

# M<sup>2</sup>Depth: Self-supervised Two-Frame Multi-camera Metric Depth Estimation

Yingshuang Zou<sup>1,2\*</sup> Yikang Ding<sup>2\*†</sup> Xi Qiu<sup>2</sup>  
Haoqian Wang<sup>1‡</sup> Haotian Zhang<sup>2‡</sup>

<sup>1</sup> Tsinghua University <sup>2</sup> Megvii Technology

**Abstract.** This paper presents a novel self-supervised two-frame multi-camera metric depth estimation network, termed M<sup>2</sup>Depth, which is designed to predict reliable scale-aware surrounding depth in autonomous driving. Unlike the previous works that use multi-view images from a single time-step or multiple time-step images from a single camera, M<sup>2</sup>Depth takes temporally adjacent two-frame images from multiple cameras as inputs and produces high-quality surrounding depth. We first construct cost volumes in spatial and temporal domains individually and propose a spatial-temporal fusion module that integrates the spatial-temporal information to yield a strong volume presentation. We additionally combine the neural prior from SAM features with internal features to reduce the ambiguity between foreground and background and strengthen the depth edges. Extensive experimental results on nuScenes and DDAD benchmarks show M<sup>2</sup>Depth achieves state-of-the-art performance. More results can be found in [project page](#).

**Keywords:** Depth Estimation · Surrounding Depth · Self-supervised Learning

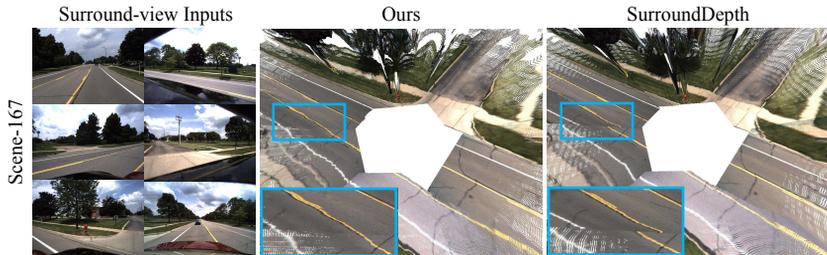
## 1 Introduction

Depth estimation aims to recover the 3D structure of the real world from 2D images, playing a fundamental role in various applications. In recent years, with the development of autonomous driving, using depth estimation methods to get the 3D representation of the driving scenes shows tremendous attraction, as replacing the expensive depth sensor (*e.g.* Lidar) with vehicle-mounted cameras is cost-effective.

Many previous works [12, 15, 37, 44] focus on estimating depth from a single RGB image. Though flexible and concise, such methods suffer from obtaining consistent scale-aware depth (*i.e.* metric depth) among multi-frame and multi-camera when applied in driving scenes. In order to simultaneously predict the surrounding depth, recent methods [16, 38] feed multiple images from 360° vehicle-mounted cameras into 2D encoder-decoder network to capture the spatial information between surround cameras. However, these methods use only

---

\* Equal Contribution. †Project Leader. ‡Co-corresponding Authors.



**Fig. 1:** Point clouds comparison on DDAD [14] dataset. By transforming the predicted depth into point clouds, we show that our method achieves more consistent and accurate estimation compared with SurroundDepth [38]. The visualized point clouds are fused using the surrounding depth of one frame, and the blue boxes highlight the challenging area spanning multiple cameras.

one frame and ignore the temporal information, still facing the challenge of predicting consistent metric depth. Some existing methods [23, 34] show that taking temporally adjacent frames as inputs could help get reliable depth under the single camera setting. Nevertheless, few works explore and leverage the spatial-temporal information to strengthen the surrounding depth estimation in driving scenes.

In this paper, we propose a novel self-supervised Two-frame *Multi-camera Metric depth* estimation network, named  $M^2$ Depth, to predict consistent scale-aware surrounding depth. The key insight of  $M^2$ Depth is that we believe combining the spatial-temporal info could boost the surrounding depth estimation, as the spatial info provides an important world scale (from calibrated extrinsic between adjacent cameras) and the temporal info benefits the depth consistency. As shown in Fig. 1,  $M^2$ Depth is able to recover the 3D point clouds with coherence between multiple cameras, while the existing method struggles with keeping multi-camera consistency. Specifically,  $M^2$ Depth determines depth by constructing 3D cost volumes within the spatial-temporal domain and applying constraints across multiple cameras, which is different from existing methods [12, 16, 21, 37, 38]. Following the classical plane-sweeping algorithm [7], we construct temporal volumes by utilizing temporal adjacent frames, while the spatial volumes are built by leveraging each view and its overlapped spatial adjacent views. Building accurate cost volumes faces several challenges. First, getting reliable relative pose and depth annotations is difficult, thus we design a pose estimation branch to predict the relative vehicle pose between two frames and train  $M^2$ Depth in a self-supervised manner. Second, the depth range in driving scenes is typically large, we consequently design a mono prior branch to estimate coarse depth to narrow down the depth search range. The initial 3D cost volumes are constructed in the spatial domain and temporal domain separately, which consist of the co-visibility information on the spatially and temporally adjacent views. To jointly use the spatial-temporal clues, we propose a novel spatial-temporal fusion (STF) module, which fuses the initial volumes with visibility-aware weights. As a result, the fused volumes integrate the space-

time correlation between multiple frames and multiple cameras, which will be then decoded to produce the final depth.

Additionally, we observe that the feature learning of M<sup>2</sup>Depth is unstable as it actually learns to simultaneously estimate the relative pose, monocular prior, and the multi-camera depth under weak supervision. Specifically, we find that the depth estimation method of constructing spatial-temporal volume through pixel matching lacks consistency within instances and discrimination between instances for features. Due to poor features that lack the discrimination between instances, the quality of the depth map decreases. Inspired by the Segment Anything Model (*a.k.a.* SAM) [22], we propose to inject the strong neural prior from pretrained SAM features into depth estimation to strengthen the feature learning. The key insight is that SAM is able to capture fine-grained inter-view and intra-view semantic information, which is critical for surrounding depth estimation. We thus design a multi-grained feature fusion (MFF) module to integrate SAM features. To the best of our knowledge, M<sup>2</sup>Depth is the first to use the SAM feature in a depth estimation task.

We train and validate M<sup>2</sup>Depth on two large-scale multi-camera depth estimation benchmarks, *i.e.* DDAD [14] and nuScenes [5], and the extensive experimental results demonstrate M<sup>2</sup>Depth achieves state-of-the-art performance in multi-camera metric depth estimation task.

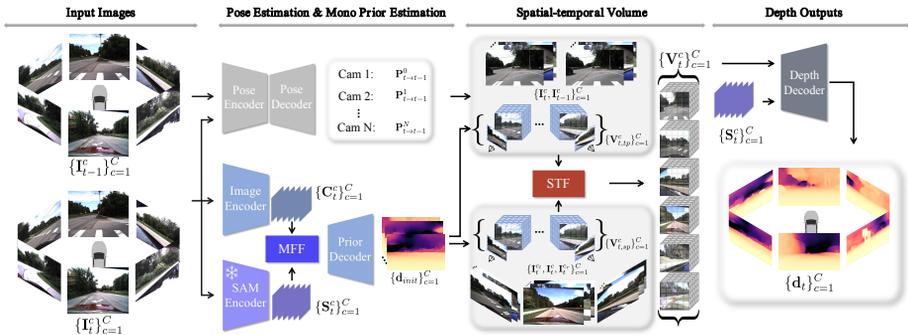
In summary, our main contributions are as follows:

- We present M<sup>2</sup>Depth, a novel self-supervised two-frame multi-camera metric depth estimation network, which achieves state-of-the-art performance on multiple surrounding depth estimation benchmarks.
- For the first time, we propose to construct spatial-temporal 3D cost volumes and design a spatial-temporal fusion (STF) module for surrounding depth estimation, which strengthens the depth accuracy by fusing the spatial-temporal information.
- We introduce the strong SAM prior into the depth estimation task and propose a multi-grained feature fusion (MFF) module to integrate SAM features with internal features for enhancing the depth quality in detail.

## 2 Related Works

### 2.1 Multi-frame Depth Estimation

Unlike the monocular depth estimation (MDE) works [2, 12, 24, 47] that solely use single frame to predict depth, multi-frame depth estimation works [1, 3, 15, 37, 39, 41] take as inputs the adjacent frames to enhance depth quality, which yields great improvements in practical applications. MonoRec [39] constructs a cost volume based on multiple frames from a single camera to estimate depth. It additionally requires a visual odometry system [41] to provide inter-frame pose and sparse depth as supervision signals. Manydepth [37] learns adaptive cost volume from input data and proposes consistency loss for moving objects. Some



**Fig. 2:** Overview of  $M^2$ Depth. Given images  $\{I_t^c\}_{c=1}^C$  and  $\{I_{t-1}^c\}_{c=1}^C$  from multiple cameras and two frames,  $M^2$ Depth first estimates the pose of the front camera  $P_{t \rightarrow t-1}^0$ , which will be used to infer the poses of all other cameras  $\{P_{t \rightarrow t-1}^c\}_{c=1}^C$ . In mono prior estimation, the multi-grained feature fusion (MFF) module aggregates the internal features  $\{C_t^c\}_{c=1}^C$  from image encoder and the SAM features  $\{S_t^c\}_{c=1}^C$  from SAM encoder to improve feature expression in multi-grained. The depth prior and constraints across multiple cameras are employed to construct 3D cost volumes  $\{V_t^c\}_{c=1}^C$  within the temporal-spatial domain, which will be then used by the spatial-temporal fusion (STF) module to strengthen the accuracy and consistency of cost volumes. Finally, the depth decoder takes as inputs the  $\{V_t^c\}_{c=1}^C$  and  $\{S_t^c\}_{c=1}^C$  to produce the surrounding depth.

methods [10, 23, 34] attempt to estimate a depth prior and then encode it into multi-frame cues.

Despite utilizing multiple frames as input, these methods still struggle to produce reliable surrounding depth when applied in multi-frame multi-camera scenarios. In contrast, our approach imposes spatial constraints via the overlaps among multiple cameras, achieving consistent scale-aware surrounding depth estimation.

