

BEYOND 2D REPRESENTATION: LEARNING 3D SCENE FIELD FOR ROBUST SELF-SUPERVISED MONOCULAR DEPTH ESTIMATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Monocular depth estimation has been extensively studied over the past few decades, yet achieving robust depth estimation in real-world scenes remains a challenge, particularly in the presence of reflections, shadow occlusions, and low-texture regions. Existing methods typically rely on extracting front-view 2D features for depth estimation, which often fail to capture those complex physical factors present in real-world scenes, leading to discontinuous, incomplete, or inconsistent depth maps. To address these issues, we turn to learning a more powerful 3D representation for robust monocular depth estimation, and propose a novel self-supervised monocular depth estimation framework based on the Three-dimensional Scene Field representation, or TSF-Depth for short. Specifically, we build our TSF-Depth framework upon an encoder-decoder architecture. The encoder extracts scene features from the input 2D image, and subsequently reshapes it as a tri-plane feature field by incorporating scene prior encoding. This tri-plane feature field is designed to implicitly model the structure and appearance of the continuous 3D scene. We then estimate a high-quality depth map from the tri-plane feature field by simulating the camera imaging process. To do this, we construct a 2D feature map with 3D geometry by sampling from the tri-plane feature field using the coordinates of points where the line of sight intersects with the scene. The aggregated multi-view geometric features are subsequently fed into the decoder for depth estimation. Extensive experiments on KITTI and NYUv2 datasets show that TSF-Depth achieves state-of-the-art performance. We also validate the generalization capability of our model on Make3D and ScanNet datasets.

1 INTRODUCTION

Monocular depth estimation is an essential computer vision task and has wide applications in autonomous driving (Geiger et al., 2013; Menze & Geiger, 2015), robot navigation (Dudek & Jenkin, 2024), and 3D reconstruction (Lyu et al., 2023; Yu et al., 2022), etc. This task aims to infer the depth of each pixel in a single image, thereby recovering the 3D scene structure. Yet, estimating depth from a single image is indeed ill-posed and inherently ambiguous, since a 3D scene can be back-projected from an infinite number of 2D images (Shao et al., 2023). Thus, the lack of sufficient 3D geometric cues in a 2D image poses a substantial challenge for monocular depth estimation.

Early monocular depth estimation (Yuan et al., 2022; Liu et al., 2023; Shao et al., 2024) worked in a supervised manner and yielded relatively accurate depth. However, depth labels are expensive to obtain. Moreover, existing physical devices can only capture sparse scene depth (Moon et al., 2023). Consequently, the sparse supervision and data scale hinder their application in practical scenarios.

Recently, self-supervised monocular depth estimation (Zhang et al., 2023a; Han et al., 2023) has attracted widespread attention. The core of such approaches is to synthesize a 2D image using the estimated depth map and minimize the photometric loss between the synthesized image and the target image (Zhao et al., 2023a). Previous efforts focused on mining effective 2D features for depth estimation by designing advanced network architectures (Lyu et al., 2021; Zhang et al., 2023a), developing more suitable loss functions (Godard et al., 2019; Liu et al., 2024), using semantic information (Casser et al., 2019), or leveraging geometric priors (Zhao et al., 2024; Sun et al., 2024).

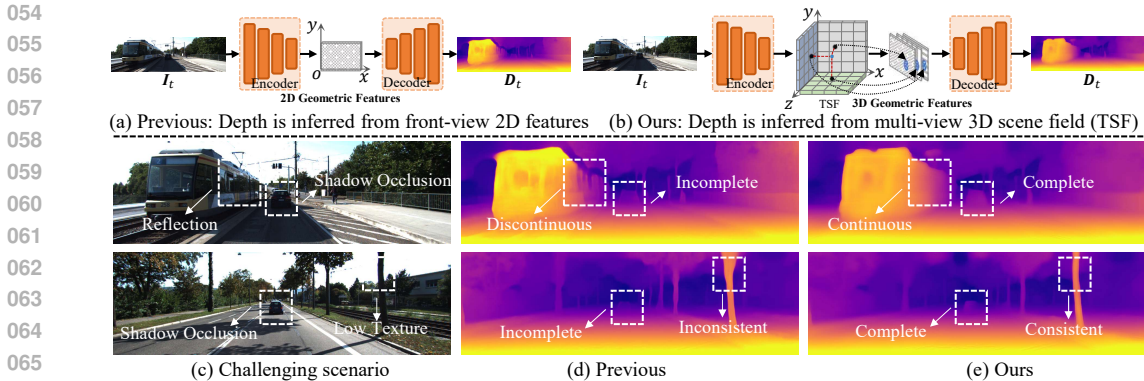


Figure 1: **An illustration of our motivation.** Real-world 3D scenes usually contain numerous challenges, including reflections, shadow occlusions, low textures and so on, which causes existing methods (*e.g.*, Lite-Mono-8M (Zhang et al., 2023a)) that typically rely on front-view 2D features to obtain discontinuous, incomplete, or inconsistent depth maps. To this end, we propose to model a multi-view 3D scene field, thereby capturing 3D geometric features for robust depth estimation.

Although these methods have shown satisfactory performance in conventional scenes, yet achieving robust depth estimation in real-world scenes remains a challenge, particularly in the presence of reflections, shadow occlusions, and low-texture regions (see Fig. 1 (c)). The main challenge posed by such scenarios is that these local regions often lack sufficient discriminative depth cues. In response, humans simulate a roughly 3D scene corresponding to the 2D image and then combine geometry cues from horizontal, vertical, and depth direction to infer the depth of a particular pixel. Even if information is lacking in one direction, geometric clues from other directions can supplement it. However, almost all existing methods typically rely on extracting front-view 2D features for depth estimation, which often fail to capture those complex physical factors present in real-world scenes. Thus, as shown in Fig. 1 (a) and (d), existing methods are limited by the paradigm of the front-view 2D representation, which often do not contain sufficient 3D geometric cues (Han et al., 2023), resulting in discontinuous, incomplete, or inconsistent depth maps.

In this paper, we propose a novel self-supervised monocular depth estimation framework based on the **Three-dimensional Scene Field** representation, or TSF-Depth for short. Unlike previous methods that employ only the front-view 2D features, we design a 3D scene field to recover the multi-view representation and then capture sufficient structure- and orientation-aware 3D geometric features from it for robust depth estimation (see Fig. 1 (b) and (e)). Specifically, we build our TSF-Depth upon an encoder-decoder architecture. The encoder extracts scene features from the input 2D image, and then are reshaped as a tri-plane feature field with three axis-aligned orthogonal feature planes by incorporating scene prior encoding. This tri-plane feature field is designed to implicitly model the structure and appearance of the continuous 3D scene. We then estimate a high-quality depth map from the tri-plane feature field by simulating the camera imaging process. To achieve this, we construct a 2D feature map with 3D geometry by sampling from the tri-plane feature field using the coordinates of the points where the line of sight intersects with the scene. The aggregated multi-view geometric feature map is then fed into the decoder for depth estimation. Extensive experiments on four datasets validate the state-of-the-art and generalization capabilities of TSF-Depth.

To summarize, the main contributions of our work are as follows:

- We propose a novel self-supervised monocular depth estimation framework based on the **Three-dimensional Scene Field** representation (TSF-Depth). To the best of our knowledge, our TSF-Depth is the first work to model 3D scene field for monocular depth estimation.
- We design a tri-plane feature field that is reshaped from hybrid features of scene content and scene prior encoding to model the multi-view representation of the continuous 3D scene.
- We attentively design a 3D-to-2D mapping strategy to sample 2D features with 3D geometry from tri-plane feature field by simulating camera image, *i.e.*, projecting the coordinates of the points where the line of sight intersects with the scene onto three orthogonal planes.
- Extensive experiments on widely used outdoor datasets (KITTI and Make3D) and indoor datasets (NYUv2 and ScanNet) show the robustness and generalization capabilities.

2 RELATED WORK

Supervised Monocular Depth Estimation. Eigen et al. (2014) first used a coarse-to-fine network for monocular depth estimation. Subsequently, numerous supervised works have been proposed. These works can be functionally classified into regression-based methods (Ranftl et al., 2021; Zhao et al., 2021; Shao et al., 2023) and classification-based methods (Bhat et al., 2021; Hu et al., 2022; Shao et al., 2024). Regression-based works use convolutional neural networks to directly learn the depth value of each pixel by minimizing the error between the prediction and ground-truth depths. However, these methods usually suffer from slow convergence and local solutions (He et al., 2022). Classification-based methods divide the depth range into different bins, and predict the probability of falling in each bin to obtain the final depth by weighted summation, which is easier to optimize. However, the high cost of data collection for training limits the wide application of these methods.

Self-Supervised Monocular Depth Estimation. Self-supervised depth estimation approaches that avoid the need for ground-truth depth during training phase have gained attention. Zhou et al. (2017) proposed a pioneering work that utilized depth network and pose network to jointly estimate depth map and camera pose, and only adopted monocular video as training data. Following this classical joint training pipeline, subsequent works improve the performance by designing robust losses (Gordon et al., 2019; Shu et al., 2020; Zhan et al., 2018), using auxiliary information during training (Watson et al., 2019; Klodt & Vedaldi, 2018), dealing with moving objects (Godard et al., 2019; Klingner et al., 2020), and adding extra geometric constraints (Yang et al., 2018; Li et al., 2021). Yet, these methods also have limitations as they infer depth from the front-view 2D feature space, which often do not contain sufficient 3D geometric cues (Han et al., 2023), and ignore the value of additional scene geometry priors in depth estimation. Instead, our TSF-Depth models 3D scene field using scene feature and scene priors to recover the multi-view representation and then capture sufficient 3D geometric features from it for robust depth estimation.

3D Scene Representation. Depth estimation using only 2D representation is a well-known ill-posed problem. To this end, learning-based multi-view stereo (MVS) methods (Yao et al., 2018; 2019; Yang et al., 2022) use the cost volume as spatial representation of the scene and then utilize 3D CNNs to continuously learn full-space 3D features. Considering the point cloud structure makes 3D feature learning more flexible compared to the cost volume, some point cloud-based MVS (Chen et al., 2019; 2020; Zhao et al., 2023b) propose to replace the cost volume with the point cloud as the spatial representation of the scene. Although these methods establish the spatial structure feature of the scene or learn the global feature of the scene, the structural attributes are not further perceived and learned. Additionally, 3D CNNs require memory cubic to the model resolution, which can be a hindrance to achieving optimal performance. Unlike existing MVS methods require complex cost volume, multi-view images and 3D CNNs, TSF-Depth models a single-view image into a tri-plane feature field as a 3D scene representation using only a 2D encoder-decoder architecture.

3 METHOD

3.1 OVERVIEW

Previous studies relied on extracting front-view 2D features for depth estimation. Although these methods are effective in conventional scenes, as discussed in Section 1, they often fail to represent those complex physical factors present in real-world scenes due to the limitations of front-view 2D representations. To this end, we propose TSF-Depth, a novel self-supervised depth estimation framework based on the 3D scene representation. The TSF-Depth is designed to model a multi-view representation of the continuous 3D scene with a tri-plane feature field. Then, we can achieve robust depth estimation using 2D features with 3D geometry sampled from the 3D scene field.

