dynamic areas. Our method achieves the best dynamic results and reconstructs more reasonable object shapes than state-of-the-art methods.

Eval	Method	Backbone	Abs Rel	Sq Rel	RMSE	RMSE <sub>log</sub>	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Overall	MonoRec [37]	Res-18	<b>0.158</b>	3.102	<b>7.553</b>	<b>0.227</b>	<b>0.854</b>	<b>0.931</b>	<b>0.961</b>
	MaGNet [1]	Effi-B5	0.208	<u>2.641</u>	10.739	0.382	0.620	0.878	0.942
	<b>Ours</b>	Res-18	<b>0.158</b>	<b>2.416</b>	<u>9.855</u>	<u>0.299</u>	<u>0.747</u>	<u>0.894</u>	<u>0.947</u>
Dynamic	MonoRec [37]	Res-18	0.544	16.703	16.116	0.482	0.460	0.667	0.798
	MaGNet [1]	Effi-B5	0.266	<u>3.982</u>	<u>11.715</u>	0.398	<u>0.462</u>	<u>0.815</u>	<u>0.917</u>
	<b>Ours</b>	Res-18	<b>0.234</b>	<b>3.611</b>	<b>11.007</b>	<b>0.331</b>	<b>0.576</b>	<b>0.835</b>	<b>0.921</b>

Table 3. **Generalizations on DDAD [16] dataset**. The best results are in **bold** and the second best results are underlined. Our method achieves a competitive overall performance with other methods while achieving the best result in dynamic areas.

Method	Mono. Err.	Final Err.	Err. Redu.
Manydepth [36]	0.212	0.222	-4.72%
Dynamicdepth [9]	0.214	0.208	2.83%
MaGNet [1]	0.153	0.141	7.84%
<b>Ours - Res.18</b>	0.149	0.118	<b>20.81%</b>
<b>Ours - Res.50</b>	0.145	0.116	<b>20.00%</b>

Table 4. **Error reduction upon monocular estimation in dynamic areas**. We show different methods’ *dynamic* Abs.Rel depth errors (Final Err.) and their reduction (Err. Redu.) compared to their individual monocular performances (Mono. Err.). Our method achieves the largest error reduction over others and exhibits a consistent error reduction when a better backbone is used.

when using a backbone with larger capacities (ResNet-50), our method achieves consistent error reduction upon the monocular network, showing its scalability and flexibility.

## 6. Conclusion

We improve multi-frame dynamic depth estimation by fusing the multi-view and monocular depth cues with the proposed cross-cue fusion. Experiments show its effectiveness as well as the generalization ability. **Limitation:** there still exists a performance gap to fill, especially in dynamic areas, with a better fusion of the cues. Meanwhile, the challenging occlusion issues for general multi-frame methods require further exploration.

**Acknowledgements** Y. Zhu, J. Sun, and Y. Zhang acknowledge the support from National Engineering Laboratory for Integrated Aero-Space-Ground-Ocean Big Data Application Technology. This work was partially supported by NSFC (No.U19B2037) and the Natural Science Basic Research Program of Shaanxi (No.2021JCW-03) to Y. Zhang; an ARC DECRA Fellowship DE230101591 to D. Gong.

## References

- [1] Gwangbin Bae, Ignas Budvytis, and Roberto Cipolla. Multi-view depth estimation by fusing single-view depth probability with multi-view geometry. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2842–2851, 2022. 2, 3, 5, 6, 7, 8
- [2] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4009–4018, 2021. 1, 5, 6
- [3] Fabian Brickwedde, Steffen Abraham, and Rudolf Mester. Mono-sf: Multi-view geometry meets single-view depth for monocular scene flow estimation of dynamic traffic scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2780–2790, 2019. 2
- [4] Zhe Cao, Abhishek Kar, Christian Hane, and Jitendra Malik. Learning independent object motion from unlabelled stereoscopic videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5594–5603, 2019. 2
- [5] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8001–8008, 2019. 2
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 6
- [7] Arda Duzceker, Silvano Galliani, Christoph Vogel, Pablo Speciale, Mihai Dusmanu, and Marc Pollefeys. Deepvideomvs: Multi-view stereo on video with recurrent spatio-temporal fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15324–15333, 2021. 2
- [8] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. 6
- [9] Ziyue Feng, Liang Yang, Longlong Jing, Haiyan Wang, YingLi Tian, and Bing Li. Disentangling object motion and occlusion for unsupervised multi-frame monocular depth. *arXiv preprint arXiv:2203.15174*, 2022. 1, 2, 3, 4, 5, 6, 7, 8
- [10] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012. 1, 3, 6, 8
- [11] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3828–3838, 2019. 2
- [12] Dong Gong, Jie Yang, Lingqiao Liu, Yanning Zhang, Ian Reid, Chunhua Shen, Anton Van Den Hengel, and Qinfeng Shi. From motion blur to motion flow: A deep learning solution for removing heterogeneous motion blur. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2319–2328, 2017. 1
- [13] Dong Gong, Frederic Z Zhang, Javen Qinfeng Shi, and Anton Van Den Hengel. Memory-augmented dynamic neural relational inference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11843–11852, 2021. 4
- [14] Ariel Gordon, Hanhan Li, Rico Jonschkowski, and Anelia Angelova. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8977–8986, 2019. 2
- [15] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuoqiuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2495–2504, 2020. 1, 2, 7
- [16] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Rantos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2485–2494, 2020. 1, 8
- [17] Vitor Guizilini, Rares Ambrus, Dian Chen, Sergey Zakharov, and Adrien Gaidon. Multi-frame self-supervised depth with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 160–170, 2022. 2
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 6
- [19] Matthias Innmann, Kihwan Kim, Jinwei Gu, Matthias Nießner, Charles Loop, Marc Stamminger, and Jan Kautz. Nrmvs: Non-rigid multi-view stereo. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2754–2763, 2020. 1
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [21] Marvin Klingner, Jan-Aike Termöhlen, Jonas Mikolajczyk, and Tim Fingscheidt. Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance. In *European Conference on Computer Vision*, pages 582–600. Springer, 2020. 2
- [22] Seokju Lee, Sunghoon Im, Stephen Lin, and In So Kweon. Learning monocular depth in dynamic scenes via instance-aware projection consistency. *arXiv preprint arXiv:2102.02629*, 2021. 2
- [23] Hanhan Li, Ariel Gordon, Hang Zhao, Vincent Casser, and Anelia Angelova. Unsupervised monocular depth learning in dynamic scenes. *arXiv preprint arXiv:2010.16404*, 2020. 2
- [24] Rui Li, Xiantuo He, Yu Zhu, Xianjun Li, Jinqiu Sun, and Yanning Zhang. Enhancing self-supervised monocular depth estimation via incorporating robust constraints. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 3108–3117, 2020. 1