## 2.2 Multi-camera Depth Estimation

Intelligent vehicles are typically equipped with multiple surrounding cameras, thus recovering the surrounding 3D environment from mounted cameras turns into a fundamental task in autonomous driving. FSM [30] pioneers the expansion of self-supervised monocular depth estimation to surrounding views by using temporal texture constraints as supervision. SurroundDepth [38] employs the cross-view transformer module to perform feature interaction in surrounding views. EGA-Depth [31] optimizes the computational cost of the attention module, enabling the utilization of higher-resolution feature maps. VFDepth [21] constructs a unified volumetric feature representation to estimate the surrounding depth and canonical vehicle pose.

Unlike the aforementioned methods that use a single frame to predict depth, we construct cost volumes using multi-camera and two frames in the spatial-temporal domain. It is noteworthy that the concurrent work R3D3 [29] uses the

SLAM algorithm [33] and a sequence of frames to estimate vehicle pose and optical flow, which will be then utilized to produce depth. Such a manner inevitably needs many frames as inputs and consumes more computation. On the contrary, our method could use only two frames to achieve comparable performance.

### 2.3 SAM and Applications

Segment Anything Model [22] (SAM) is an effective, promptable transformer-based model for image segmentation, trained on the SA-1B dataset [22], which comprises over 1 billion masks on 11M licensed and privacy respecting images. Its exceptional performance in fine-grained semantic segmentation establishes SAM as a prominent player in numerous tasks, such as tracking [6], image inpainting [43], video text spotting [18], medical image domain [4, 27, 46]. SAMFeat [40] applies SAM to segmentation-independent visual tasks and improve local feature description by using feature-level distillation.

To the best of our knowledge, we are the first to apply SAM in the depth estimation task. By leveraging the SAM feature, we extract fine-grained semantic information from and enhance the accuracy of surrounding depth.

## 3 Methods

### 3.1 Problem Formulation

This paper focuses on the surrounding depth estimation task in autonomous driving, where the cameras are mounted on ego vehicles and provide 360° visual observations. We define that there are  $C$  cameras with known intrinsics  $\{\mathbf{K}^c\}_{c=1}^C$  and extrinsics  $\{\mathbf{T}^c\}_{c=1}^C$ , which associate the cameras with the ego vehicle. Given images of the current frame  $\{\mathbf{I}_t^c\}_{c=1}^C$  and the previous frame  $\{\mathbf{I}_{t-1}^c\}_{c=1}^C$  from surrounding cameras, M<sup>2</sup>Depth produces the scale-aware surrounding depth  $\{\mathbf{d}_t^c\}_{c=1}^C$  at the current timestamp  $t$ . It is noteworthy that the ground truth relative pose of the vehicle between two frames is not required, M<sup>2</sup>Depth is able to predict the relative pose with scale, which is inherently stored in  $\{\mathbf{T}^c\}_{c=1}^C$  and the overlap between adjacent views.

### 3.2 Network Overview

The overall architecture of M<sup>2</sup>Depth is illustrated in Fig. 2. The input images  $\{\mathbf{I}_t^c\}_{c=1}^C$  and  $\{\mathbf{I}_{t-1}^c\}_{c=1}^C$  are first used to perform pose estimation and mono prior estimation. Specifically, the pose encoder takes the images of the front view ( $\mathbf{I}_t^0$  and  $\mathbf{I}_{t-1}^0$ ) as input and learns to predict the 6-Dof relative pose between  $\mathbf{I}_t^0$  and  $\mathbf{I}_{t-1}^0$ . Unlike previous methods that concatenate surrounding views to directly predict the ego pose [38] or construct 4D volumes to estimate optical flow and calculate the ego pose [29], our approach simplifies the ego pose estimation problem by estimating the front camera’s pose  $\mathbf{P}_{t \rightarrow t-1}^0$ , thus the ego pose  $\mathbf{P}_{t \rightarrow t-1}$  can be derived by:

$$\mathbf{P}_{t \rightarrow t-1} = (\mathbf{T}^0)^{-1} \mathbf{P}_{t \rightarrow t-1}^0 \mathbf{T}^0, \quad (1)$$

where  $(\mathbf{T}^0)^{-1}$  indicates the inverse matrix of  $\mathbf{T}^0$ , and  $\mathbf{T}^0$  is the extrinsic matrix between the front camera and ego vehicle.

In monocular prior estimation (Sec. 3.3), the  $\mathbf{I}_t^0$  is fed into a trainable image encoder and a frozen SAM encoder, where the output features are then fused in multi-grained feature fusion (MFF) module. MFF aims to integrate the fine-grained semantic information in SAM features with the depth cues in internal features, which helps M<sup>2</sup>Depth understand the 3D environment. As a result, the prior decoder takes fused features as inputs and produces the mono prior depth  $\{\mathbf{d}_{t\text{prior}}^c\}_{c=1}^C$ , which plays the role of depth guidance in volume construction.

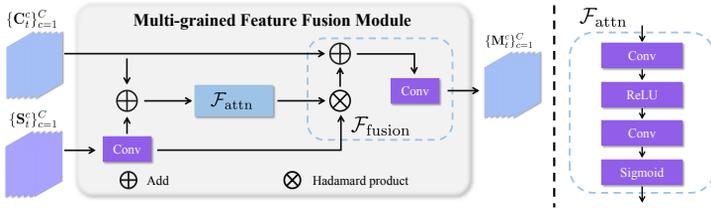
After obtaining the relative vehicle pose  $\mathbf{P}_{t \rightarrow t-1}$  and the mono prior  $\{\mathbf{d}_{t\text{prior}}^c\}_{c=1}^C$ , M<sup>2</sup>Depth constructs the spatial-temporal cost volumes (Sec. 3.4). Specifically, we use the plane-sweeping algorithm [7] to construct the initial cost volume. We first employ a feature pyramid network [25] to extract matching features  $\{\mathbf{F}_t^c\}_{c=1}^C$  from  $\{\mathbf{I}_t^c\}_{c=1}^C$ . As for the temporal domain, we warp the feature  $\mathbf{F}_{t-1}^c$ , which is decoded from  $\mathbf{I}_{t-1}^c$ , to  $\mathbf{F}_t^c$  according to the sampled depth values  $d_t^c$  to get the temporal volume  $\mathbf{V}_{t,tp}^c$ , where  $d_t^c$  is appointed based on  $\mathbf{d}_{t\text{prior}}^c$ . Similarly, the spatial volumes  $\mathbf{V}_{t,sp}^c$  are constructed using spatial adjacent views  $\mathbf{F}_t^{c_l}$  and  $\mathbf{F}_t^{c_r}$ , where the  $c_l$  and  $c_r$  represent the adjacent left and right camera of the reference camera. The initial volumes  $\{\mathbf{V}_{t,tp}^c\}_{c=1}^C$  and  $\{\mathbf{V}_{t,sp}^c\}_{c=1}^C$  will be fed into the spatial-temporal fusion (STF) module, which fuses the spatial-temporal information and yields the final volumes  $\{\mathbf{V}_t^c\}_{c=1}^C$ . Subsequently, we decode the  $\{\mathbf{V}_t^c\}_{c=1}^C$  into depth probability distribution volumes, and produce the estimated depth  $\{\mathbf{d}_t^c\}_{c=1}^C$  by calculating the depth expectation.

### 3.3 Mono Prior Estimation

Following the plane sweep paradigm, the two-frame depth estimation problem can be transformed into a feature matching task [8, 13], where the depth samples would significantly affect final depth quality. Unfortunately, the total depth range in driving scenarios is typically large. As a result, it requires a lot of depth samples to predict precise depth, which would cost tremendous computation in multi-camera settings. To handle this problem, we use a monocular prior estimation branch to produce coarse guidance for cost volume construction.

Specifically, given surrounding views in the current timestamp  $\{\mathbf{I}_t^c\}_{c=1}^C$ , we first use a CNN encoder and a SAM encoder [22] to extract internal features  $\{\mathbf{C}_t^c\}_{c=1}^C$  and SAM priors  $\{\mathbf{S}_t^c\}_{c=1}^C$ , respectively. Then we use the Multi-grained Feature Fusion (MFF) module to fuse  $\{\mathbf{C}_t^c\}_{c=1}^C$  and  $\{\mathbf{S}_t^c\}_{c=1}^C$ , and finally use a CNN decoder to generate mono depth prior  $\{\mathbf{d}_{t\text{prior}}^c\}_{c=1}^C$ .

**Multi-grained Feature Fusion** The detailed structure of the MFF module is illustrated in Fig. 3, which is the key part of depth prior estimation. Given an internal feature  $\mathbf{C}_t^c \in \mathbb{R}^{H \times W \times C}$  and a SAM feature  $\mathbf{S}_t^c \in \mathbb{R}^{H \times W \times C'}$  of a certain camera, we first use convolution layers to align the dimension of  $\mathbf{S}_t^c$  with  $\mathbf{C}_t^c$ , which will be then combined and fed into the  $\mathcal{F}_{\text{attn}}$  to yield attention weights



**Fig. 3:** Details of the multi-grained feature fusion (MFF) module. MFF takes the internal features  $\{C_t^c\}_{c=1}^C$  and the SAM features  $\{S_t^c\}_{c=1}^C$  as inputs, and utilizes a  $\mathcal{F}_{\text{attn}}$  block to yield the weight map, which fetches the complementary info between  $\{C_t^c\}_{c=1}^C$  and  $\{S_t^c\}_{c=1}^C$ , and will be used in  $\mathcal{F}_{\text{fusion}}$  block to produce the fused feature  $\{M_t^c\}_{c=1}^C$ .

$\mathbf{W}_t^c$ :

$$\mathbf{W}_t^c = \mathcal{F}_{\text{attn}}(\mathbf{C}_t^c, \mathbf{S}_t^c) = \sigma(f^{3 \times 3}(\delta(f^{3 \times 3}(\mathbf{C}_t^c + \mathbf{S}_t^c)))), \quad (2)$$

where the  $\sigma$  refers to the *sigmoid* function,  $\delta$  refers to the *ReLU* [28] function, and  $f^{3 \times 3}$  denotes the convolution layer with kernel size of  $3 \times 3$ . Intuitively, the  $\mathbf{W}_t^c$  fetches the complementary info between  $\mathbf{S}_t^c$  and  $\mathbf{C}_t^c$ , and performs as a feature guidance in feature fusion:

$$\mathbf{M}_t^c = \mathcal{F}_{\text{fusion}}(\mathbf{C}_t^c, \mathbf{S}_t^c, \mathbf{W}_t^c) = f^{3 \times 3}(\mathbf{C} + \mathbf{W}_t^c \odot \mathbf{S}_t^c), \quad (3)$$

where  $\mathbf{M}_t^c$  represents the fused features and  $\odot$  denotes the *Hadamard* product. As introducing SAM features into mono prior estimation helps get better performance, we further utilize SAM features in the depth decoding phase in a similar manner to mono prior estimation, more details can be found in supplementary materials, and experimental evaluation is conducted in Sec. 4.