The proposed pipeline, illustrated in Fig. 2, consists of two essential steps: modeling 3D scene with tri-plane feature field and depth estimation with 3D scene field. Give a target image $\mathbf{I}_t \in \mathbb{R}^{H \times W \times 3}$, and image coordinates $\{\mathbf{C}^s \in \mathbb{R}^{H/2^s \times W/2^s \times 3}\}_{s=1}^S$ at multiple scales, an encoder is used to extract multi-scale scene features $\{\mathbf{F}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C}\}_{s=1}^S$ from the former, and positional encoding is performed to obtain multi-scale scene priors $\{\mathbf{E}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C}\}_{s=1}^S$ from the latter. The H and W denote the height and width of image, while S and C denotes the feature scale and channel. The scene priors not only initialize 3D scene structure, but also provide spatial cues. Then, the scene

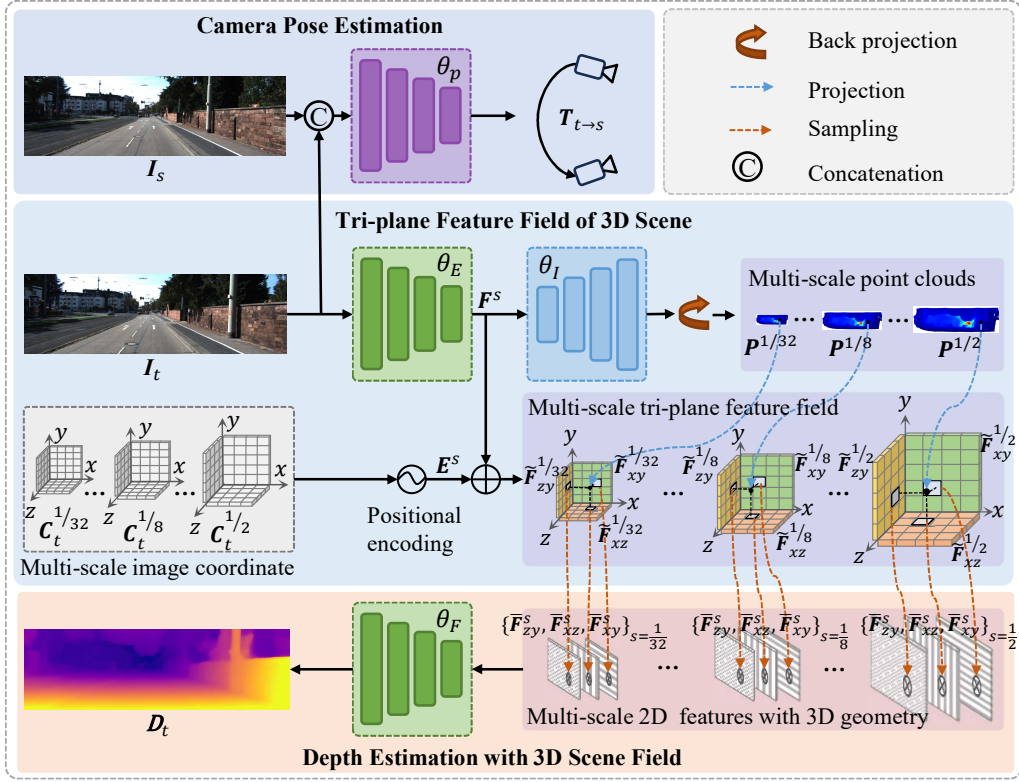


Figure 2: **Overview of the proposed TSF-Depth.** Given a target image and image coordinates at multiple scales, the scene features and scene prior features are first extracted and summed to generate multi-scale hybrid features. The hybrid features are reshaped into multi-scale tri-plane feature fields to implicitly model the multi-view representation of the 3D scene. We then recover the coordinates of the points where the line of sight intersects with the scene, and project all points onto three orthogonal planes to retrieve the 2D feature with 3D geometry for robust depth estimation.

features are split into three $\frac{C}{3}$ -channel feature maps in the channel dimension, and incorporated into scene priors to generate multi-scale hybrid features $\{\tilde{F}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C}\}_{s=1}^S$ with scene structure awareness. Subsequently, the hybrid features are reshaped as axis-aligned orthogonal tri-plane feature field $\{\tilde{F}_{xy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}, \tilde{F}_{xz}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}, \tilde{F}_{zy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}\}_{s=1}^S$, thereby implicitly modeling the multi-view representation of the continuous 3D scene. Meanwhile, we again use intermediate semantic features to recover the coordinates of the points where the line of sight from each pixel intersects with the scene. Finally, we project all 3D points P onto three orthogonal feature planes, retrieving the corresponding feature $\{\bar{F}_{xy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}, \bar{F}_{xz}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}, \bar{F}_{zy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}\}_{s=1}^S$ via bilinear interpolation, and aggregating these multi-view geometric feature map via concatenation to predict the depth map. During training phase, we follow classic self-supervised depth estimation approaches to simultaneously learning a PoseNet to predict relative pose, which will be combined with depth to construct the optimized object.

3.2 TRI-PLANE FEATURE FIELD OF 3D SCENE

Inspired by humans combine depth cues from horizontal, vertical, and depth directions to infer the depth of a particular pixel, we thus implicitly model a three-plane feature field to recover the multi-view representation. In addition, due to the ill-posed depth estimation task, relying solely on image content to infer depth is limited. Observing that the pixel depth in an image is closely related to its relative spatial position, we thus explore this regularity and incorporate it into our 3D scene Field.

Scene Feature Extracting. Given target image $I_t \in \mathbb{R}^{H \times W \times 3}$, we employ a 2D encoder to capture multi-scale scene features:

$$F^s = \theta_E(I_t), \quad (1)$$

where $\mathbf{F}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C}$. Benefiting from the powerful learning and representation capabilities of neural networks, it is feasible to learn features of 3D scenes in different directions from the input image. To this end, we split the semantic feature maps in the channel dimension and form three $\frac{C}{3}$ -channel feature maps $\{\mathbf{F}_{xy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times \frac{C}{3}}, \mathbf{F}_{xz}^s \in \mathbb{R}^{H/2^s \times W/2^s \times \frac{C}{3}}, \mathbf{F}_{yz}^s \in \mathbb{R}^{H/2^s \times W/2^s \times \frac{C}{3}}\}_s$ as the preliminary representation of the three orthogonal views of the 3D scene.

Scene Prior Encoding. In order to reasonably incorporate the scene prior, *i.e.*, relative spatial position, instead of directly passing the image coordinate into network, we introduce a positional encoding to map the image coordinate to a high-dimensional feature vector. More details about scene prior are discussed in Appendix B. Formally, the position encoding function is defined as:

$$\gamma(c) = (\sin(2^0\pi c), \cos(2^0\pi c), \dots, \sin(2^{L-1}\pi c), \cos(2^{L-1}\pi c)), \quad (2)$$

where c is the stored value of coordinate, L is the number of encoding frequencies, and $\gamma(c)$ denotes the mapping of c from \mathbb{R} into a higher dimensional space \mathbb{R}^{2L} . Thus, given the homogeneous coordinates $\{\mathbf{C}_{xy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times 3}, \mathbf{C}_{xz}^s \in \mathbb{R}^{H/2^s \times W/2^s \times 3}, \mathbf{C}_{zy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times 3}\}$ of three orthogonal views at multiple scales, the multi-scale and multi-view scene prior can be obtained by:

$$\begin{cases} \mathbf{E}_{xy}^s(u, v) = \Xi(\gamma(u'), \gamma(v'), \gamma(1)), \\ \mathbf{E}_{xz}^s(u, v) = \Xi(\gamma(u'), \gamma(1), \gamma(v')), \\ \mathbf{E}_{zy}^s(u, v) = \Xi(\gamma(1), \gamma(v'), \gamma(u')). \end{cases} \quad (3)$$

where $u' = \frac{2u}{W/2^s - 1} - 1$, $v' = \frac{2v}{W/2^s - 1} - 1$, and $\Xi[\cdot]$ is the concatenation operator. Here, u' and v' are normalized to $[-1, 1]$, which ensure scale consistency of scene prior and numerical stability.

Since the channel dimension of the generated scene prior features is controlled by L , it may not match the scene feature. We further introduce two 1×1 convolutions to adjust the channel dimension of as $\{\mathbf{E}_{xy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}, \mathbf{E}_{xz}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}, \mathbf{E}_{zy}^s \in \mathbb{R}^{H/2^s \times W/2^s \times C/3}\}_{s=1}^S$.

Multi-scale Tri-plane Feature Field. After the above two steps, we obtain the scene features that represent the semantic details of 3D scene, and the scene prior features that initialize 3D scene structure, but also provide spatial cues. Subsequently, we incorporate the scene semantic features into scene prior features to obtain multi-scale hybrid features with scene structure awareness:

$$\begin{cases} \tilde{\mathbf{F}}_{xy}^s = \mathbf{F}_{xy}^s + \mathbf{E}_{xy}^s, \\ \tilde{\mathbf{F}}_{xz}^s = \mathbf{F}_{xz}^s + \mathbf{E}_{xz}^s, \\ \tilde{\mathbf{F}}_{zy}^s = \mathbf{F}_{zy}^s + \mathbf{E}_{zy}^s. \end{cases} \quad (4)$$

The three hybrid feature planes are axis-aligned orthogonal planes, which are defined as our tri-plane feature fields. The $\tilde{\mathbf{F}}_{xy}^s$ perceives the continuous change of depth in z-direction, the $\tilde{\mathbf{F}}_{xz}^s$ perceives the consistency of depth in the vertical direction and the $\tilde{\mathbf{F}}_{zy}^s$ perceives the similarity of depth in the horizontal direction. Therefore, this tri-plane feature field is designed to implicitly model the structure, orientation, appearance of the continuous 3D scene. Moreover, benefiting from the design of our multi-scale and multi-view tri-plane feature fields, our method can effectively perceive the 3D scene scale, thereby alleviating the ambiguity of monocular depth estimation.

3.3 DEPTH ESTIMATION WITH 3D SCENE FIELD

After model the 3D scene filed with tri-plane feature field, we aim to capture the 2D feature with 3D geometric from it for depth estimation. However, for each pixel in the image space, we do not know its corresponding position in the 3D scene field. To address this issue, we simulating the inverse process of camera imaging, *i.e.*, we need to recover the coordinates of the points where the line of sight intersects the scene. In other words, we need to estimate the point clouds of the input image.

Point clouds are often used as structural representations of 3D scenes. Many works use 3D CNNs to directly estimate 3D point clouds, but they are costly in terms of efficiency and storage. In this work, we employ the depth map as intermediate representation to efficiently reconstruct the point clouds. To this end, we use a decoder θ_I to predict the initial depth map $\mathbf{D}_{I,t}^s = \theta_I(\Xi[\mathbf{F}_{xy}^s, \mathbf{F}_{xz}^s, \mathbf{F}_{zy}^s])$.

After the above step, we can obtain multi-scales depth maps $\mathbf{D}_{I,t}^s \in \mathbb{R}^{\frac{H}{2^s} \times \frac{W}{2^s}}$. Since the scale and shift coefficients of the predicted depth maps are unknown, the reconstructed 3D structure from them is likely to be distorted from inappropriate affine changes. We thus use all scales of depth map to reconstruct multi-scale point clouds from coarse to fine. Given the estimated initial depth map $\mathbf{D}_{I,t}^s$ with scale s and camera intrinsics, we can reconstruct the point cloud from the pixel coordinate

based on the pinhole camera model. Specifically, for a 2D point $\mathbf{p}_i = [u, v]^T$ in the pixel coordinate system, it can be reprojected back to a 3D point $\mathbf{P}_i = [X, Y, Z]^T$ in the camera system by:

$$[X, Y, Z]^T = \mathbf{D}_{I,t}^s(u, v) \mathbf{K}^{-1}[u, v, 1]^T, \quad (5)$$

where $\mathbf{K} \in \mathbb{R}^{3 \times 3}$ denotes the camera intrinsic. Therefore, using Eq. 5, we can explicitly reconstruct multi-scale point clouds $\{\mathbf{P}^s \in \mathbb{R}^{\frac{H}{2^s} \times \frac{W}{2^s} \times 3}\}_s^S$. It should be noted that for different scale depths, the coordinate range will also change accordingly, *i.e.*, $u \in [0, \frac{W}{2^s} - 1]$, $v \in [0, \frac{H}{2^s} - 1]$.

After the implicitly and explicitly modeling 3D scene phase, we obtain the dense tri-plane feature field and sparse point clouds, respectively. Although the former learns from the 2D image, they have the ability to perceive the 3D scene structure and provide direction-aware multi-view features. The later learn from 2D feature space and lacks of 3D scene awareness, but they can provide approximate spatial structure. To obtain the final precise and robust depth prediction, we design a 3D-to-2D mapping strategy to sample 2D features with 3D geometry by combining the advantage of both.