- [25] Rui Li, Danna Xue, Shaolin Su, Xiantuo He, Qing Mao, Yu Zhu, Jinqiu Sun, and Yanning Zhang. Learning depth via leveraging semantics: Self-supervised monocular depth estimation with both implicit and explicit semantic guidance. *Pattern Recognition*, page 109297, 2023. [1](#)
- [26] Rui Li, Danna Xue, Yu Zhu, Hao Wu, Jinqiu Sun, and Yanning Zhang. Self-supervised monocular depth estimation with frequency-based recurrent refinement. *IEEE Transactions on Multimedia*, 2022. [1](#)
- [27] Yin hao Li, Zheng Ge, Guanyi Yu, Jinrong Yang, Zengran Wang, Yukang Shi, Jianjian Sun, and Zeming Li. Bevdepth: Acquisition of reliable depth for multi-view 3d object detection. *arXiv preprint arXiv:2206.10092*, 2022. [1](#)
- [28] Zhenxing Mi, Chang Di, and Dan Xu. Generalized binary search network for highly-efficient multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12991–13000, 2022. [2](#)
- [29] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. [6](#)
- [30] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. [6](#)
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. [4](#)
- [32] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, and Marc Pollefeys. Itermv: Iterative probability estimation for efficient multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8606–8615, 2022. [2](#)
- [33] Tai Wang, Jiangmiao Pang, and Dahua Lin. Monocular 3d object detection with depth from motion. *arXiv preprint arXiv:2207.12988*, 2022. [1](#)
- [34] Xiaofeng Wang, Zheng Zhu, Guan Huang, Fangbo Qin, Yun Ye, Yijia He, Xu Chi, and Xingang Wang. Mvster: epipolar transformer for efficient multi-view stereo. In *European Conference on Computer Vision*, pages 573–591. Springer, 2022. [1](#), [4](#)
- [35] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. [4](#)
- [36] Jamie Watson, Oisín Mac Aodha, Victor Prisacariu, Gabriel Brostow, and Michael Firman. The temporal opportunist: Self-supervised multi-frame monocular depth. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1164–1174, 2021. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [37] Felix Wimbauer, Nan Yang, Lukas Von Stumberg, Niclas Zeller, and Daniel Cremers. Monorec: Semi-supervised dense reconstruction in dynamic environments from a single moving camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6112–6122, 2021. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [38] Qingsen Yan, Dong Gong, Javen Qinfeng Shi, Anton van den Hengel, Chunhua Shen, Ian Reid, and Yanning Zhang. Dual-attention-guided network for ghost-free high dynamic range imaging. *International Journal of Computer Vision*, pages 1–19, 2022. [4](#)
- [39] Zhenpei Yang, Zhile Ren, Qi Shan, and Qixing Huang. Mvs2d: Efficient multi-view stereo via attention-driven 2d convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8574–8584, 2022. [2](#)
- [40] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. [1](#), [2](#), [4](#), [7](#)
- [41] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5525–5534, 2019. [2](#)
- [42] Wei Yin, Yifan Liu, and Chunhua Shen. Virtual normal: Enforcing geometric constraints for accurate and robust depth prediction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. [1](#), [6](#)
- [43] 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*, pages 5684–5693, 2019. [6](#)
- [44] Wei Yin, Jianming Zhang, Oliver Wang, Simon Niklaus, Simon Chen, Yifan Liu, and Chunhua Shen. Towards accurate reconstruction of 3d scene shape from a single monocular image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. [1](#)
- [45] Wei Yin, Jianming Zhang, Oliver Wang, Simon Niklaus, Long Mai, Simon Chen, and Chunhua Shen. Learning to recover 3d scene shape from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 204–213, 2021. [1](#)
- [46] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. New crfs: Neural window fully-connected crfs for monocular depth estimation. *arXiv preprint arXiv:2203.01502*, 2022. [1](#)