### 3.4 Spatial-temporal Cost Volume

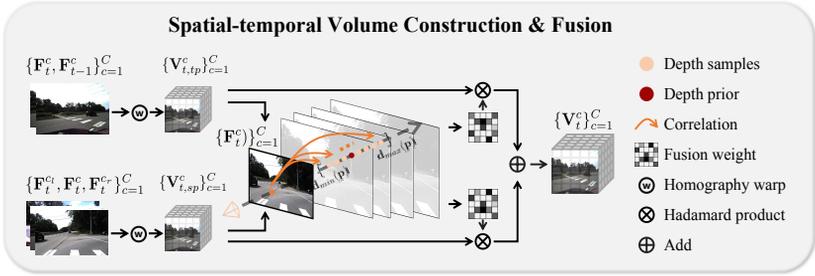
Taking  $\mathbf{I}_t^c \in \mathbb{R}^{H \times W \times 3}$  of a certain camera for example, we denote its spatial adjacent views as  $\mathbf{I}_t^{c_l}, \mathbf{I}_t^{c_r}$  and its last frame as  $\mathbf{I}_{t-1}^c$ . In this section, we use the estimated relative pose  $\mathbf{P}_{t \rightarrow t-1}$  pose and the known camera extrinsics  $\{\mathbf{T}^c\}_{c=1}^C$  to construct the spatial-temporal cost volume  $\{\mathbf{V}_t^c\}_{c=1}^C$ .

**Initial Volume Construction** Taking the aforementioned images as input, we first employ a Feature Pyramid Network (FPN) [25] to extract image features  $\mathbf{F}_t^c \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 3}$ . As for spatial cost volume, we warp  $\mathbf{F}_t^{c_l}$  and  $\mathbf{F}_t^{c_r}$  into the current camera’s frustum according to the camera intrinsics, extrinsics, and the depth samples. The warping operation between a pixel  $\mathbf{p}$  in reference view  $c$  and its corresponding pixel  $\hat{\mathbf{p}}$  in adjacent views  $c' \in \{c_l, c_r\}$  under depth sample  $d$  is defined as:

$$\hat{\mathbf{p}} = \mathbf{K}^{c'} \cdot [\mathbf{R}^{c \rightarrow c'} \cdot ((\mathbf{K}^c)^{-1} \cdot \mathbf{p} \cdot d) + \mathbf{t}^{c \rightarrow c'}], \quad (4)$$

where the transformation matrix  $\mathbf{P}^{c \rightarrow c'} = [\mathbf{R} | \mathbf{t}]^{c \rightarrow c'}$  can be written as:

$$\mathbf{P}^{c \rightarrow c'} = (\mathbf{T}^{c'})^{-1} \cdot \mathbf{T}^c, \quad (5)$$



**Fig. 4:** Overview of the volume construction and STF module. Given the reference image feature  $\mathbf{F}_t^c$  and its temporal-spatial adjacent features, we first warp the adjacent features to reference view to form the initial volumes  $\mathbf{V}_{t,sp}^c$  and  $\mathbf{V}_{t,tp}^c$  in spatial domain and temporal domain respectively. After that, STF fuses the initial volumes by computing the correlation between  $\mathbf{F}_t^c$  and  $\mathbf{V}_{t,sp}^c$ ,  $\mathbf{V}_{t,tp}^c$  and produces the weight maps  $\mathbf{W}_{t,sp}^c$ ,  $\mathbf{W}_{t,tp}^c$ , which will be used as fusion weights to guide the volume fusion.

Afterward, we combine the warped features to form the initial spatial volumes  $\{\mathbf{V}_{t,sp}^c\}_{c=1}^C$ . The temporal volume  $\{\mathbf{V}_{t,tp}^c\}_{c=1}^C$  can be constructed similarly. Specifically, the warping operation in the temporal domain is defined as:

$$\hat{\mathbf{p}} = \mathbf{K}^c \cdot [\mathbf{R}_{t \rightarrow t'}^c \cdot ((\mathbf{K}^c)^{-1} \cdot \mathbf{p} \cdot d) + \mathbf{t}_{t \rightarrow t'}^c], \quad (6)$$

where the  $\hat{\mathbf{p}}$  denotes the corresponding pixel in the previous frame, and the  $\mathbf{P}_{t \rightarrow t'}^c = [\mathbf{R}^c | \mathbf{t}^c]_{t \rightarrow t'}$  is obtained according to Eq. (1):

$$\mathbf{P}_{t \rightarrow t'}^c = (\mathbf{T}^c)^{-1} \cdot \mathbf{P}_{t \rightarrow t}^c \cdot \mathbf{T}^c. \quad (7)$$

During the construction of the initial volumes, we employ  $\{\mathbf{d}_{prior}\}_{c=1}^C$  as the depth guidance and appoint the depth samples in an adaptive range, more details of the depth samples can be found in supplementary materials.

**Spatial-temporal Volume Fusion** By constructing the initial cost volumes, we approach depth estimation as a feature-matching task, where we try to find the best matching feature among the sampled volume features. To strengthen the matching quality using the spatial-temporal information, we propose a spatial-temporal fusion module, which is depicted in Fig. 4.

Given the reference feature  $\mathbf{F}_t^c$  and the initial volumes  $\mathbf{V}_{t,sp}^c$  and  $\mathbf{V}_{t,tp}^c$ , STF produces the fusion weight maps by computing the group-wise correlation [17]. Let  $\mathbf{p}$  denote the pixel in the reference feature, to compute the correlation between  $\mathbf{F}_t^c(\mathbf{p})$  and  $\mathbf{V}_{t,sp}^c(\mathbf{p})$ , we first divide  $C$  feature channels evenly into  $G$  groups,  $\mathbf{F}_t^c(\mathbf{p})^g$  and  $\mathbf{V}_{t,sp}^c(\mathbf{p})^g$ , thus the  $g$ -th group correlation  $\mathbf{Cr}_{t,sp}^c(\mathbf{p})^g$  can be computed as:

$$\mathbf{Cr}_{t,sp}^c(\mathbf{p})^g = \frac{G}{C} \langle \mathbf{F}_t^c(\mathbf{p})^g, \mathbf{V}_{t,sp}^c(\mathbf{p})^g \rangle, \quad (8)$$

where  $\langle a, b \rangle$  indicates the inner product of  $a$  and  $b$ , and the group correlation  $\mathbf{Cr}_{t,tp}^c(\mathbf{p})^g$  between  $\mathbf{F}_t^c(\mathbf{p})^g$  and  $\mathbf{V}_{t,tp}^c(\mathbf{p})^g$  can be obtained with the same manner.

The maximum correlation along group dimension will serve as the final fuse weight  $\mathbf{W}_{t,sp}^c$  and  $\mathbf{W}_{t,tp}^c$ , which will be used to fuse the final spatial-temporal volume  $\mathbf{V}_t^c$ :

$$\mathbf{V}_t^c = \mathbf{W}_{t,sp}^c \mathbf{V}_{t,sp}^c + \mathbf{W}_{t,tp}^c \mathbf{V}_{t,tp}^c. \quad (9)$$

**Depth Prediction** The constructed  $\{\mathbf{V}_t^c\}_{c=1}^C$  are finally input to a depth encoder to produce the final depth  $\{\mathbf{d}_t^c\}_{c=1}^C$ , which also takes the SAM features  $\{\mathbf{S}_t^c\}_{c=1}^C$  as context feature. The key operation is that the depth decoder first transforms the  $\{\mathbf{V}_t^c\}_{c=1}^C$  into probability volumes  $\{\mathbf{P}_t^c\}_{c=1}^C$ , where the  $\mathbf{P}_t^c(\mathbf{p})$  represents the probability distribution among sampled depth of each pixel  $\mathbf{p}$ . The final depth can be obtained by:

$$\mathbf{d}_t^c(\mathbf{p}) = \sum_{i=1}^D d_i \cdot \mathbf{P}_t^c(\mathbf{p}, i), \quad (10)$$

where the  $d_i$  indicates the  $i$ -th sampled depth and the  $\mathbf{P}_t^c(\mathbf{p}, i)$  is the probability of  $\mathbf{p}$  at  $i$ -th depth. Please refer to the supplementary materials for more details about the architecture of the depth decoder.

### 3.5 Loss Function

*Photometric Loss* Following the common practice in self-supervised monocular depth estimation works [12, 31, 38], we optimize M<sup>2</sup>Depth by using the per-pixel photometric error  $\mathcal{L}_{\text{photo}}$  as:

$$\mathcal{L}_{\text{photo}} = \frac{\alpha}{2}(1 - \text{SSIM}(\mathbf{I}_a, \mathbf{I}_b)) + (1 - \alpha)\|\mathbf{I}_a - \mathbf{I}_b\|_1, \quad (11)$$

where SSIM is the structural similarity between two images [36],  $\mathbf{I}_a$  and  $\mathbf{I}_b$  indicate the ground truth image and the reconstructed image respectively. It is noteworthy that M<sup>2</sup>Depth uses  $\mathcal{L}_{\text{photo}}$  in both spatial domain and temporal domain, where the spatial photometric error additionally provides the important scale information.

*Depth Smoothness Loss* As in previous works [11, 12, 38], we use the edge-aware smoothness loss [11] to prevent estimated depth from shrinking:

$$\mathcal{L}_{\text{smooth}} = |\partial_x \mathbf{d}|e^{-|\partial_x \mathbf{I}|} + |\partial_y \mathbf{d}|e^{-|\partial_y \mathbf{I}|}. \quad (12)$$

*Depth Edge Loss* Inspired by [32], we employ edge information derived from images to enhance the quality of depth edges. Given an RGB image  $\mathbf{I}$  and its depth map  $\mathbf{d}$ , we utilize a pre-trained edge detection model [19] to extract the edge map  $\mathbf{E}_{\text{img}}$  from  $\mathbf{I}$ . Subsequently, the edge map  $\mathbf{E}_{\text{depth}}$  of  $\mathbf{d}$  can be calculated by depth gradient. The depth edge loss is defined as:

$$\mathcal{L}_{\text{edge}} = \text{FL}(\mathbf{E}_{\text{img}}, \mathbf{E}_{\text{depth}}) \quad (13)$$

where FL denote the focal loss [26].