Since there are different spatial scales between each 3D point in the point cloud and each orthogonal planes of tri-plane feature fields, they need to be aligned in the same space. Compare to back-projecting the each orthogonal planes of tri-plane feature field from the 2D space to 3D space, projecting point cloud from 3D space to 2D space is more efficient. With the multi-scale point cloud and tri-plane feature field, we perform an orthographic projection onto the three axis-aligned planes:

$$[x, y, z]^T = \mathbf{K}[X, Y, Z]^T. \quad (6)$$

Generally, the value ranges of point $[x, y, z]$ in three axis are different, *i.e.*, $x \in [0, W - 1]$, $y \in [0, H - 1]$ and $z \in [0, M]$, where M is the maximum depth. To ensure that the projected 2D points are aligned with each planes without going out of bounds, they need to be normalized to $[-1, 1]$:

$$\begin{cases} \bar{x} = (x/(W-1) - 0.5) \times 2, \\ \bar{y} = (y/(H-1) - 0.5) \times 2, \\ \bar{z} = (z/M - 0.5) \times 2. \end{cases} \quad (7)$$

Ultimately, we obtain the 2D points of the point cloud projected onto each orthogonal plane, *i.e.*, the $p_{xy}^s = [\bar{x}, \bar{y}]$ is the projected 2D point located at $\tilde{\mathbf{F}}_{xy}^s$, the $p_{xz}^s = [\bar{x}, \bar{z}]$ is the projected 2D point located at $\tilde{\mathbf{F}}_{xz}^s$ and the $p_{zy}^s = [\bar{z}, \bar{y}]$ is the projected 2D point located at $\tilde{\mathbf{F}}_{zy}^s$. Then, we sample 2D feature map with 3D geometric from each orthogonal plane of tri-planes feature fields using differentiable bilinear sampling operator:

$$\begin{cases} \bar{\mathbf{F}}_{xy}^s = \tilde{\mathbf{F}}_{xy}^s \langle p_{xy}^s \rangle, \\ \bar{\mathbf{F}}_{xz}^s = \tilde{\mathbf{F}}_{xz}^s \langle p_{xz}^s \rangle, \\ \bar{\mathbf{F}}_{zy}^s = \tilde{\mathbf{F}}_{zy}^s \langle p_{zy}^s \rangle. \end{cases} \quad (8)$$

where $\langle \cdot \rangle$ is the sampling operator (Jaderberg et al., 2015). Finally, the multi-view geometric features aggregated through concatenation, and a depth decoder is used to predict the high-quality depth map:

$$\mathbf{D}_{F,t} = \theta_F(\Xi[\bar{\mathbf{F}}_{xy}, \bar{\mathbf{F}}_{xz}, \bar{\mathbf{F}}_{zy}]). \quad (9)$$

3.4 SELF-SUPERVISED LEARNING

Following monodepth2 (Godard et al., 2019), we use a target frame \mathbf{I}_t and two adjacent frames \mathbf{I}_a ($a \in \{t-1, t+1\}$) to jointly train a DepthNet ($\theta_E, \theta_I, \theta_F$) and a PoseNet θ_p . During training, \mathbf{I}_t is fed into the DepthNet to get the depth $\mathbf{D}_{F,t}$, and $(\mathbf{I}_t, \mathbf{I}_a)$ are put into PoseNet to get the relative camera pose $\mathbf{T}_{t \rightarrow a}$. Then, we can synthesize a target frame $\mathbf{I}_{F,a \rightarrow t}$ by warping the source frame \mathbf{I}_t : $\mathbf{I}_{F,a \rightarrow t} = \mathbf{I}_a \langle \text{proj}(\mathbf{D}_{F,t}, \mathbf{T}_{t \rightarrow a}, \mathbf{K}) \rangle$, where $\text{proj}(\cdot)$ is the coordinate projection operation (Zhou et al., 2017). The photometric error between $\mathbf{I}_{F,t}$ and \mathbf{I}_t , consisting of L_1 and $SSIM$ weighted by α , is calculated as: $pe(\mathbf{I}_t, \mathbf{I}_{F,a \rightarrow t}) = \frac{\alpha}{2}(1 - SSIM(\mathbf{I}_t, \mathbf{I}_{F,a \rightarrow t})) + (1 - \alpha)\|\mathbf{I}_t - \mathbf{I}_{F,a \rightarrow t}\|_1$. Following Monodepth2, we adopt the per-pixel minimum photometric loss as our reprojection loss:

$$\mathcal{L}_{ph}(\mathbf{I}_t, \mathbf{I}_{F,a \rightarrow t}) = \min_a pe(\mathbf{I}_t, \mathbf{I}_{F,a \rightarrow t}). \quad (10)$$

We also use the edge-aware smoothness loss to encourage locally smooth depth maps:

$$\mathcal{L}_{sm}(\mathbf{D}_{F,t}, \mathbf{I}_t) = |\partial_x \mathbf{D}_{F,t}^*| e^{-|\partial_x \mathbf{I}_t|} + |\partial_y \mathbf{D}_{F,t}^*| e^{-|\partial_y \mathbf{I}_t|}. \quad (11)$$

where ∂_x and ∂_y are the gradients in the horizontal and vertical direction respectively. Besides, $\mathbf{D}_{F,t}^* = \mathbf{D}_{F,t} / \bar{\mathbf{D}}_{F,t}$ is the mean-normalized inverse depth from Monodepth2 to discourage shrinking of the estimated depth. Therefore, the total loss for final depth map is defined as:

$$\mathcal{L}_F = \beta \mathcal{L}_{ph}(\mathbf{I}_t, \mathbf{I}_{F,a \rightarrow t}) + \gamma \mathcal{L}_{sm}(\mathbf{D}_{F,t}, \mathbf{I}_t). \quad (12)$$

Table 1: **Depth estimation results on KITTI** (Geiger et al., 2013). We divide compared methods into three categories. S: stereo training. M: monocular training. MS: stereo and monocular training.

Method	Train	Error Metric (\downarrow)				Accuracy Metric (\uparrow)		
		Sq Rel	Abs Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Monodepth (Godard et al., 2017)	S	1.344	0.148	5.927	0.247	0.803	0.922	0.964
3Net (Poggi et al., 2018)	S	1.201	0.119	5.888	0.208	0.844	0.941	0.978
Monodepth2 (Godard et al., 2019)	S	0.873	0.109	4.960	0.208	0.864	0.948	0.975
DepthHints (Watson et al., 2019)	S	0.780	0.106	4.695	0.193	0.875	0.958	0.980
BRNet (Han et al., 2022)	S	0.792	0.103	4.716	0.197	0.876	0.954	0.978
Monodepth2 (Godard et al., 2019)	MS	0.818	0.106	4.750	0.196	0.874	0.957	0.979
DepthHints (Watson et al., 2019)	MS	0.769	0.105	4.627	0.189	0.875	0.959	0.982
HR-Depth (Lyu et al., 2021)	MS	0.785	0.107	4.612	0.185	0.887	0.962	0.982
R-MSFM6 (Zhou et al., 2021b)	MS	0.787	0.111	4.625	0.189	0.882	0.961	0.981
GeoNet (Yin & Shi, 2018)	M	1.060	0.149	5.567	0.226	0.796	0.935	0.975
Monodepth2 (Godard et al., 2019)	M	0.903	0.115	4.863	0.193	0.877	0.959	0.971
DepthHints (Watson et al., 2019)	M	0.845	0.109	4.800	0.196	0.870	0.956	0.980
S ³ Net (Cheng et al., 2020)	M	0.826	0.124	4.981	0.200	0.846	0.955	0.982
HR-Depth (Lyu et al., 2021)	M	0.792	0.109	4.632	0.185	0.884	0.962	0.983
CADepth-Net (Yan et al., 2021)	M	0.769	0.105	4.535	0.181	0.892	0.964	0.983
DIFFNet (Zhou et al., 2021a)	M	0.764	0.102	4.483	0.180	0.896	0.965	0.983
DynaDepth (Zhang et al., 2022a)	M	0.761	0.108	4.608	0.187	0.883	0.962	0.982
MonoFormer (Bae et al., 2023)	M	0.846	0.104	4.580	0.183	0.891	0.962	0.982
SC-DepthV3 (Sun et al., 2023)	M	0.756	0.118	4.709	0.188	0.864	0.960	0.984
Zhang et al. (Zhang et al., 2023b)	M	0.786	0.105	4.572	0.182	0.890	0.964	0.983
Lite-Mono (Zhang et al., 2023a)	M	0.765	0.107	4.561	0.183	0.886	0.963	0.983
Lite-Mono-8M (Zhang et al., 2023a)	M	0.729	0.101	4.454	0.178	0.897	0.965	0.983
Zhao et al. Zhao et al. (2024)	M	0.809	0.110	4.616	0.185	-	-	-
Xiong et al. (Xiong et al., 2024)	M	0.868	0.122	4.986	0.200	0.857	0.953	0.980
ShuffleMono (Feng et al., 2024)	M	0.850	0.114	4.821	0.193	0.872	0.957	0.980
Liu et al. (Liu et al., 2024)	M	0.747	0.114	4.724	0.187	0.863	0.960	0.984
Dynamo-Depth (Sun et al., 2024)	M	0.758	0.112	4.505	0.183	0.873	0.959	0.984
TSF-Depth	M	0.692	0.096	4.335	0.173	0.903	0.967	0.984

To reconstruct accuracy point cloud, we apply a photometric loss to the initial depth:

$$\mathcal{L}_I = \frac{1}{S} \sum_{s=1}^S \varphi \mathcal{L}_{ph}(\mathbf{D}_{F,t}, \mathbb{U}(\mathbf{D}_{I,t}^s)). \quad (13)$$

where \mathbb{U} is the upsampling operation. Our overall loss function can be formulated as:

$$\mathcal{L} = \mathcal{L}_I + \mathcal{L}_F. \quad (14)$$

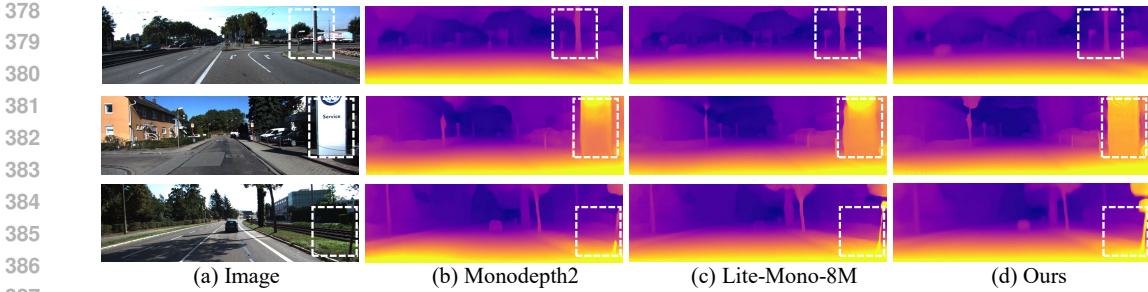
4 EXPERIMENTS

4.1 DATASETS AND EVALUATION METRICS

Outdoor Datasets. **KITTI** (Geiger et al., 2013) is an outdoor benchmark with a resolution of 1242×375 . Following Zhou et al. (2017), we use 39810, 4424, and 697 images for training, validation, and testing, respectively. **Make3D** (Saxena et al., 2008) contains 134 outdoor image-depth pairs with a resolution of 1704×2272 , which is used in this work as a generalization test.

Indoor Datasets. **NYUv2** (Silberman et al., 2012) is an indoor benchmark with a resolution of 640×480 . Following StructDepth (Li et al., 2021), we use the official training and validation splits which include 302 and 33 sequences, and use officially provided 654 images for testing. **ScanNet** (Dai et al., 2017) contains 1513 indoor scenes with a resolution of 1296×968 . We use official data split and use this dataset as a generalization test for indoor scenes.

Evaluation Metrics. We follow Monodepth2 (Godard et al., 2019) using relative squared error (**Sq Rel**), absolute relative error (**Abs Rel**), root mean squared error (**RMSE**), root mean squared logarithmic error (**RMSE log**) and threshold accuracy ($\sigma < 1.25$, $\sigma < 1.25^2$ and $\sigma < 1.25^3$).

Figure 3: **Qualitative comparison on KITTI** (Geiger et al., 2013). We highlight challenging areas.Table 2: **Generalization on Make3D** (Saxena et al., 2008). All methods are trained on KITTI.

Method	Error Metric (\downarrow)			
	Sq Rel	Abs Rel	RMSE	RMSE log
Monodepth2 (Godard et al., 2019)	3.589	0.322	7.418	0.163
HR-Depth (Lyu et al., 2021)	3.208	0.315	7.024	0.159
DIFFNet (Zhou et al., 2021a)	3.313	0.309	7.008	0.155
DynaDepth (Zhang et al., 2022a)	3.311	0.334	7.463	0.169
Chen et al. (Chen et al., 2023)	3.610	0.370	7.133	-
Zhang et al. (Zhang et al., 2023b)	3.485	0.314	7.188	-
Lite-Mono (Zhang et al., 2023a)	3.060	0.305	6.981	0.158
Zhao et al. (Zhao et al., 2024)	3.200	0.316	7.095	0.158
Xiong et al. (Xiong et al., 2024)	3.102	0.319	7.005	0.161
TSF-Depth	2.925	0.292	6.744	0.150

4.2 IMPLEMENTATION DETAILS

We implement TSF-Depth in PyTorch (Paszke et al., 2017), training it for 20 epochs on the outdoor dataset and 40 epochs on the indoor dataset by using AdamW (Loshchilov & Hutter, 2017) optimizer on a single RTX 3090 GPU. The batch size is set to 12 for the outdoor dataset and 16 for the indoor dataset. The initial learning rate for PoseNet and depth decoder is 1×10^{-4} , while the depth encoder is trained with an initial learning rate of 5×10^{-5} . For the PoseNet, we use the same architecture as Monodepth2 (Godard et al., 2019). For the encoder θ_E , initial depth decoder θ_I and final depth decoder θ_F of DepthNet, we choose mpvit, the decoder of Monodepth2 (Godard et al., 2019) and HRDecoder (He et al., 2022), respectively. The hyper-parameters S , L , α , β , γ , and φ are set to 5, 10, 0.85, 1.0, 0.001, and 0.5 respectively. More implementation details are reported in Appendix C.

4.3 COMPARISON ON OUTDOOR SCENE

KITTI. Table 1 presents the quantitative comparison at resolution of 640×192 on the outdoor benchmark, *i.e.*, KITTI dataset (Geiger et al., 2013). Compared to existing methods trained on monocular videos, our method outperforms all these works by significant margins, and also outperforms counterparts trained with additional stereo pairs. In particular, our method relatively outperforms the SOTA method Lite-Mono-8M by 5.1% and by 5.0% in terms of Sq Rel and Abs Rel, respectively. We also compare the qualitative performance with the classic work Monodepth2 (Godard et al., 2019) and the SOTA work Lite-Mono-8M (Zhang et al., 2023a). Fig. 3 presents three visual samples and highlights challenging areas. We observed that for the first two examples, the traditional CNN-based MonoDepth2 and the attention-based Lite-Mono-8M, which extract only front-view 2D feature for depth estimation, both obtain inconsistent depth map. For the third challenging example with low texture, the compared methods obtained discontinuous results. In contrast, our method generates superior visual results due to our multi-view representation to capture the 3D geometric cues for robust depth estimation. More visualization results are shown in Appendix D. In addition, the quantitative results at 1280×384 and 1024×320 resolutions are reported in Appendix E.

Make3D. To show the generalization capability, we further test our proposed method on the Make3D dataset (Saxena et al., 2008). Following the evaluation strategy in (Zhou et al., 2017), our model is training on the KITTI dataset (Geiger et al., 2013) without any fine-tuning on the Make3D dataset. As shown in Table 2, our method outperforms all other self-supervised monocular depth estimation approaches, which demonstrates our models can be well generalized to unseen outdoor scenes.

Table 3: Depth estimation results on NYUv2 (Silberman et al., 2012).

Method	Error Metric (\downarrow)		Accuracy Metric (\uparrow)		
	Abs Rel	RMSE	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
MovingIndoor (Zhou et al., 2019)	0.208	0.712	0.674	0.900	0.968
Monodepth2 (Godard et al., 2019)	0.160	0.601	0.767	0.949	0.988
TrainFlow (Zhao et al., 2020)	0.189	0.686	0.701	0.912	0.978
P ² Net (Yu et al., 2020)	0.159	0.599	0.772	0.942	0.984
Bian et al. (Bian et al., 2020)	0.147	0.536	0.804	0.950	0.986
SC-DepthV1 (Bian et al., 2021)	0.159	0.639	0.734	0.937	0.983
PLNet (Jiang et al., 2021)	0.151	0.562	0.790	0.953	0.989
StructDepth (Li et al., 2021)	0.142	0.540	0.813	0.954	0.988
Zhang et al. (Zhang et al., 2022b)	0.177	0.634	0.733	0.936	-
ADPDepth (Song et al., 2023)	0.165	0.592	0.753	0.934	0.981
F ² Depth (Guo et al., 2024a)	0.153	0.569	0.787	0.950	0.987
Guo et al. (Guo et al., 2024b)	0.152	0.567	0.792	0.950	0.988
TSF-Depth	0.129	0.527	0.846	0.966	0.991

Table 4: Generalization results on ScanNet (Dai et al., 2017). All methods are trained on NYUv2.

Method	Error Metric (\downarrow)		Accuracy Metric (\uparrow)		
	Abs Rel	RMSE	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
MovingIndoor (Zhou et al., 2019)	0.212	0.483	0.650	0.905	0.976
Monodepth2 (Godard et al., 2019)	0.200	0.458	0.672	0.922	0.981
TrainFlow (Zhao et al., 2020)	0.179	0.415	0.726	0.927	0.980
P ² Net (Yu et al., 2020)	0.175	0.420	0.740	0.932	0.982
PLNet (Jiang et al., 2021)	0.176	0.414	0.735	0.939	0.985
IFMNet (Wei et al., 2021)	0.170	0.402	0.758	0.940	0.989
SC-Depthv1 (Bian et al., 2021)	0.169	0.392	0.749	0.938	0.983
StructDepth (Li et al., 2021)	0.165	0.400	0.754	0.939	0.985
TSF-Depth	0.157	0.390	0.775	0.954	0.988

4.4 COMPARISON ON INDOOR SCENE

NYUv2. Table 3 presents a performance comparison between our approach and state-of-the-art methods on the NYUv2 dataset (Silberman et al., 2012). For smaller indoor scenes, TSF-Depth significantly outperforms all previous self-supervised methods compared to larger outdoor scenes. This shows that building 3D scene fields is effective and can be easily done in small scenes.

ScanNet. We further validate the generalization ability in indoor scenes. All methods are trained on NYUv2 (Silberman et al., 2012) and tested on ScanNet (Dai et al., 2017). The results in Table 4 demonstrates that TSF-Depth has excellent generalization ability for unseen indoor scene.

4.5 ABLATION STUDY

To investigate the main contributions and key designs of TSF-Depth, a series of ablation experiments on the KITTI (Geiger et al., 2013) dataset are conducted. The pipeline of the baseline is reported in Fig. 1 (a). In addition, the ablation studies for indoor scene are presented in Appendix F.

Effects of Multi-scale Scene Priors. We first analyze the impact of incorporating multi-scale scene prior into model by positional encoding. See Table 5 (a), (b), (d) and (e), training either with single-scale or full-scale scene prior encoding, all significantly improves the depth accuracy over the baseline without it, and the combination it with tri-plane feature scenes at different scales yields an additional improvement. Based on the above analysis, the scene prior we introduce is effective.

Effects of Multi-scale Tri-plane Feature Fields. We further analyze the impact of modeling 3D scene field using multi-scale tri-plane feature fields. The results are shown in Table 5 (c) and (e). When building a single-scale tri-plane feature field, the depth accuracy is improved. Better performance is achieved by using multi-scale tri-plane feature field or incorporating scene priors, suggesting that modeling multi-view representation is effective for robust depth estimation.

Effects of 3D Scene Field on 3D Geometric Representation. We finally evaluate the effectiveness of 3D scene fields for 3D geometric representation. The results is reported in Table 6. We remove the decoder θ_I branch and positional encoding branch of TSF-Depth as a baseline and then train it.

Table 5: Ablation results for each component of our method on KITTI (Geiger et al., 2013). SP^i : Incorporate scene prior encoding with resolution $\frac{H}{2^i} \times \frac{W}{2^i}$. TP^i : Model the 3D scene using tri-plane feature field with resolution $\frac{H}{2^i} \times \frac{W}{2^i}$. SP^{All}/TP^{All} : Use all resolution SP/TP .

Exp Setting	SP TP	Error Metric (\downarrow)				Accuracy Metric (\uparrow)		
		Sq Rel	Abs Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
(a) Baseline		0.746	0.102	4.464	0.176	0.897	0.965	0.983
SP^1	✓	0.733	0.098	4.388	0.174	0.900	0.967	0.984
SP^2	✓	0.729	0.098	4.385	0.174	0.902	0.967	0.984
(b) SP^3	✓	0.722	0.098	4.386	0.175	0.901	0.967	0.984
SP^4	✓	0.745	0.098	4.419	0.174	0.899	0.967	0.984
SP^5	✓	0.737	0.100	4.384	0.175	0.901	0.967	0.984
SP^{All}	✓	0.736	0.099	4.385	0.175	0.901	0.967	0.984
TP^1	✓	0.734	0.098	4.385	0.174	0.899	0.967	0.984
TP^2	✓	0.738	0.098	4.387	0.174	0.901	0.967	0.984
(c) TP^3	✓	0.719	0.098	4.395	0.174	0.901	0.967	0.984
TP^4	✓	0.743	0.099	4.393	0.175	0.900	0.967	0.984
TP^5	✓	0.760	0.099	4.425	0.175	0.898	0.967	0.984
TP^{All}	✓	0.742	0.099	4.384	0.175	0.901	0.967	0.984
$SP^{All} + TP^1$	✓ ✓	0.722	0.099	4.387	0.174	0.901	0.967	0.984
$SP^{All} + TP^2$	✓ ✓	0.719	0.098	4.385	0.174	0.902	0.967	0.984
(d) $SP^{All} + TP^3$	✓ ✓	0.750	0.099	4.427	0.175	0.902	0.967	0.984
$SP^{All} + TP^4$	✓ ✓	0.734	0.099	4.385	0.175	0.902	0.966	0.984
$SP^{All} + TP^5$	✓ ✓	0.748	0.100	4.406	0.175	0.902	0.967	0.984
$TP^{All} + SP^1$	✓ ✓	0.734	0.099	4.385	0.174	0.902	0.967	0.984
$TP^{All} + SP^2$	✓ ✓	0.752	0.099	4.397	0.175	0.902	0.967	0.984
(e) $TP^{All} + SP^3$	✓ ✓	0.740	0.099	4.435	0.175	0.901	0.966	0.984
$TP^{All} + SP^4$	✓ ✓	0.720	0.098	4.383	0.174	0.902	0.967	0.984
$TP^{All} + SP^5$	✓ ✓	0.733	0.099	4.425	0.175	0.902	0.967	0.984
(f) TSF-Depth ($SP^{All} + TP^{All}$)	✓ ✓	0.692	0.096	4.335	0.173	0.903	0.967	0.984

Table 6: Ablation results of the ability of 3D scene fields to represent 3D geometry on KITTI (Geiger et al., 2013). θ_E^{PT} is an encoder that is frozen after being pre-trained in the baseline.

Setting	θ_E	θ_E^{PT}	θ_I	θ_F	Error Metric (\downarrow)				Accuracy Metric (\uparrow)		
					Sq Rel	Abs Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Baseline (2D geometry)	✓		✓		0.746	0.102	4.464	0.176	0.897	0.965	0.983
TSF-Depth (2D geometry)	✓	✓	✓	✓	0.755	0.103	4.475	0.181	0.892	0.963	0.982
TSF-Depth (3D geometry)	✓		✓	✓	0.692	0.096	4.335	0.173	0.903	0.967	0.984

Thus, the baseline relies on 2D features for depth estimation. Note that we denote the encoder of baseline as θ_E^{PT} . Then, we replace the θ_E of TSF-Depth with the trained encoder θ_E^{PT} , and then train other modules. In this setup, the Tri-plane feature field is constructed using 2D representation. As for the last row, it uses our complete training. Compared to the baseline and the complete TSF-Depth, we observe that the results of TSF-Depth based on 2D geometry features are worse than both them, due to the 3D-to-2D mapping requiring 3D geometric representation rather than 2D representation. Thus, our TSF-Depth can learn 3D geometric representations for robust depth estimation.