*SfM Loss* Although the  $\mathcal{L}_{\text{photo}}$  in spatial domains could constrain the scale of estimated depth and pose, it relies on good initialization. Following the previous works [16, 38], we use SfM to generate sparse depth between spatially adjacent views to endow the network with an initial rudimentary estimation of scale-aware depth. More details of SfM loss  $\mathcal{L}_{\text{sfm}}$  can be found in supplementary materials.

*Total Loss* The overall loss function can be written as:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{photo}} + \lambda_2 \mathcal{L}_{\text{smooth}} + \lambda_3 \mathcal{L}_{\text{edge}} + \lambda_4 \mathcal{L}_{\text{sfm}}, \quad (14)$$

where the  $\lambda_1, \lambda_2, \lambda_3, \lambda_4$  are weights of different losses, we set  $\lambda_1 = 1.0, \lambda_2 = 1.0e-3, \lambda_3 = 1.0e-2, \lambda_4 = 1.0e-2$ , and the  $\lambda_4$  is set to 0 after initialization.

## 4 Experiments

### 4.1 Implementation Details

We implement M<sup>2</sup>Depth using PyTorch and train the model using Adam as optimizer with a learning rate set to  $10^{-4}$ ,  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . For the pose estimation, we employ a ResNet-34 [20] model to predict the axis-angle and translation of the front camera. For the depth prior estimation branch, we employ the ResNet-34 [20] model to predict internal features and use the frozen SAM encoder provided by MobileSAM [45] for saving memory. All of our experiments are conducted using 8 NVIDIA V100 GPUs.

### 4.2 Dataset

We train and evaluate M<sup>2</sup>Depth on two public datasets including DDAD [14] and nuScenes [5].

The dense depth for autonomous driving (DDAD) dataset [14] is an autonomous driving benchmark that consists of 150 training and 50 validation scenes in complex and diverse urban environments. Following the previous work [29, 38], we downsample the images from their initial resolution of  $1216 \times 1936$  to  $384 \times 640$  and evaluate depth up to 200m averaged across all cameras.

The nuScenes dataset [5] comprises 700 training, 150 validation, and 150 testing urban scenes. Following the previous work [38], we downsample the images from the initial resolution of  $900 \times 1600$  to  $352 \times 640$  and evaluate depth up to 80m averaged across all cameras.

### 4.3 Experimental Results

This paper focuses on scale-aware surrounding depth estimation task, thus we only report the scale-aware results and mainly compare M<sup>2</sup>Depth with the recent self-supervised surrounding depth estimation methods, including FSM [16], SurroundDepth [38], VFDepth [21] and R3D3 [29], without comparing with the numerous MDE methods [2, 12, 35, 37].

**Table 1:** Quantitative results on DDAD dataset [14] (evaluate depth up to 200m) and nuScenes dataset [5] (evaluate depth up to 80m). We present the mean accuracy across all views using the metrics from [9]. The *Frame* stands for the number of frames in the training\testing phase. FSM\* indicates the implementation from [21]. R3D3\* denotes the results using the official code and the same frame setting with us. (**Bold** figures indicate the best and underlined figures indicate the second best)

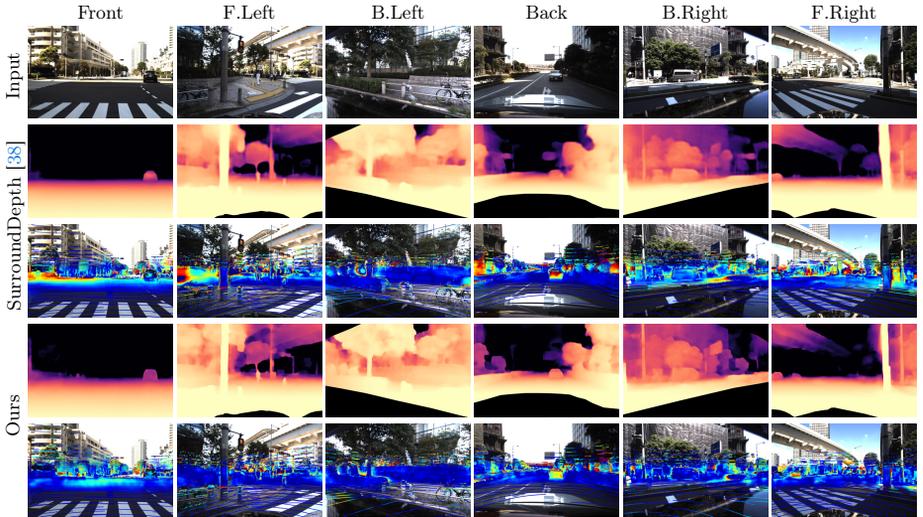
Method	Dataset	Frame	Abs. Rel. ↓	Sq. Rel. ↓	RMSE ↓	RMSE log ↓	$\delta < 1.25 \uparrow$	$\delta < 1.25^2 \uparrow$	$\delta < 1.25^3 \uparrow$
R3D3 [29]		6\5	0.169	3.041	11.372	-	0.809	-	-
FSM [16]	DDAD [14]	3\1	<u>0.201</u>	-	-	-	-	-	-
FSM* [16]		3\1	0.228	4.409	13.433	0.342	0.687	0.870	0.932
VFDepth [21]		3\1	0.218	3.660	13.327	0.339	0.674	0.862	0.932
SurroundDepth [38]		3\1	0.208	<u>3.371</u>	<u>12.977</u>	<u>0.330</u>	<u>0.693</u>	<u>0.871</u>	<u>0.934</u>
R3D3* [29]		3\2	0.311	5.473	14.094	0.385	0.604	0.814	0.903
Ours		3\2	<b>0.183</b>	<b>2.920</b>	<b>11.963</b>	<b>0.299</b>	<b>0.756</b>	<b>0.897</b>	<b>0.947</b>
R3D3 [29]		6\5	0.253	4.759	7.150	-	0.729	-	-
FSM [16]	nuScenes [5]	3\1	<u>0.297</u>	-	-	-	-	-	-
FSM* [16]		3\1	0.319	7.534	7.860	0.362	<u>0.716</u>	<u>0.874</u>	<u>0.931</u>
VFDepth [21]		3\1	0.289	5.718	7.551	<u>0.348</u>	0.709	<b>0.876</b>	<b>0.932</b>
SurroundDepth [38]		3\1	0.280	<b>4.401</b>	<u>7.467</u>	0.364	0.661	0.844	0.917
R3D3* [29]		3\2	0.498	5.489	11.740	0.746	0.155	0.375	0.613
Ours		3\2	<b>0.259</b>	<u>4.599</u>	<b>6.898</b>	<b>0.332</b>	<b>0.734</b>	0.871	0.928

**Table 2:** Per-camera evaluation on DDAD dataset [14]. SD is the abbreviation of SurroundDepth [38]. R3D3\* indicates the results using its official code and the same frame setting with us. Our method achieves superior overall performance across multiple cameras to existing works. According to the memory and computation analysis, M<sup>2</sup>Depth achieves a good balance between overall performance and computational efficiency.

Method	Abs.Rel. ↓							Memory & Computation		
	Front	F.Left	F.Right	B.Left	B.Right	Back	Avg.	Memory(MB)	Flops(G)	Time(s)
FSM [16]	<b>0.130</b>	<u>0.201</u>	<u>0.224</u>	0.229	0.240	<u>0.186</u>	<u>0.201</u>	-	-	-
SD [38]	0.152	0.207	0.230	<u>0.220</u>	<u>0.239</u>	0.200	0.208	<b>3042</b>	<b>237.106</b>	<b>0.215</b>
R3D3* [29]	0.234	0.284	0.355	0.347	0.392	0.255	0.311	8371	2621.738	0.378
Ours	<u>0.146</u>	<b>0.182</b>	<b>0.200</b>	<b>0.198</b>	<b>0.203</b>	<b>0.169</b>	<b>0.183</b>	<u>5546</u>	<u>866.019</u>	<u>0.295</u>

*Results on DDAD* Following the common practice, we report the quantitative results of *Abs.Rel.*, *Sq.Rel.*, *RMSE*, *RMSE log* and  $\delta$  as shown in Tab. 1. The specific definition of the evaluation metrics can be found in supplementary materials. Previous works [21, 38] typically use three frames  $[t - 1, t, t + 1]$  for training, where the  $t + 1$  frame is only used in computing loss, we follow this paradigm to train M<sup>2</sup>Depth. It is noteworthy that the R3D3 [29] takes sequence frames as input, we test its results with 2 frames using their official code for a fair comparison.

As shown in Tab. 1, M<sup>2</sup>Depth achieves significant improvement on all metrics compared with existing methods when tested in a similar setting. To be specific, our method outperforms the SOTA method of single-frame surrounding depth estimation, SurroundDepth [38], by 12.02% on *Abs. Rel.* and 13.38% on *Sq. Rel.*, indicating that our usage of spatial-temporal volumes substantially improves the depth quality. We also compare the visualization results of M<sup>2</sup>Depth and



**Fig. 5:** Qualitative comparison of predicted surrounding depth on DDAD dataset [14]. Given the input surrounding images (the top row), we show the visualized depth maps and depth errors of SurroundDepth [38] and  $M^2$ Depth. The depth maps are visualized in the range of  $[0, 50m]$ . Our method is able to produce more accurate depth with less error and sharper depth edge across multiple cameras.

SurroundDepth in Fig. 5, where the estimated surrounding depth and depth errors show our method produces more accurate and consistent depth predictions among multiple cameras in challenging scenarios. In Tab. 2, we show the per-camera evaluation results on DDAD. In terms of *Abs.Rel.*, our method is able to predict more accurate depth in nearly all cameras, demonstrating the superior performance of  $M^2$ Depth.

*Results on nuScenes* In Tab. 1, we evaluate the proposed  $M^2$ Depth on the evaluation set of nuScenes dataset [5], where the quantitative results show our method significantly outperforms the existing method in terms of multiple metrics in similar setting. Compared with R3D3 [29] that use 5 frames as input, our method utilizes only 2 frames and achieves comparable performance on *Abs. Rel.* and superior performance on other metrics. As the test data in nuScenes dataset [5] is more challenging than DDAD dataset [14], the aforementioned results indicate that  $M^2$ Depth achieves state-of-the-art overall performance.