5 CONCLUSION

In this paper, we propose a novel self-supervised monocular depth estimation framework based on the Three-dimensional Scene Field representation (TSF-Depth). Unlike previous methods typically rely on extracting front-view 2D features, we turn to learning a more powerful 3D representation for robust depth estimation. We design a tri-plane feature field that is reshaped from hybrid features of scene content and scene prior to implicitly model the multi-view representation of the continuous 3D scene. Then, we attentively design a 3D-to-2D mapping strategy to sample 2D features with 3D geometry for depth estimation. Extensive experiments on widely-used outdoor datasets (KITTI and Make3D) and indoor datasets (NYUv2 and ScanNet) show the robustness and generalization ability.

REFERENCES

- 540
541
542 Jinwoo Bae, Sungho Moon, and Sunghoon Im. Deep digging into the generalization of self-
543 supervised monocular depth estimation. In *Proceedings of the AAAI conference on artificial*
544 *intelligence*, volume 37, pp. 187–196, 2023.
- 545 Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adap-
546 tive bins. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-*
547 *tion*, pp. 4009–4018, 2021.
- 548 Jia-Wang Bian, Huangying Zhan, Naiyan Wang, Tat-Jun Chin, Chunhua Shen, and Ian Reid. Unsu-
549 pervised depth learning in challenging indoor video: Weak rectification to rescue. *arXiv preprint*
550 *arXiv:2006.02708*, 2(5):7, 2020.
- 551 Jia-Wang Bian, Huangying Zhan, Naiyan Wang, Zhichao Li, Le Zhang, Chunhua Shen, Ming-Ming
552 Cheng, and Ian Reid. Unsupervised scale-consistent depth learning from video. *International*
553 *Journal of Computer Vision (IJCV)*, 2021.
- 554 Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Unsupervised monocular
555 depth and ego-motion learning with structure and semantics. In *Proceedings of the IEEE/CVF*
556 *Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0, 2019.
- 557 Rui Chen, Songfang Han, Jing Xu, and Hao Su. Point-based multi-view stereo network. In *Pro-*
558 *ceedings of the IEEE/CVF international conference on computer vision*, pp. 1538–1547, 2019.
- 559 Rui Chen, Songfang Han, Jing Xu, and Hao Su. Visibility-aware point-based multi-view stereo
560 network. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3695–3708,
561 2020.
- 562 Xingyu Chen, Thomas H Li, Ruonan Zhang, and Ge Li. Frequency-aware self-supervised monoc-
563 ular depth estimation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of*
564 *Computer Vision*, pp. 5808–5817, 2023.
- 565 Bin Cheng, Inderjot Singh Saggu, Raunak Shah, Gaurav Bansal, and Dinesh Bharadia. S 3 net:
566 Semantic-aware self-supervised depth estimation with monocular videos and synthetic data. In
567 *European Conference on Computer Vision*, pp. 52–69. Springer, 2020.
- 568 Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias
569 Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the*
570 *IEEE conference on computer vision and pattern recognition*, pp. 5828–5839, 2017.
- 571 Gregory Dudek and Michael Jenkin. *Computational principles of mobile robotics*. Cambridge
572 university press, 2024.
- 573 David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using
574 a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014.
- 575 Yingwei Feng, Zhiyong Hong, Liping Xiong, Zhiqiang Zeng, and Jingmin Li. Shufflemono: Re-
576 thinking lightweight network for self-supervised monocular depth estimation. *Journal of Artificial*
577 *Intelligence and Soft Computing Research*, 14(3):191–205, 2024.
- 578 Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The
579 kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- 580 Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estima-
581 tion with left-right consistency. In *Proceedings of the IEEE conference on computer vision and*
582 *pattern recognition*, pp. 270–279, 2017.
- 583 Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-
584 supervised monocular depth estimation. In *Proceedings of the IEEE/CVF international confer-*
585 *ence on computer vision*, pp. 3828–3838, 2019.
- 586 Ariel Gordon, Hanhan Li, Rico Jonschkowski, and Anelia Angelova. Depth from videos in the
587 wild: Unsupervised monocular depth learning from unknown cameras. In *Proceedings of the*
588 *IEEE/CVF International Conference on Computer Vision*, pp. 8977–8986, 2019.
- 589
590
591
592
593

- 594 Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raventos, and Adrien Gaidon. 3d packing
595 for self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF conference on*
596 *computer vision and pattern recognition*, pp. 2485–2494, 2020.
- 597 Xiaotong Guo, Huijie Zhao, Shuwei Shao, Xudong Li, and Baochang Zhang. F²depth: Self-
598 supervised indoor monocular depth estimation via optical flow consistency and feature map syn-
599 thesis. *Engineering Applications of Artificial Intelligence*, 133:108391, 2024a.
- 600 Xiaotong Guo, Huijie Zhao, Shuwei Shao, Xudong Li, Baochang Zhang, and Na Li. Sim-
601 multidepth: Self-supervised indoor monocular multi-frame depth estimation based on texture-
602 aware masking. *Remote Sensing*, 16(12):2221, 2024b.
- 603 Wencheng Han, Junbo Yin, Xiaogang Jin, Xiangdong Dai, and Jianbing Shen. Brnet: Exploring
604 comprehensive features for monocular depth estimation. In *European Conference on Computer*
605 *Vision*, pp. 586–602. Springer, 2022.
- 606 Wencheng Han, Junbo Yin, and Jianbing Shen. Self-supervised monocular depth estimation by
607 direction-aware cumulative convolution network. In *Proceedings of the IEEE/CVF International*
608 *Conference on Computer Vision*, pp. 8613–8623, 2023.
- 609 Mu He, Le Hui, Yikai Bian, Jian Ren, Jin Xie, and Jian Yang. Ra-depth: Resolution adaptive
610 self-supervised monocular depth estimation. In *European Conference on Computer Vision*, pp.
611 565–581. Springer, 2022.
- 612 Dongting Hu, Lihua Peng, Tingjin Chu, Xiaoxing Zhang, Yinian Mao, Howard Bondell, and Ming-
613 ming Gong. Uncertainty quantification in depth estimation via constrained ordinal regression. In
614 *European Conference on Computer Vision*, pp. 237–256. Springer, 2022.
- 615 Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances*
616 *in neural information processing systems*, 28, 2015.
- 617 Hualie Jiang, Laiyan Ding, Junjie Hu, and Rui Huang. Plnet: Plane and line priors for unsupervised
618 indoor depth estimation. In *2021 International Conference on 3D Vision (3DV)*, pp. 741–750.
619 IEEE, 2021.
- 620 Marvin Klingner, Jan-Aike Termöhlen, Jonas Mikolajczyk, and Tim Fingscheidt. Self-supervised
621 monocular depth estimation: Solving the dynamic object problem by semantic guidance. In
622 *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020,*
623 *Proceedings, Part XX 16*, pp. 582–600. Springer, 2020.
- 624 Maria Klodt and Andrea Vedaldi. Supervising the new with the old: learning sfm from sfm. In
625 *Proceedings of the European conference on computer vision (ECCV)*, pp. 698–713, 2018.
- 626 Boying Li, Yuan Huang, Zeyu Liu, Danping Zou, and Wenxian Yu. Structdepth: Leveraging
627 the structural regularities for self-supervised indoor depth estimation. In *Proceedings of the*
628 *IEEE/CVF International Conference on Computer Vision*, pp. 12663–12673, 2021.
- 629 Ce Liu, Suryansh Kumar, Shuhang Gu, Radu Timofte, and Luc Van Gool. Va-depthnet: A variational
630 approach to single image depth prediction. *arXiv preprint arXiv:2302.06556*, 2023.
- 631 Runze Liu, Dongchen Zhu, Guanghui Zhang, Yue Xu, Wenjun Shi, Xiaolin Zhang, Lei Wang, and
632 Jiamao Li. Unsupervised monocular depth estimation based on hierarchical feature-guided diffu-
633 sion. *arXiv preprint arXiv:2406.09782*, 2024.
- 634 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
635 *arXiv:1711.05101*, 2017.
- 636 Xiaoyang Lyu, Liang Liu, Mengmeng Wang, Xin Kong, Lina Liu, Yong Liu, Xinxin Chen, and
637 Yi Yuan. Hr-depth: High resolution self-supervised monocular depth estimation. In *Proceedings*
638 *of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 2294–2301, 2021.
- 639 Xiaoyang Lyu, Peng Dai, Zizhang Li, Dongyu Yan, Yi Lin, Yifan Peng, and Xiaojuan Qi. Learning
640 a room with the occ-sdf hybrid: Signed distance function mingled with occupancy aids scene
641 representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
642 pp. 8940–8950, 2023.

- 648 Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of*
649 *the IEEE conference on computer vision and pattern recognition*, pp. 3061–3070, 2015.
- 650
- 651 Jaeho Moon, Juan Luis Gonzalez Bello, Byeongjun Kwon, and Munchurl Kim. From-ground-
652 to-objects: Coarse-to-fine self-supervised monocular depth estimation of dynamic objects with
653 ground contact prior. *arXiv preprint arXiv:2312.10118*, 2023.
- 654 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
655 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in
656 pytorch. 2017.
- 657
- 658 Matteo Poggi, Fabio Tosi, and Stefano Mattoccia. Learning monocular depth estimation with un-
659 supervised trinocular assumptions. In *2018 International conference on 3d vision (3DV)*, pp.
660 324–333. IEEE, 2018.
- 661 René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction.
662 In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 12179–12188,
663 2021.
- 664 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
665 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
666 recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 667
- 668 Ashutosh Saxena, Min Sun, and Andrew Y Ng. Make3d: Learning 3d scene structure from a single
669 still image. *IEEE transactions on pattern analysis and machine intelligence*, 31(5):824–840,
670 2008.
- 671 Shuwei Shao, Zhongcai Pei, Weihai Chen, Xingming Wu, and Zhengguo Li. Nddepth: Normal-
672 distance assisted monocular depth estimation. In *Proceedings of the IEEE/CVF International*
673 *Conference on Computer Vision*, pp. 7931–7940, 2023.
- 674
- 675 Shuwei Shao, Zhongcai Pei, Xingming Wu, Zhong Liu, Weihai Chen, and Zhengguo Li. Iebins:
676 Iterative elastic bins for monocular depth estimation. *Advances in Neural Information Processing*
677 *Systems*, 36, 2024.
- 678 Chang Shu, Kun Yu, Zhixiang Duan, and Kuiyuan Yang. Feature-metric loss for self-supervised
679 learning of depth and egomotion. In *European Conference on Computer Vision*, pp. 572–588.
680 Springer, 2020.
- 681
- 682 Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and sup-
683 port inference from rgbd images. In *Computer Vision—ECCV 2012: 12th European Conference*
684 *on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part V 12*, pp. 746–760.
685 Springer, 2012.
- 686 Xiaogang Song, Haoyue Hu, Li Liang, Weiwei Shi, Guo Xie, Xiaofeng Lu, and Xinhong Hei.
687 Unsupervised monocular estimation of depth and visual odometry using attention and depth-
688 pose consistency loss. *IEEE Transactions on Multimedia*, 2023.
- 689
- 690 Libo Sun, Jia-Wang Bian, Huangying Zhan, Wei Yin, Ian Reid, and Chunhua Shen. Sc-depthv3:
691 Robust self-supervised monocular depth estimation for dynamic scenes. *IEEE Transactions on*
692 *Pattern Analysis and Machine Intelligence*, 2023.
- 693 Yihong Sun, Hariharan Bharath Hariharan, Bharath, and Bharath Hariharan. Dynamo-depth: Fixing
694 unsupervised depth estimation for dynamical scenes. *Advances in Neural Information Processing*
695 *Systems*, 36, 2024.
- 696
- 697 Jamie Watson, Michael Firman, Gabriel J Brostow, and Daniyar Turmukhambetov. Self-supervised
698 monocular depth hints. In *Proceedings of the IEEE/CVF International Conference on Computer*
699 *Vision*, pp. 2162–2171, 2019.
- 700 Yi Wei, Hengkai Guo, Jiwen Lu, and Jie Zhou. Iterative feature matching for self-supervised indoor
701 depth estimation. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(6):3839–
3852, 2021.

- 702 Mingkang Xiong, Zhenghong Zhang, Jiyuan Liu, Tao Zhang, and Huilin Xiong. Monocular depth
703 estimation using self-supervised learning with more effective geometric constraints. *Engineering*
704 *Applications of Artificial Intelligence*, 128:107489, 2024.
- 705
- 706 Jiaxing Yan, Hong Zhao, Penghui Bu, and YuSheng Jin. Channel-wise attention-based network
707 for self-supervised monocular depth estimation. In *2021 International Conference on 3D vision*
708 *(3DV)*, pp. 464–473. IEEE, 2021.
- 709
- 710 Jiayu Yang, Jose M Alvarez, and Miaomiao Liu. Non-parametric depth distribution modelling based
711 depth inference for multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer*
712 *Vision and Pattern Recognition*, pp. 8626–8634, 2022.
- 713
- 714 Zhenheng Yang, Peng Wang, Wei Xu, Liang Zhao, and Ramakant Nevatia. Unsupervised learning
715 of geometry from videos with edge-aware depth-normal consistency. In *Proceedings of the AAAI*
716 *Conference on Artificial Intelligence*, volume 32, 2018.
- 717
- 718 Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstruc-
719 tured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*,
720 pp. 767–783, 2018.
- 721
- 722 Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for
723 high-resolution multi-view stereo depth inference. In *Proceedings of the IEEE/CVF conference*
724 *on computer vision and pattern recognition*, pp. 5525–5534, 2019.
- 725
- 726 Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and
727 camera pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
728 pp. 1983–1992, 2018.
- 729
- 730 Zehao Yu, Lei Jin, and Shenghua Gao. P²net: Patch-match and plane-regularization for unsu-
731 pervised indoor depth estimation. In *European Conference on Computer Vision*, pp. 206–222.
732 Springer, 2020.
- 733
- 734 Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Ex-
735 ploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural*
736 *information processing systems*, 35:25018–25032, 2022.
- 737
- 738 Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. Neural window fully-connected
739 crfs for monocular depth estimation. In *Proceedings of the IEEE/CVF conference on computer*
740 *vision and pattern recognition*, pp. 3916–3925, 2022.
- 741
- 742 Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian Reid.
743 Unsupervised learning of monocular depth estimation and visual odometry with deep feature re-
744 construction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
745 pp. 340–349, 2018.
- 746
- 747 Ning Zhang, Francesco Nex, George Vosselman, and Norman Kerle. Lite-mono: A lightweight
748 cnn and transformer architecture for self-supervised monocular depth estimation. In *Proceedings*
749 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18537–18546,
750 2023a.
- 751
- 752 Sen Zhang, Jing Zhang, and Dacheng Tao. Towards scale-aware, robust, and generalizable unsuper-
753 vised monocular depth estimation by integrating imu motion dynamics. In *European Conference*
754 *on Computer Vision*, pp. 143–160. Springer, 2022a.
- 755
- 756 Yourun Zhang, Maoguo Gong, Jianzhao Li, Mingyang Zhang, Fenlong Jiang, and Hongyu Zhao.
757 Self-supervised monocular depth estimation with multiscale perception. *IEEE transactions on*
758 *image processing*, 31:3251–3266, 2022b.
- 759
- 760 Yourun Zhang, Maoguo Gong, Mingyang Zhang, and Jianzhao Li. Self-supervised monocular depth
761 estimation with self-perceptual anomaly handling. *IEEE Transactions on Neural Networks and*
762 *Learning Systems*, 2023b.

- 756 Chaoqiang Zhao, Matteo Poggi, Fabio Tosi, Lei Zhou, Qiyu Sun, Yang Tang, and Stefano Mattoccia.
757 Gasmono: Geometry-aided self-supervised monocular depth estimation for indoor scenes. In
758 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16209–16220,
759 2023a.
- 760 Hailiang Zhao, Yongyi Kong, Chonghao Zhang, Haoji Zhang, and Jiansen Zhao. Learning effec-
761 tive geometry representation from videos for self-supervised monocular depth estimation. *ISPRS*
762 *International Journal of Geo-Information*, 13(6):193, 2024.
- 764 Jiawei Zhao, Ke Yan, Yifan Zhao, Xiaowei Guo, Feiyue Huang, and Jia Li. Transformer-based dual
765 relation graph for multi-label image recognition. In *Proceedings of the IEEE/CVF international*
766 *conference on computer vision*, pp. 163–172, 2021.
- 767 Rong Zhao, Xie Han, Xindong Guo, Liqun Kuang, Xiaowen Yang, and Fusheng Sun. Exploring the
768 point feature relation on point cloud for multi-view stereo. *IEEE Transactions on Circuits and*
769 *Systems for Video Technology*, 33(11):6747–6763, 2023b.
- 771 Wang Zhao, Shaohui Liu, Yezhi Shu, and Yong-Jin Liu. Towards better generalization: Joint depth-
772 pose learning without posenet. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
773 *and Pattern Recognition*, pp. 9151–9161, 2020.
- 774 Hang Zhou, David Greenwood, and Sarah Taylor. Self-supervised monocular depth estimation with
775 internal feature fusion. *arXiv preprint arXiv:2110.09482*, 2021a.
- 777 Junsheng Zhou, Yuwang Wang, Kaihuai Qin, and Wenjun Zeng. Moving indoor: Unsupervised
778 video depth learning in challenging environments. In *Proceedings of the IEEE/CVF international*
779 *conference on computer vision*, pp. 8618–8627, 2019.
- 780 Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth
781 and ego-motion from video. In *Proceedings of the IEEE conference on computer vision and*
782 *pattern recognition*, pp. 1851–1858, 2017.
- 784 Zhongkai Zhou, Xinnan Fan, Pengfei Shi, and Yuanxue Xin. R-msfm: Recurrent multi-scale fea-
785 ture modulation for monocular depth estimating. In *Proceedings of the IEEE/CVF international*
786 *conference on computer vision*, pp. 12777–12786, 2021b.

788 A OVERVIEW

790 The appendix document supplements the method details and additional experimental results. In Sec-
791 tion **B**, we discuss the scene priors in detail. In Section **C**, we supplement the implementation details.
792 In Section **D**, we present more visualization results on the KITTI (Geiger et al., 2013), Make3D Sax-
793 ena et al. (2008), NYUv2 Silberman et al. (2012), and ScanNet Dai et al. (2017) datasets. In Section
794 **E**, we provide more quantitative comparisons with previous state-of-the-art methods at other image
795 resolutions. In Section **F**, we report the additional ablation study on the indoor scene. In Section
796 **G**, we present the complexity of the model and the speed of inference. In Section **H**, we present
797 a challenging sample. In Section **I**, we present the qualitative results on challenging samples. In
798 Section **J**, we present the visualization of depth maps and reconstructed point clouds. In Section **K**,
799 we present the qualitative results of cropped image. In Section **L**, we discuss the limitations of our
800 method.

802 B DETAILS ABOUT SCENE PRIOR

804 When humans understand the real-world or infer the depth from the 3D scene including indoor or
805 outdoor scene, they will employ specific prior knowledge about the physical world. For the indoor
806 scene, the floor and ceiling are located in the lower and the upper parts respectively, and the room
807 is surrounded by flat walls perpendicular to the floor and ceiling. As for the outdoor scene such as
808 driving scene, the road and sky appear in the lower and upper parts respectively, and other objects
809 such as people, vehicles, houses, etc. are connected and located above the road. Besides, objects
in the middle area of the image often have greater depth than objects in other areas. Inspired by

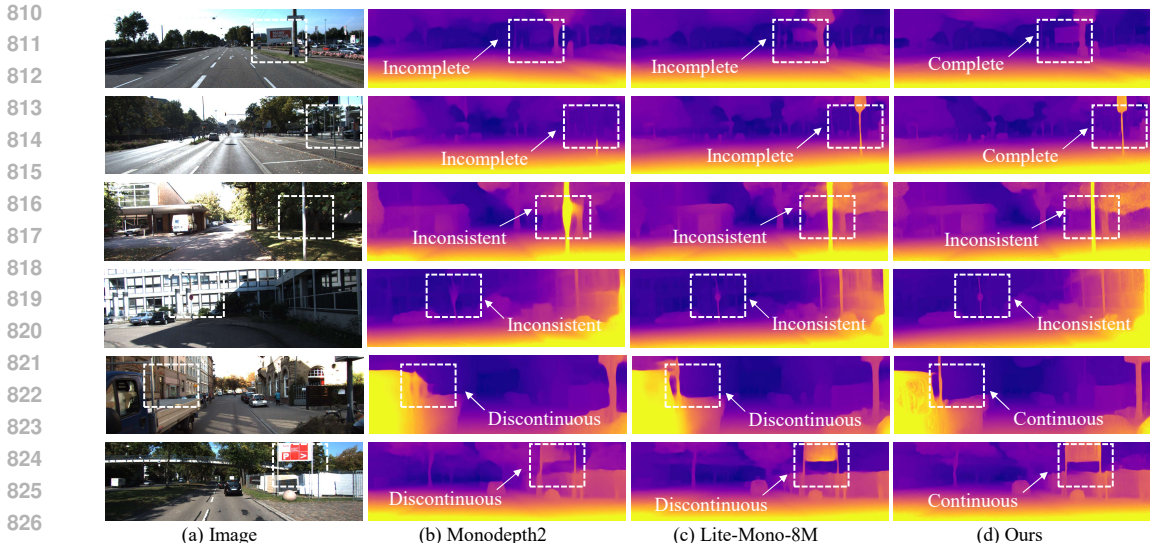


Figure 4: **More qualitative comparison results with a resolution of 640×192 on the KITTI (Geiger et al., 2013) dataset**. We highlight challenging areas with white dashed boxes. Compared with the classical method Monodepth2 (Godard et al., 2019) and the latest state-of-the-art method Lite-Mono-8M (Zhang et al., 2023a), our proposed method generates superior visual results.

these observations, we intuitively argue that the rough depth range can be inferred when the scene nature are known and the relative positions of objects in the image are give. To this end, we mine an additional scene prior, *i.e.*, relative spatial position, for depth estimation model to enhance its perception of spatial structure.

C MORE IMPLEMENTATION DETAILS

Following Zhou et al. (2021a); He et al. (2022), we use the weights pre-trained on ImageNet (Russakovsky et al., 2015) to initialize the encoders of depth network and pose network. In order to fully verify the effectiveness of the proposed framework, we train and validate outdoor scenes at three resolutions, including 640×192 , 1024×320 and 1280×384 . And for indoor scenes, we train and validate only at 320×256 resolution. To improve the training speed, we only output a single-scale depth for the final depth decoder and compute the loss on single-scale depth map. Following existing evaluation rules (Godard et al., 2019; Zhang et al., 2023a), we adopt the same median scaling on the depth results for all methods.

D MORE QUALITATIVE RESULTS

For a clear comparison between TSF-Depth and related works, Fig. 4, Fig. 5, and Fig. 14 present more qualitative comparison results on the KITTI and Make3D datasets. In addition, Fig. 6 and Fig. 7 present more qualitative comparison results on the NYUv2 and ScanNet datasets. We highlight challenging areas with white dashed boxes. As we can see from this figure, our proposed TSF-Depth generates more accurate depth maps, particularly in these region with low texture, slender structure, shadow occlusion and reflections.