*Memory and Computation Analysis.* As shown in Tab. 2, compared with R3D3 [29] and SurroundDepth [38],  $M^2$ Depth achieves a good balance between overall performance and computational efficiency. According to the results, our method consumes much less *memory* and *FLOPs* than R3D3 [29] while achieving competitive performance.

**Table 3:** Quantitative results of the ablation study on DDAD dataset [14]. M.P. stands for mono prior, S.Vol. and T.Vol. indicate the spatial volume and temporal volume, the STF is the proposed spatial-temporal fusion module. Jointly constructing spatial-temporal cost volumes significantly improves the depth quality compared with the mono prior depth, and the STF further increases the capabilities of M<sup>2</sup>Depth on nearly all metrics.

No.	M.P.	S.Vol.	T.Vol.	STF	Abs. Rel.	Sq. Rel.	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
1	✓				0.216	3.758	13.200	0.338	0.686	0.863	0.929
2	✓	✓			0.212	3.662	12.959	0.326	0.696	0.872	0.936
3	✓	✓	✓		<u>0.197</u>	<u>3.379</u>	<u>12.341</u>	<u>0.313</u>	<u>0.738</u>	<b>0.886</b>	<b>0.941</b>
4	✓	✓	✓	✓	<b>0.194</b>	<b>3.331</b>	<u>12.347</u>	<b>0.311</b>	<b>0.741</b>	<b>0.886</b>	<b>0.941</b>

#### 4.4 Ablation Study

In this section, we conduct ablation studies on DDAD [14] to analyze the effectiveness of each module of M<sup>2</sup>Depth.

*Spatial-temporal Volume* In Tab. 3, we evaluate the performance using different cost volumes, where the base model uses none of the spatial volume (*S. Vol.*), temporal volume (*T. Vol.*) and STF module. By fusing the spatially adjacent views, *S. Vol.* improves 2.55% on the *Sq. Rel.* metric and 3.55% on the *RMSE log* metric. When further injecting the temporal information into cost volumes, the *T. Vol.* achieves more than 7% improvement on *Abs. Rel* and *Sq. Rel.* The aforementioned results show that integrating the spatial-temporal information is able to significantly strengthen the depth quality.

*Spatial-temporal Fusion* Compared with directly using features that warp from previous frames or adjacent views, the proposed STF module fuses the volume features within the spatial-temporal domain, where the updated feature integrates the global information with spatial- and temporal- features and consequently strengthens the feature expressiveness. Compared to only using information from a single domain, our volumes achieve better performance.

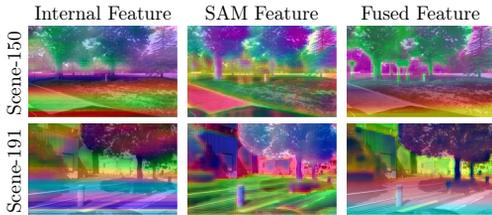
*Multi-grained Feature Fusion* As shown in Tab. 4, we conduct the ablation study to evaluate the effectiveness of the MFF module, which enhances feature learning by combining SAM features with internal features. According to the quantitative results, introducing MFF into mono prior estimation achieves improvement in nearly all metrics. We also show the visualized features in Fig. 6, where the internal features represent the geometric distance infor-

**Table 4:** Ablation studies of  $\mathcal{L}_{\text{edge}}$ , MFF and D. D. on DDAD dataset [14], where  $\mathcal{L}_{\text{edge}}$  indicates the depth edge loss, MFF stands for multi-grained feature fusion module, D. D. represents the depth decoding with SAM features.

$\mathcal{L}_{\text{edge}}$	MFF	D. D.	Abs. Rel.	Sq. Rel.	RMSE	$\delta < 1.25$
			0.194	3.331	12.347	0.741
✓			0.192	3.224	12.447	0.741
✓	✓		<u>0.188</u>	<u>3.032</u>	12.213	<u>0.748</u>
✓		✓	0.191	3.262	<u>12.175</u>	<u>0.748</u>
✓	✓	✓	<b>0.183</b>	<b>2.920</b>	<b>11.963</b>	<b>0.756</b>

mation and the SAM features contain semantic instance information. By combining the internal features and SAM features in latent space, the fused features derive a comprehensive understanding of the surrounding environment.

As mentioned in Sec. 3.3, we further integrate SAM priors in depth decoding and conduct experimental comparison in Tab. 4, where the results show that utilizing SAM features in both MFF and depth decoding achieves the best performance of M<sup>2</sup>Depth.



**Fig. 6:** Visualization results of different features in M<sup>2</sup>Depth on DDAD dataset [14]. The internal feature is from the internal image encoder, the SAM feature is from the frozen SAM encoder [45], and the fused feature is produced by the MFF module.

*Edge Loss and Others* The ablation study of  $\mathcal{L}_{\text{edge}}$  is performed in Tab. 4, where the results indicate  $\mathcal{L}_{\text{edge}}$  is able to improve the *Abs. Rel.* and *Sq. Rel.*. We also perform ablation studies for other hyper-parameters and candidate designs, please refer to supplementary materials for more results and analysis.

## 5 Conclusion

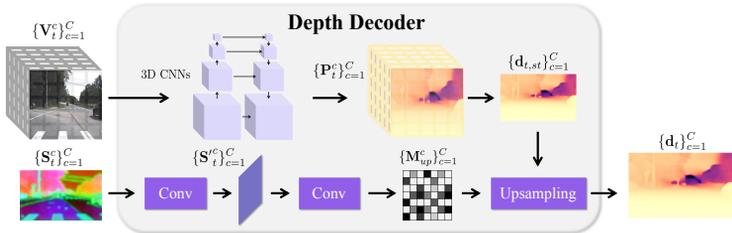
*Limitation* Currently, M<sup>2</sup>Depth constructs as many volumes as the number of cameras, which consumes a lot of memory when increasing the cameras. In the future, we’d like to build a unified cost volume to represent the surrounding environment.

*Conclusion* In this paper, we propose M<sup>2</sup>Depth which is designed for the self-supervised two-frame multi-camera metric depth estimation task in autonomous driving. Different from the previous methods that use single frame or single camera, M<sup>2</sup>Depth takes two-frame from multi-camera as inputs and learns to construct spatial-temporal cost volumes, which is the first method to exploit spatial-temporal fusion in constructing cost volumes. We additionally propose a novel multi-grained feature fusion module to combine the SAM priors with internal features. Experimental results on two public benchmarks indicate that M<sup>2</sup>Depth achieves state-of-the-art performance.

## A Implementation Details

### A.1 Depth Decoder

The detailed structure of the depth decoder is illustrated in Fig. 7. Given the spatial-temporal volume  $\{\mathbf{V}_t^c\}_{c=1}^C$  and the SAM feature  $\{\mathbf{S}_t^c\}_{c=1}^C$  from SAM encoder [45], we first transform  $\{\mathbf{V}_t^c\}_{c=1}^C$  into probability volumes  $\{\mathbf{P}_t^c\}_{c=1}^C$  by 3D CNNs. Then, we calculate the spatial-temporal depth  $\{\mathbf{d}_{t,st}^c\}_{c=1}^C$  using depth samples. Subsequently, we utilize  $\{\mathbf{S}_t^c\}_{c=1}^C$  as context features to compute the upsampling mask  $\{\mathbf{M}_{up}^c\}_{c=1}^C$ . Finally, by integrating  $\{\mathbf{M}_{up}^c\}_{c=1}^C$  and  $\{\mathbf{d}_{t,st}^c\}_{c=1}^C$ , we can obtain the final depth  $\{\mathbf{d}_t^c\}_{c=1}^C$ .



**Fig. 7:** Overview of Depth decoder. Given the spatial-temporal volume  $\{\mathbf{V}_t^c\}_{c=1}^C$  and the SAM feature  $\{\mathbf{S}_t^c\}_{c=1}^C$  as inputs, we initially compute the spatial-temporal depth  $\{\mathbf{d}_{t,st}^c\}_{c=1}^C$ . Subsequently, the  $\{\mathbf{d}_{t,st}^c\}_{c=1}^C$  is upsampled with the mask  $\{\mathbf{M}_{up}^c\}_{c=1}^C$  which are calculated from  $\{\mathbf{S}_t^c\}_{c=1}^C$  to procure the final depth  $\{\mathbf{d}_t^c\}_{c=1}^C$ .

### A.2 Adaptive Depth Sample

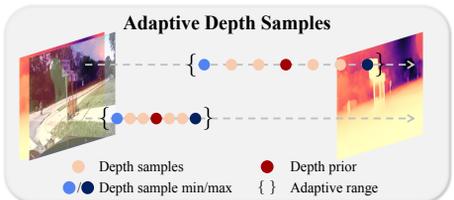
Following the plane sweep paradigm, the selection of depth samples directly affects the depth quality. Previous methods [8, 13, 42] usually adopt a wide-range sampling strategy for the entire scene, which improves the accuracy of depth estimation to some extent, but also brings a huge computational burden.

To solve this problem, we propose utilizing the mono depth estimation result as prior information and conducting adaptive sampling in the vicinity of the prior depth. This method not only significantly reduces the computational complexity, but also improves the efficiency of depth estimation.

The method of adaptive depth sampling is shown in Fig. 8. Specifically, we determine the range of depth sampling  $[\mathbf{d}_{\min}(\mathbf{p}), \mathbf{d}_{\max}(\mathbf{p})]$  for each pixel  $\mathbf{p}$  based on the given depth  $\mathbf{d}_{\text{init}}$  and scaling factor  $\alpha$  as follow:

$$\mathbf{d}_{\min}(\mathbf{p}) = \mathbf{d}_{\text{init}}(\mathbf{p}) \div (1 + \alpha), \quad (15)$$

$$\mathbf{d}_{\max}(\mathbf{p}) = \mathbf{d}_{\text{init}}(\mathbf{p}) \times (1 + \alpha), \quad (16)$$



**Fig. 8:** We illustrate the examples of the adaptive depth sample, where the depth range increases for pixels at a farther distance, and conversely, decreases for pixels at a closer proximity.

It is evident from this formula that the depth range varies with the depth. When the  $\mathbf{d}_{\text{init}}(\mathbf{p})$  is large, that is, the object is farther away, the range of depth sampling will increase accordingly; conversely, when the  $\mathbf{d}_{\text{init}}(\mathbf{p})$  is small, the range of depth sampling will decrease. This adaptive depth sampling strategy is more in line with the depth distribution of actual scenes, thus effectively improving the quality of depth.