E MORE QUANTITATIVE RESULTS

In order to make more comparisons with previous state-of-the-art methods, we also report the performance comparisons at other resolutions on the KITTI (Geiger et al., 2013) dataset, as shown in Table 7. Note that all methods were trained on both 1280×384 and 1024×320 resolutions and then evaluated at the corresponding resolutions. As can be seen, our proposed TSF-Depth significantly outperforms previous self-supervised monocular depth estimation approaches in both resolutions. At a resolution of 1280×384 , almost all errors are reduced by about 3%. Specifically, the Sq Rel,

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

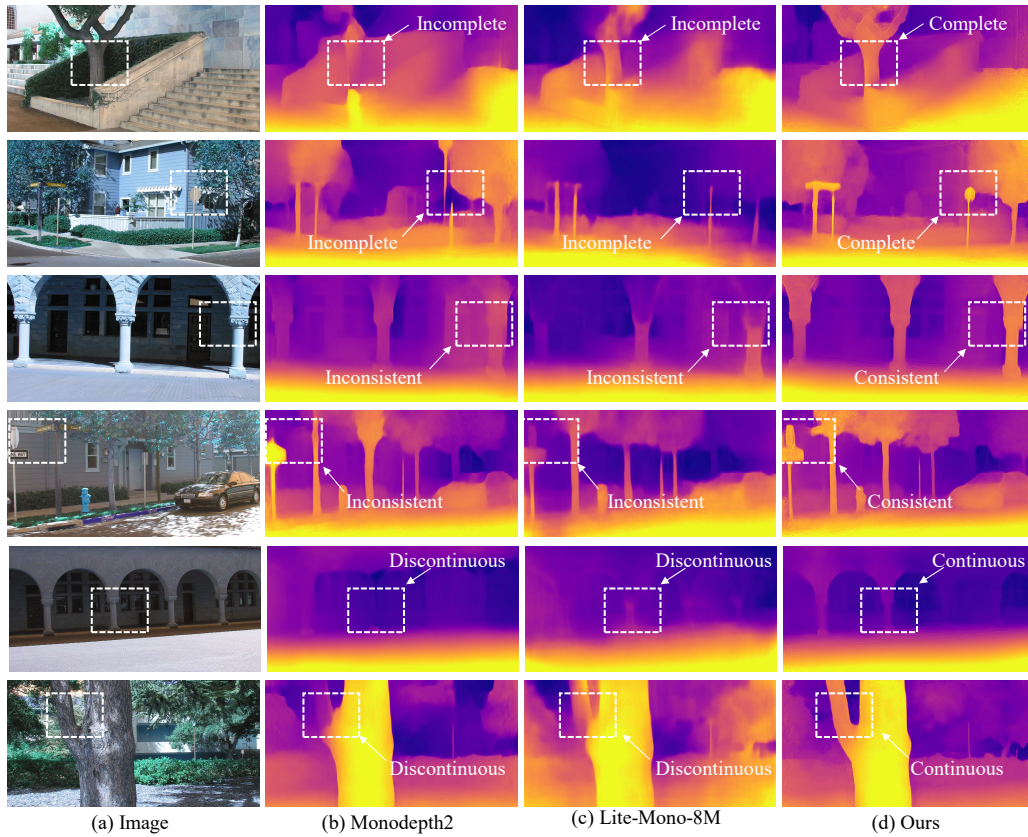


Figure 5: **Qualitative comparison results with a resolution of 640×192 on the Make3D dataset.** We highlight challenging areas. Compared with the classical method Monodepth2 (Godard et al., 2019) and the latest state-of-the-art method Lite-Mono-8M (Zhang et al., 2023a), our proposed method generates superior visual results.

918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

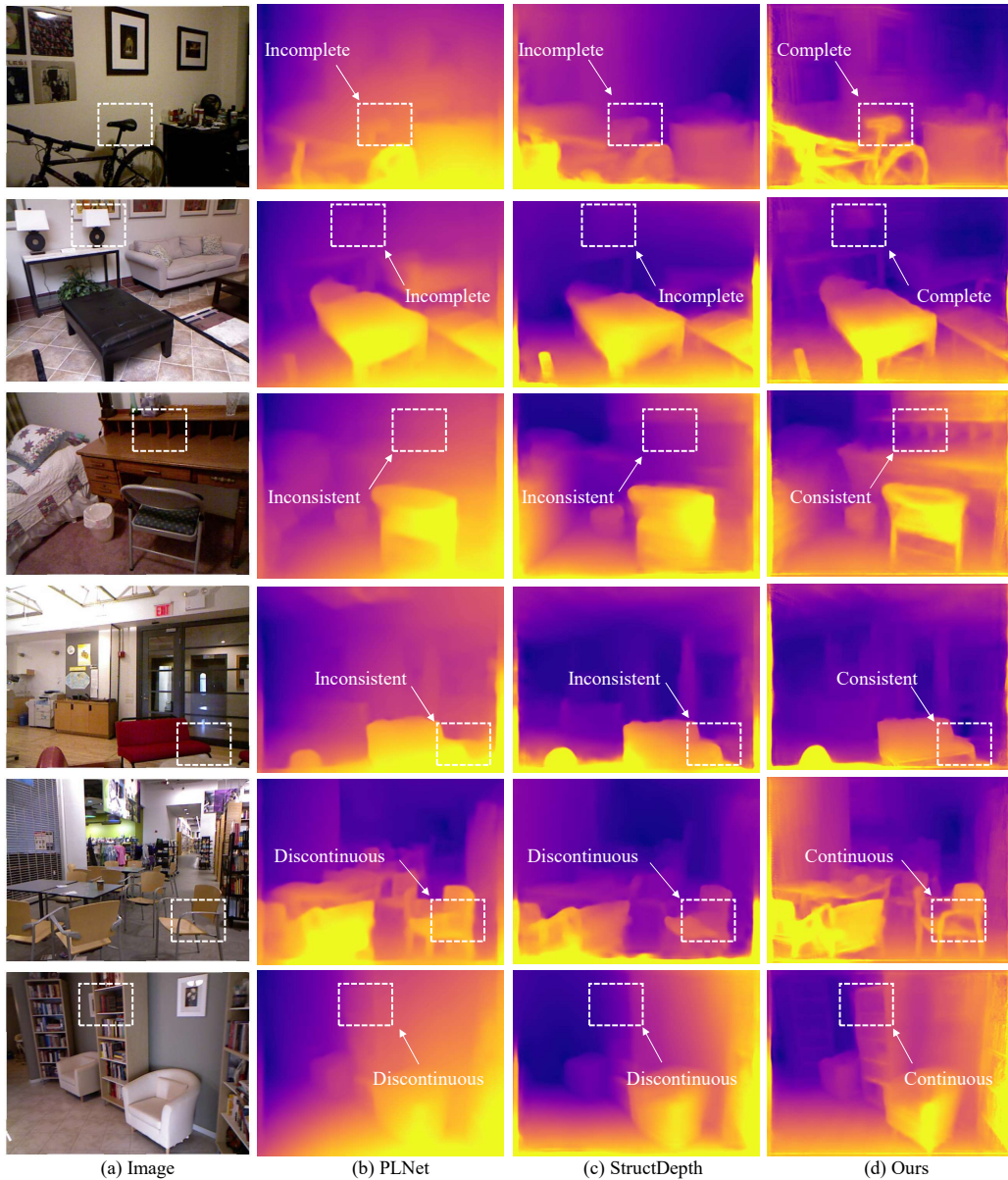


Figure 6: **Qualitative comparison results with a resolution of 640×192 on the NYUv2 dataset.** We highlight challenging areas. Compared with the classical method PLNet (Jiang et al., 2021) and the latest state-of-the-art method StructDepth (Li et al., 2021), our proposed method generates superior visual results.

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

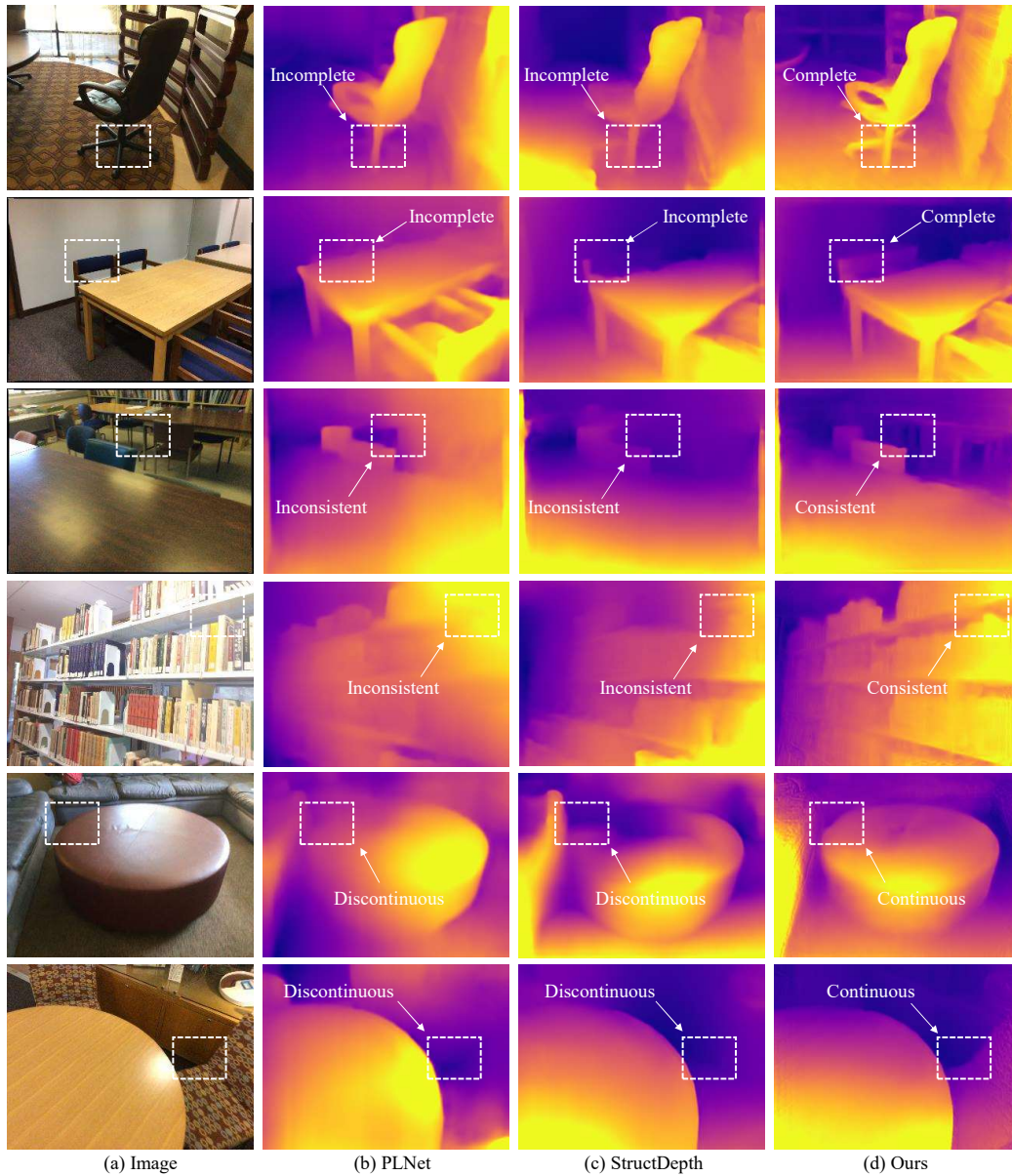


Figure 7: **Qualitative comparison results with a resolution of 640×192 on the ScanNet dataset.** We highlight challenging areas. Compared with PLNet (Jiang et al., 2021) and the latest state-of-the-art method StructDepth (Li et al., 2021), our proposed method generates superior visual results.

Table 7: **Depth estimation results at other resolutions on the Eigen split of KITTI (Geiger et al., 2013) dataset.** The best results are marked in bold.

Method	Resolution	Error Metric (\downarrow)				Accuracy Metric (\uparrow)		
		Sq Rel	Abs Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
PackNet (Guizilini et al., 2020)	1280×384	0.758	0.104	4.384	0.182	0.895	0.964	0.982
SGDepth Klingner et al. (2020)	1280×384	0.768	0.107	4.468	0.186	0.891	0.963	0.982
HR-Depth (Lyu et al., 2021)	1280×384	0.727	0.104	4.410	0.179	0.894	0.966	0.984
CADepth (Zhou et al., 2021a)	1280×384	0.715	0.102	4.312	0.176	0.900	0.968	0.984
TSF-Depth	1280×384	0.679	0.093	4.188	0.170	0.912	0.970	0.985
Monodepth2 (Godard et al., 2019)	1024×320	0.882	0.115	4.701	0.190	0.879	0.961	0.982
HR-Depth (Lyu et al., 2021)	1024×320	0.755	0.106	4.472	0.181	0.892	0.966	0.984
Lite-Mono (Zhang et al., 2023a)	1024×320	0.746	0.102	4.444	0.179	0.896	0.965	0.983
Zhao et al. (Zhao et al., 2024)	1024×320	0.731	0.105	4.412	0.181	0.891	0.965	0.983
TSF-Depth	1024×320	0.672	0.093	4.179	0.169	0.911	0.969	0.984

Abs Rel, RMSE and RMSE log errors are decreased by 8.8%, 5.0%, 2.9%, 3.4%. For the resolution of 1024×320 , almost all errors have a more obvious drop and are reduced by about 5%. The Sq Rel, Abs Rel, RMSE and RMSE log errors are decreased by 8.0%, 11.4%, 5.3%, 6.6%. Therefore, the quantitative results demonstrate that our designed multi-scale multi-view 3D scene field is more robust to depth estimation at different image resolutions.