### A.3 Structure-from-Motion Loss

Through self-supervised photometric loss  $\mathcal{L}_{\text{photo}}$ , we can effectively supervise the estimated depth and pose. However, during the initial phase of training, obtaining valid projection results is challenging due to insufficient overlap between adjacent cameras, which ultimately renders supervision ineffective. To address this issue, we follow previous methods [16, 38] and obtain scale-aware depth through triangulation of adjacent cameras utilizing their camera extrinsics, which serves as pseudo labels for effective supervision. By doing so, we successfully enhance the accuracy of depth and pose estimation by leveraging information from neighboring cameras and extrinsics.

The calculation for  $\mathcal{L}_{\text{sfm}}$  is as follows:

$$\mathcal{L}_{\text{sfm}} = \frac{1}{|\mathbb{M}|} \sum_{\mathbf{p} \in \mathbb{M}} |\mathbf{d}(\mathbf{p}) - \mathbf{d}_{\text{sfm}}(\mathbf{p})|_1, \quad (17)$$

where  $\mathbb{M}$  represents the set of valid pixel  $\mathbf{p}$  in pseudo depth labels  $\mathbf{d}_{\text{sfm}}$ .

### A.4 Evaluation Metrics

Following in previous work [16, 38], the description of the evaluation metrics we used is as follows:

$$\text{Abs.Rel.} := \frac{1}{|\mathbb{N}|} \sum_{\mathbf{p} \in \mathbb{N}} \frac{|\mathbf{d}(\mathbf{p}) - \mathbf{d}^*(\mathbf{p})|}{\mathbf{d}^*(\mathbf{p})}, \quad (18)$$

$$\text{Sq. Rel.} := \frac{1}{|\mathbb{N}|} \sum_{\mathbf{p} \in \mathbb{N}} \frac{\|\mathbf{d}(\mathbf{p}) - \mathbf{d}^*(\mathbf{p})\|^2}{\mathbf{d}^*(\mathbf{p})}, \quad (19)$$

$$\text{RMSE} := \frac{1}{|\mathbb{N}|} \sqrt{\sum_{\mathbf{p} \in \mathbb{N}} \|\mathbf{d}(\mathbf{p}) - \mathbf{d}^*(\mathbf{p})\|^2}, \quad (20)$$

$$\text{RMSE log} := \frac{1}{|\mathbb{N}|} \sqrt{\sum_{\mathbf{p} \in \mathbb{N}} \|\log \mathbf{d}(\mathbf{p}) - \log \mathbf{d}^*(\mathbf{p})\|^2}, \quad (21)$$

$$\delta < n: \text{fraction of } d \in \mathbf{d} \text{ for which } \max\left(\frac{d}{d^*}, \frac{d^*}{d}\right) < n, \quad (22)$$

where  $\mathbf{d}$  and  $\mathbf{d}^*$  indicate the predicted depth and ground-truth depth respectively.  $\mathbb{N}$  indicates the all valid pixels  $\mathbf{p}$  in  $\mathbf{d}^*$ .

## B Computation Analysis

In Tab. Tab. 5, we show the computation cost of each module. It can be observed that the cost volume construction and fusion occupy a high proportion of memory and time, as the grid sample operation is well known to be time-consuming. Reducing the runtime in V.C.F is an important future work.

**Table 5:** Computation analysis of each module: Pose Branch (Pose), Image Encoder (I.E.), SAM Encoder (S.E.), Prior Decoder (P.D.), Volume Construct & Fusion (V.C.F.), Depth Decoder (D.D.). Experiments are performed on V100.

	Pose	I.E.	S.E.	MFF	P.D.	V.C.F.	D.D.
Memory(MB)	139.20	139.07	173.03	51.10	105.39	397.12	196.33
Percent(%)	11.59%	11.58%	14.40%	4.25%	8.77%	33.06%	16.34%
Time(ms)	39.33	3.35	20.65	3.58	1.39	216.35	2.34
Percent(%)	13.71%	1.17%	7.20%	1.25%	0.48%	75.39%	0.81%

## C Ablation Study

*Design of Pose Estimation* Tab. 6 shows that the *Front Camera (F. Cam.)* can achieve better results. We take the previous method [38] which concatenates surrounding views to directly predict the ego pose as the baseline *Concat Camera (C. Cam.)*. Experiments indicate that the method *F. Cam.*, which predict the pose of front-view camera  $\mathbf{P}_{t \rightarrow t-1}^0$  and then derive the ego pose  $\mathbf{P}_{t \rightarrow t-1}$ , is more effective.

*Design of Multi-grained Feature Fusion Module* In Tab. 7, we evaluate the performance of different feature fusion methods in mono prior estimation. Specifically, we compare the base model, which does not utilize the MFF module, against the multi-grained feature fusion (MFF) module and the vanilla feature fusion (VFF) module that blends SAM features with internal features through simple addition. The results presented in Tab. 7 demonstrate that the incorporation of SAM features notably elevates the quality of depth estimation outcomes. Comparing the MFF module with the VFF module, our multi-grained feature fusion module exhibits superior performance in fusing internal features with fine-grained semantic information, thereby further augmenting the precision of depth estimation.

*Design of Depth Decoder* For Tab. 8, we train two variants of our depth decoder: Vanilla Refine (*V. Refine*) and SAM Refine (*S. Refine*). The former utilizes context features from FPN [25], whereas the latter employs context features from the SAM encoder [22]. Through evaluation on the DDAD dataset, *S. Refine* attains superior results. The results show that the network necessitates the integration of more fine-grained information to enhance depth refinement. When compared

**Table 6:** Ablation study on the design of pose estimation module comparison. Experiments demonstrate that the method, which utilizes the front-view camera to estimate the front-view pose and subsequently infer the ego pose, is well-suited for our depth estimation network and embodies its effectiveness. (**Bold** figures indicate the best and underlined figures indicate the second best)

Method	Abs. Rel.	Sq. Rel.	RMSE	RMSE	log $\delta < 1.25$
C. Cam.	<u>0.189</u>	<u>2.942</u>	<u>12.239</u>	<u>0.309</u>	<u>0.732</u>
F. Cam.	<b>0.183</b>	<b>2.920</b>	<b>11.963</b>	<b>0.299</b>	<b>0.756</b>

**Table 8:** Designs of depth decoder comparison. We train SAM Refine (*S. Refine*) as described in the main paper and train Vanilla Refine (*V. Refine*) using the context feature from FPN [25]. We evaluate both the network on DDAD and the experiments show that SAM Refine effectively enhances depth quality. (**Bold** figures indicate the best and underlined figures indicate the second best)

Method	Abs. Rel.	Sq. Rel.	RMSE	RMSE	log $\delta < 1.25$
Base	<u>0.192</u>	<b>3.224</b>	12.447	<u>0.312</u>	<u>0.741</u>
V. Refine	0.196	3.313	<u>12.366</u>	0.313	0.734
S. Refine	<b>0.191</b>	<u>3.262</u>	<b>12.175</b>	<b>0.305</b>	<b>0.748</b>

**Table 7:** Designs of feature fusion module comparison. We train MFF as described in the main paper and train the VFF module which fuses the internal feature and SAM feature through direct addition. Experimental results demonstrate that our design effectively integrates diverse-grained features, thereby significantly enhancing the quality of depth estimation. (**Bold** figures indicate the best and underlined figures indicate the second best)

Method	Abs. Rel.	Sq. Rel.	RMSE	RMSE	log $\delta < 1.25$
Base	0.191	3.262	<u>12.175</u>	<u>0.305</u>	<u>0.748</u>
VFF	<u>0.185</u>	<u>3.044</u>	12.209	0.307	0.746
MFF	<b>0.183</b>	<b>2.920</b>	<b>11.963</b>	<b>0.299</b>	<b>0.756</b>

**Table 9:** Designs of depth sample comparison. We train Adaptive Sample (*A. Sample*), Vanilla Sample (*V. Sample*) and Fixed Sample (*F. Sample*) with 16 samples. We evaluate both the network on DDAD and the experiments show that using adaptive methods yields better results. (**Bold** figures indicate the best and underlined figures indicate the second best)

Method	Abs. Rel.	Sq. Rel.	RMSE	RMSE	log $\delta < 1.25$
V. Sample	0.362	5.932	14.891	0.422	0.534
F. Sample	<u>0.195</u>	<u>3.054</u>	<u>12.362</u>	<u>0.309</u>	<u>0.721</u>
A. Sample	<b>0.183</b>	<b>2.920</b>	<b>11.963</b>	<b>0.299</b>	<b>0.756</b>

to FPN features, which encompass feature-matching information, SAM features are deemed more suitable.

*Adaptive Depth Sample* In Tab. 9, we perform a comparison between the adaptive depth samples as described in the main paper (*A. Sample*), the fixed depth samples within a fixed depth sampling range (*F. Sample*), the vanilla depth sample within the entire space (*V. Sample*). The experimental results consistently show that the adaptive method yields better outcomes.

*Number of Bins* We conduct an ablation study against the number of bins on DDAD [14] dataset, and the results are shown in Tab. 10. Our results demonstrate that increasing the quantity of bins does not significantly enhance the quality of depth. This indicates that the utilization of adaptive depth samples effectively contributes to improving computational efficiency.

**Table 10:** Ablation study on number of bins. We compare the influence of the different number of bins used to train the network. (**Bold** figures indicate the best and underlined figures indicate the second best)

Bins	Abs. Rel.	Sq. Rel.	RMSE	$\delta < 1.25$	Memory(MB)
8	<u>0.195</u>	<u>3.316</u>	<u>12.349</u>	<u>0.740</u>	<b>3483</b>
16	<b>0.194</b>	3.331	<b>12.347</b>	<b>0.741</b>	<u>3853</u>
32	0.200	<b>3.264</b>	12.491	0.724	4751

**Table 11:** Ablation study on number of frames. The experimental results demonstrate that our method achieves highly competitive results with just two frames. (**Bold** figures indicate the best and underlined figures indicate the second best)

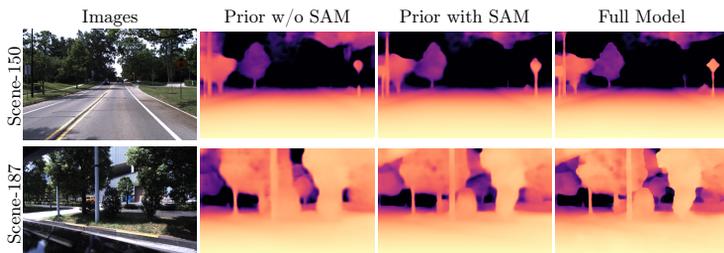
Frames	Abs. Rel.	Sq. Rel.	RMSE	RMSE	$\log \delta < 1.25$
(-1, 0)	<b>0.183</b>	<u>2.920</u>	<b>11.963</b>	<b>0.299</b>	<b>0.756</b>
(-2, -1, 0)	<u>0.185</u>	2.956	<u>12.100</u>	<u>0.301</u>	<u>0.747</u>
(-3, -2, -1, 0)	0.186	<b>2.911</b>	12.185	0.303	0.740

*More Frames* We conduct a multi frames experiment using multiple frames (2 frames, 3 frames, 4 frames) as inputs for depth estimation. Tab. 11 reveals that increasing the number of frames does not necessarily improve depth accuracy. As our method is not specifically designed to handle sequence data, increasing the input frames does not effectively contribute new information. Notably, employing just two frames is sufficient to produce commendable results.

## D Visualized

### D.1 SAM Feature Enhanced Depth

As shown in Fig. 9, integrating SAM features gets a notable enhancement in both the depth prior and the final depth, particularly evident at the edges of the instance.



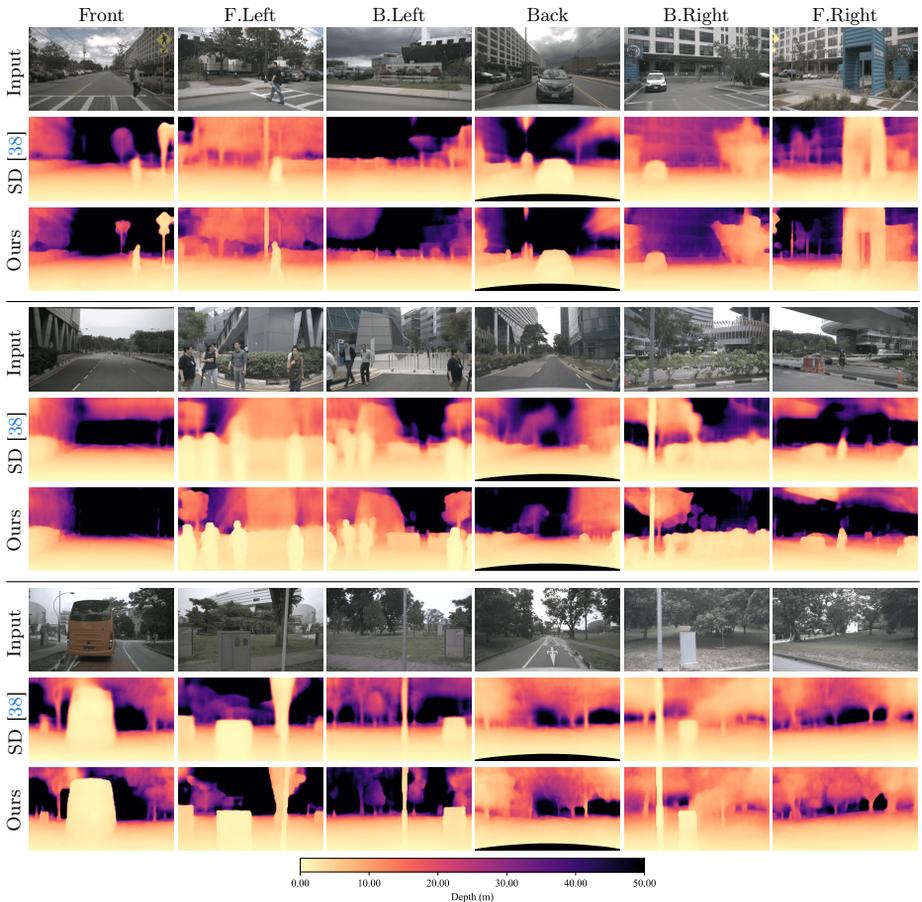
**Fig. 9:** Visualization of produced depth results on DDAD dataset [14]. It can be observed clearly that consistency within instances and discrimination between different instances for both depths has improved.

### D.2 More Depth Results

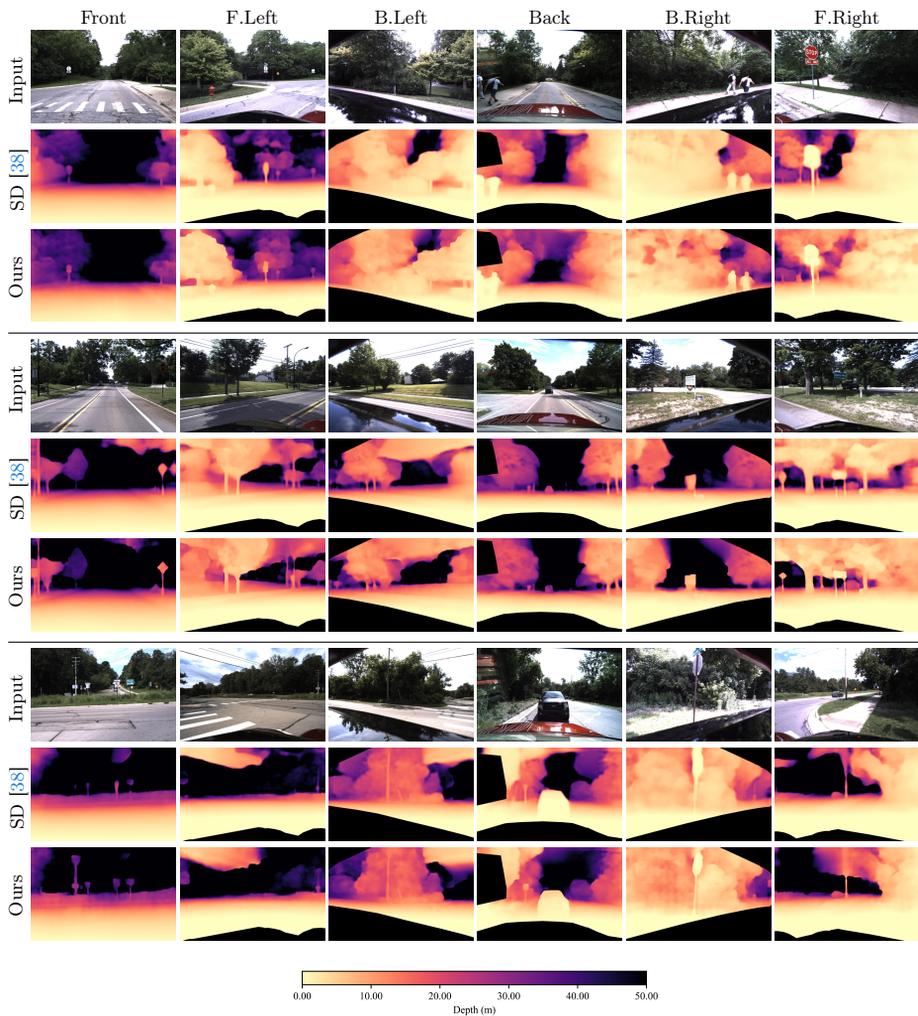
We visualize more depth results in Nuscenes [5] and DDAD [14] dataset. In Fig. 10 and Fig. 11, our M<sup>2</sup>Depth consistently exhibits robustness and effectiveness across diverse scenes. Notably, at the object edges, our method produces sharper depth predictions.

### D.3 More Depth Error Results

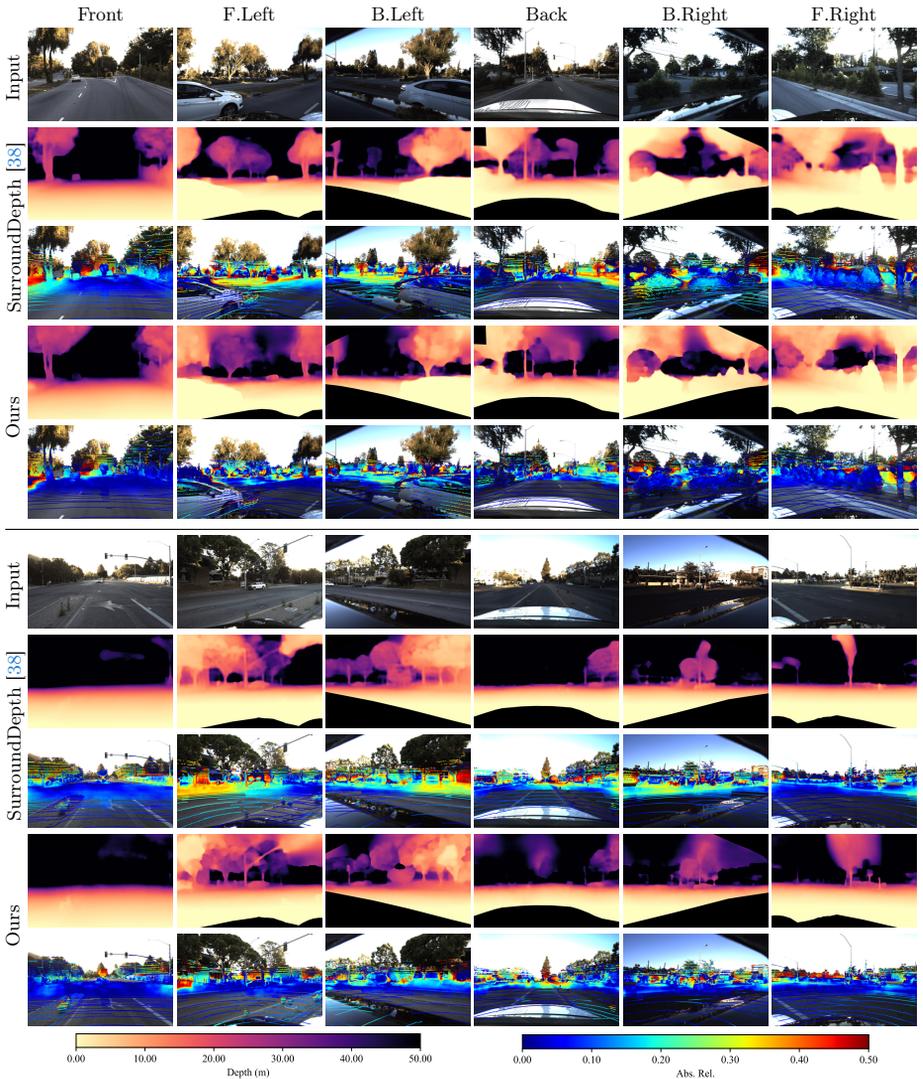
In Fig. 12, we qualitatively compare our method with existing works in terms of scale-aware depth estimation in DDAD. It can be observed that our method achieves better results at the overlapping between adjacent views.



**Fig. 10:** Qualitative comparison of predicted surrounding depth on NuScenes [5]. We show a comparison of depth maps from our method to the depth maps of the state-of-the-art approach SurroundDepth [38]. We observe that our method produces significantly sharper and more accurate depth predictions, particularly in fine details.



**Fig. 11:** Qualitative comparison of predicted surrounding depth on DDAD [14]. We show a comparison of depth maps from  $M^2$ Depth to the depth maps of the state-of-the-art approach SurroundDepth [38]. We observe that our method produces significantly sharper and more accurate depth predictions, particularly in fine details.



**Fig. 12:** Qualitative comparison of predicted surrounding depth on DDAD dataset [14]. Given the input surrounding images (the top row), we show the visualized depth maps and depth errors of SurroundDepth [38] and M<sup>2</sup>Depth. Our method is able to produce more accurate depth with less error and sharper depth edge across multiple cameras.

## References

1. Bae, G., Budvytis, I., Cipolla, R.: Multi-view depth estimation by fusing single-view depth probability with multi-view geometry. In: CVPR. pp. 2842–2851 (2022) [3](#)
2. Bhat, S.F., Alhashim, I., Wonka, P.: Adabins: Depth estimation using adaptive bins. In: CVPR. pp. 4009–4018 (2021) [3](#), [10](#)
3. Bian, J., Li, Z., Wang, N., Zhan, H., Shen, C., Cheng, M.M., Reid, I.: Unsupervised scale-consistent depth and ego-motion learning from monocular video. *Advances in neural information processing systems* **32** (2019) [3](#)
4. Bui, N.T., Hoang, D.H., Tran, M.T., Le, N.: Sam3d: Segment anything model in volumetric medical images. arXiv preprint arXiv:2309.03493 (2023) [5](#)
5. Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. In: CVPR. pp. 11621–11631 (2020) [3](#), [10](#), [11](#), [12](#), [5](#), [7](#)
6. Cheng, Y., Li, L., Xu, Y., Li, X., Yang, Z., Wang, W., Yang, Y.: Segment and track anything. arXiv preprint arXiv:2305.06558 (2023) [5](#)
7. Collins, R.T.: A space-sweep approach to true multi-image matching. In: Proceedings CVPR IEEE computer society conference on computer vision and pattern recognition. pp. 358–363. Ieee (1996) [2](#), [6](#)
8. Ding, Y., Yuan, W., Zhu, Q., Zhang, H., Liu, X., Wang, Y., Liu, X.: Transmvsnet: Global context-aware multi-view stereo network with transformers. In: CVPR. pp. 8585–8594 (2022) [6](#), [1](#)
9. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems* **27** (2014) [11](#)
10. Feng, Z., Yang, L., Jing, L., Wang, H., Tian, Y., Li, B.: Disentangling object motion and occlusion for unsupervised multi-frame monocular depth. In: ECCV. pp. 228–244 (2022) [4](#)
11. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: CVPR. pp. 270–279 (2017) [9](#)
12. Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth prediction. In: ICCV. pp. 3828–3838 (2019) [1](#), [2](#), [3](#), [9](#), [10](#)
13. Gu, X., Fan, Z., Zhu, S., Dai, Z., Tan, F., Tan, P.: Cascade cost volume for high-resolution multi-view stereo and stereo matching. In: CVPR. pp. 2495–2504 (2020) [6](#), [1](#)
14. Guizilini, V., Ambrus, R., Pillai, S., Raventos, A., Gaidon, A.: 3d packing for self-supervised monocular depth estimation. In: CVPR. pp. 2485–2494 (2020) [2](#), [3](#), [10](#), [11](#), [12](#), [13](#), [14](#), [4](#), [5](#), [8](#), [9](#)
15. Guizilini, V., Ambrus, R., Chen, D., Zakharov, S., Gaidon, A.: Multi-frame self-supervised depth with transformers. In: CVPR. pp. 160–170 (2022) [1](#), [3](#)
16. Guizilini, V., Vasiljevic, I., Ambrus, R., Shakhnarovich, G., Gaidon, A.: Full surround monodepth from multiple cameras. *IEEE Robotics and Automation Letters (RA-L)* pp. 5397–5404 (2022) [1](#), [2](#), [10](#), [11](#)
17. Guo, X., Yang, K., Yang, W., Wang, X., Li, H.: Group-wise correlation stereo network. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 3273–3282 (2019) [8](#)
18. He, H., Zhang, J., Xu, M., Liu, J., Du, B., Tao, D.: Scalable mask annotation for video text spotting. arXiv preprint arXiv:2305.01443 (2023) [5](#)

19. He, J., Zhang, S., Yang, M., Shan, Y., Huang, T.: Bi-directional cascade network for perceptual edge detection. In: CVPR. pp. 3828–3837 (2019) [9](#)
20. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016) [10](#)
21. Kim, J.H., Hur, J., Nguyen, T.P., Jeong, S.G.: Self-supervised surround-view depth estimation with volumetric feature fusion. In: NeurIPS. pp. 4032–4045 (2022) [2](#), [4](#), [10](#), [11](#)
22. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. In: ICCV. pp. 4015–4026 (2023) [3](#), [5](#), [6](#)
23. Li, R., Gong, D., Yin, W., Chen, H., Zhu, Y., Wang, K., Chen, X., Sun, J., Zhang, Y.: Learning to fuse monocular and multi-view cues for multi-frame depth estimation in dynamic scenes. In: CVPR. pp. 21539–21548 (2023) [2](#), [4](#)
24. Li, Z., Wang, X., Liu, X., Jiang, J.: Binsformer: Revisiting adaptive bins for monocular depth estimation. arXiv preprint arXiv:2204.00987 (2022) [3](#)
25. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: CVPR. pp. 2117–2125 (2017) [6](#), [7](#), [3](#), [4](#)
26. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV. pp. 2980–2988 (2017) [9](#)
27. Ma, J., Wang, B.: Segment anything in medical images. arXiv preprint arXiv:2304.12306 (2023) [5](#)
28. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: International Conference on Machine Learning (ICML). pp. 807–814 (2010) [7](#)
29. Schmied, A., Fischer, T., Danelljan, M., Pollefeys, M., Yu, F.: R3d3: Dense 3d reconstruction of dynamic scenes from multiple cameras. In: ICCV. pp. 3216–3226 (2023) [4](#), [5](#), [10](#), [11](#), [12](#)
30. Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: CVPR. pp. 4104–4113 (2016) [4](#)
31. Shi, Y., Cai, H., Ansari, A., Porikli, F.: Ega-depth: Efficient guided attention for self-supervised multi-camera depth estimation. In: CVPRW. pp. 119–129 (2023) [4](#), [9](#)
32. Talker, L., Cohen, A., Yosef, E., Dana, A., Dinerstein, M.: Mind the edge: Refining depth edges in sparsely-supervised monocular depth estimation. arXiv preprint arXiv:2212.05315 (2022) [9](#)
33. Teed, Z., Deng, J.: Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. *Advances in neural information processing systems* **34**, 16558–16569 (2021) [5](#)
34. Wang, X., Zhu, Z., Huang, G., Chi, X., Ye, Y., Chen, Z., Wang, X.: Crafting monocular cues and velocity guidance for self-supervised multi-frame depth learning. In: AAAI. pp. 2689–2697 (2023) [2](#), [4](#)
35. Wang, Y., Liang, Y., Xu, H., Jiao, S., Yu, H.: Ssqldepth: Generalizable self-supervised fine-structured monocular depth estimation. arXiv preprint arXiv:2309.00526 (2023) [10](#)
36. 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* **13**(4), 600–612 (2004) [9](#)
37. Watson, J., Mac Aodha, O., Prisacariu, V., Brostow, G., Firman, M.: The temporal opportunist: Self-supervised multi-frame monocular depth. In: CVPR. pp. 1164–1174 (2021) [1](#), [2](#), [3](#), [10](#)

38. Wei, Y., Zhao, L., Zheng, W., Zhu, Z., Rao, Y., Huang, G., Lu, J., Zhou, J.: Surrounddepth: Entangling surrounding views for self-supervised multi-camera depth estimation. In: Conference on Robot Learning (CoRL). pp. 539–549 (2022) [1](#), [2](#), [4](#), [5](#), [9](#), [10](#), [11](#), [12](#), [3](#), [7](#), [8](#)
39. Wimbauer, F., Yang, N., Von Stumberg, L., Zeller, N., Cremers, D.: Monorec: Semi-supervised dense reconstruction in dynamic environments from a single moving camera. In: CVPR. pp. 6112–6122 (2021) [3](#)
40. Wu, J., Xu, R., Wood-Doughty, Z., Wang, C.: Segment anything model is a good teacher for local feature learning. arXiv preprint arXiv:2309.16992 (2023) [5](#)
41. Yang, N., Wang, R., Stuckler, J., Cremers, D.: Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In: ECCV. pp. 817–833 (2018) [3](#)
42. Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. In: ECCV. pp. 767–783 (2018) [1](#)
43. Yu, T., Feng, R., Feng, R., Liu, J., Jin, X., Zeng, W., Chen, Z.: Inpaint anything: Segment anything meets image inpainting. arXiv preprint arXiv:2304.06790 (2023) [5](#)
44. Yuan, W., Gu, X., Dai, Z., Zhu, S., Tan, P.: New crfs: Neural window fully-connected crfs for monocular depth estimation. arXiv preprint arXiv:2203.01502 (2022) [1](#)
45. Zhang, C., Han, D., Qiao, Y., Kim, J.U., Bae, S.H., Lee, S., Hong, C.S.: Faster segment anything: Towards lightweight sam for mobile applications. arXiv preprint arXiv:2306.14289 (2023) [10](#), [14](#), [1](#)
46. Zhang, K., Liu, D.: Customized segment anything model for medical image segmentation. arXiv preprint arXiv:2304.13785 (2023) [5](#)
47. Zhang, N., Nex, F., Vosselman, G., Kerle, N.: Lite-mono: A lightweight cnn and transformer architecture for self-supervised monocular depth estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18537–18546 (2023) [3](#)