F ADDITIONAL ABLATION STUDY RESULTS

To fully investigate the main contributions and key designs of TSF-Depth, a series of ablation experiments on the NYUv2 (Silberman et al., 2012) dataset are conducted. Compared to outdoor scenes, the challenges in indoor scenes lies in addressing highly diverse environments and near-field clutter with arbitrarily arranged objects. The pipeline for the baseline is reported in Fig. 1 (a), which follows existing self-supervised monocular depth estimation frameworks and employ only front-view 2D geometric feature for depth estimation.

Effects of Multi-scale Scene Priors. We first analyze the impact of incorporating multi-scale scene prior into model by positional encoding. As shown in Table 8 (a) and (b), training with either single-scale or full-scale scene prior encoding significantly improves the depth prediction accuracy over the baseline without it. In addition, see Table 8 (d), while combining the multi-scale scene priors with tri-plane feature scenes at different scales does not yield significant improvements, the best depth quality was achieved when combining with multi-scale tri-plane feature fields, as shown in 8 (f). Based on the above analysis, the scene prior we introduce is effective for depth estimation even for monocular challenging indoor scenes.

Effects of Multi-scale Tri-plane Feature Fields. We further analyze the impact of modeling 3D scene field using multi-scale tri-plane feature fields. The results are shown in Table 8 (c) and (e). Compared to the baseline, when a single-scale tri-plane feature field is constructed, the depth accuracy is improved. Better performance is achieved by using multi-scale tri-plane or combining scene priors. These ablation results illustrate that our proposed TSF-Depth modeling multi-view 3D scene field representation is effective for robust depth estimation.

G MODEL COMPLEXITY AND SPEED EVALUATION

Table 9 reports the parameter complexity (#Params), computation complexity (GLOPs), and inference speed on the KITTI (Geiger et al., 2013) dataset. We perform inference at a resolution of 640×1920 , and set the batch size to 16. The models for all comparison methods were inferred on the same platform with NVIDIA RTX 3090 GPU. As can be seen from this table, our model has a similar number of parameters as most existing depth estimation methods, such as Monodepth2-R50 (Godard et al., 2019), DynaDepth (Zhang et al., 2022a) and Zhang et al. (2023b), which allows it to be used on edge devices. Moreover, our model achieve similar computation complexity and inference speed as the existing state-of-the-art model Lite-Mono-8M (Zhang et al., 2023a). Compared

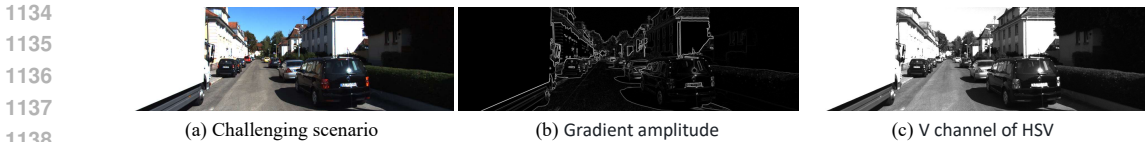
Table 8: **Ablation results for each component of our method on NYUv2 (Silberman et al., 2012).** SP^i : Incorporate scene prior encoding with resolution $\frac{H}{2^i} \times \frac{W}{2^i}$. TP^i : Model the 3D scene using tri-plane feature field with resolution $\frac{H}{2^i} \times \frac{W}{2^i}$. SP^{All}/TP^{All} : Use all resolution SP/TP .

Exp Setting	SP TP	Error Metric ↓				Accuracy Metric ↑		
		Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
(a) Baseline		0.147	0.123	0.574	0.189	0.795	0.956	0.990
SP^1	✓	0.136	0.112	0.544	0.178	0.830	0.968	0.991
SP^2	✓	0.138	0.109	0.541	0.177	0.822	0.964	0.992
SP^3	✓	0.138	0.116	0.559	0.177	0.823	0.965	0.991
SP^4	✓	0.135	0.110	0.543	0.174	0.830	0.965	0.992
(b) SP^5	✓	0.136	0.106	0.536	0.176	0.830	0.966	0.992
SP^{All}	✓	0.137	0.109	0.539	0.175	0.830	0.966	0.992
TP^1	✓	0.139	0.120	0.558	0.178	0.826	0.964	0.991
TP^2	✓	0.143	0.113	0.551	0.183	0.812	0.959	0.991
TP^3	✓	0.142	0.120	0.559	0.180	0.818	0.963	0.991
TP^4	✓	0.136	0.107	0.536	0.174	0.830	0.965	0.992
(c) TP^5	✓	0.138	0.115	0.548	0.177	0.829	0.963	0.991
TP^{All}	✓	0.136	0.112	0.549	0.182	0.812	0.961	0.991
$SP^{All} + TP^1$	✓ ✓	0.139	0.113	0.550	0.178	0.822	0.963	0.991
$SP^{All} + TP^2$	✓ ✓	0.137	0.112	0.549	0.176	0.826	0.965	0.991
$SP^{All} + TP^3$	✓ ✓	0.139	0.118	0.556	0.178	0.825	0.963	0.991
(d) $SP^{All} + TP^4$	✓ ✓	0.135	0.111	0.546	0.175	0.830	0.966	0.991
$SP^{All} + TP^5$	✓ ✓	0.135	0.109	0.536	0.176	0.829	0.964	0.991
$TP^{All} + SP^1$	✓ ✓	0.140	0.111	0.548	0.179	0.816	0.963	0.992
$TP^{All} + SP^2$	✓ ✓	0.139	0.112	0.548	0.178	0.823	0.962	0.991
$TP^{All} + SP^3$	✓ ✓	0.141	0.115	0.554	0.180	0.816	0.961	0.991
(e) $TP^{All} + SP^4$	✓ ✓	0.141	0.115	0.558	0.181	0.817	0.962	0.991
$TP^{All} + SP^5$	✓ ✓	0.135	0.108	0.537	0.173	0.830	0.967	0.992
(f) TSF-Depth ($SP^{All} + TP^{All}$)	✓ ✓	0.134	0.103	0.530	0.172	0.831	0.967	0.992

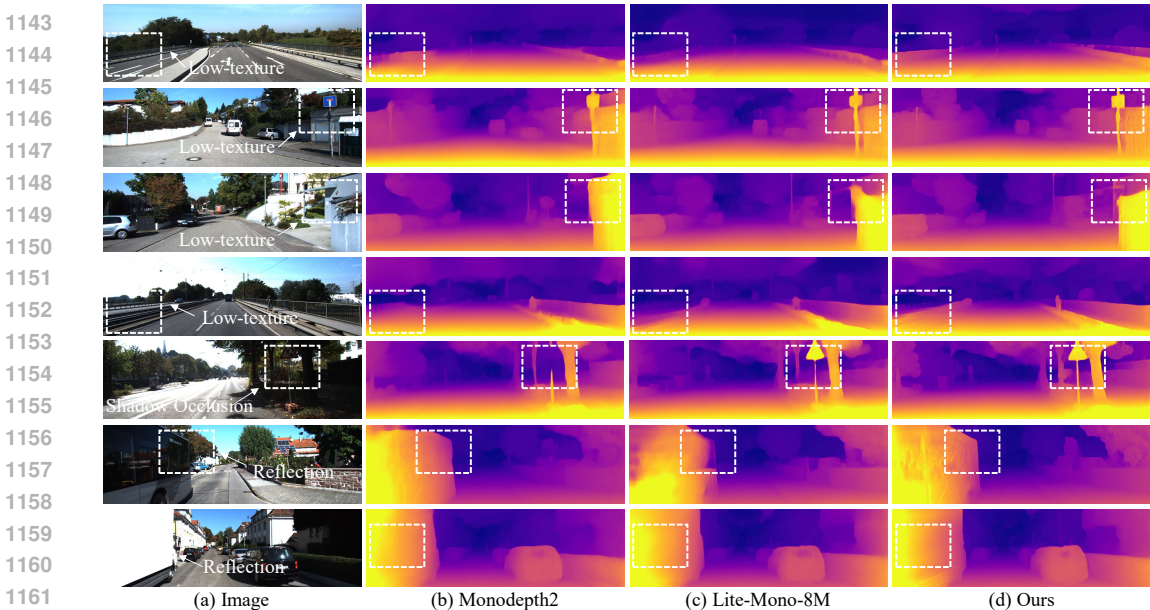
Table 9: **Model complexity and speed evaluation.** We compare parameters (#Params), giga floating-point operations per second (GFLOPS), and inference speed on the KITTI (Geiger et al., 2013) dataset. The input size is 640×192 , and the batch size is 16. All models are inferred on the same platform with NVIDIA RTX 3090 GPU. “-” indicates that the method is not open source code and we cannot make inferences.

Method	#Params	GFLOPs	Speed	Error Metric (↓)			
				Sq Rel	Abs Rel	RMSE	RMSE log
GeoNet (Yin & Shi, 2018)	31.6M	-	-	1.060	0.149	5.567	0.226
Monodepth2-R18 (Godard et al., 2019)	14.3M	8.04G	1.8ms	0.903	0.115	4.863	0.193
Monodepth2-R50 (Godard et al., 2019)	32.5M	16.7G	2.0ms	0.831	0.110	4.642	0.187
DynaDepth (Zhang et al., 2022a)	32.5M	16.7G	3.4ms	0.761	0.108	4.608	0.187
Zhang et al. (Zhang et al., 2023b)	32.6M	-	-	0.786	0.105	4.572	0.182
Lite-Mono-8M (Zhang et al., 2023a)	8.70M	11.2G	6.3ms	0.729	0.101	4.454	0.178
Dynamo-Depth (MD2) (Sun et al., 2024)	14.3M	8.04G	1.4ms	0.864	0.120	4.850	0.195
Dynamo-Depth (Sun et al., 2024)	8.77M	11.2G	4.7ms	0.758	0.112	4.505	0.183
Baseline	27.2M	39.9G	7.7ms	0.746	0.102	4.464	0.176
TSF-Depth	29.8M	44.2G	9.8ms	0.692	0.096	4.335	0.173

to the baseline model, our method scales up to state-of-the-art depth performance with only a few additional parameters. The additional parameters are brought by the initial depth network, however it does not produce any features information to be incorporated into the final depth estimation, only providing sampling points. Therefore, our proposed TSF-Depth is able to be practically applied.



1139 Figure 8: **A challenging sample.** (a) A scenario with extensive low-texture areas (*e.g.*, roads, buildings, bushes, etc.) and reflective areas (*e.g.*, cars, buildings). (b) The gradient amplitude computed using the Sobel operator to assess the texture distribution. (c) The V channel in HSV color space used to analyze brightness levels.



1162 Figure 9: **Qualitative comparison of challenging test subsets in the KITTI dataset.** We highlight challenging areas with low-texture areas, shadow occlusion, and reflective surfaces. Our method generates robust depth maps in these challenging conditions.

1166 H CHALLENGING SAMPLE

1167
1168 In general, low-texture regions exhibit smooth variations with gradient values typically close to zero, while reflective regions often display extremely high luminance values near saturation. We present a challenging sample containing extensive low-texture areas (*e.g.*, roads, buildings, bushes, etc.) and reflective areas (*e.g.*, cars, buildings) in Fig. 8. To analyze these characteristics, we illustrate its gradient magnitude (Fig. 8 (b)) to evaluate the texture level and the V channel in HSV space (Fig. 8 (c)) to examine brightness. The observations align well with common perceptions.

1175 I QUALITATIVE RESULTS ON CHALLENGING SAMPLES

1176
1177 We showcase depth estimation results on subsets with low-texture areas, shadow occlusion, and reflective surfaces, as shown in Fig. 9. The highlighted regions (dotted boxes) emphasize the differences among Monodepth2, Lite-Mono-8M, and our method. Our approach demonstrates more accurate and robust depth predictions in these challenging conditions.

1182 J VISUALIZATION OF DEPTH MAPS AND RECONSTRUCTED POINT CLOUDS

1183
1184 As illustrated in Fig. 10 (b), we present the visualization results of the depth map at six scales (2st column) alongside the corresponding point clouds without (3rd column) and with (4rd column) color values of the target image. Additionally, we show the reconstructed results (6th and 7th columns) obtained after upsampling the predicted depth map to match the real image resolution, as depicted in Fig. 10 (c). Furthermore, we provide reconstructed results of Fig. 10 from the multiple

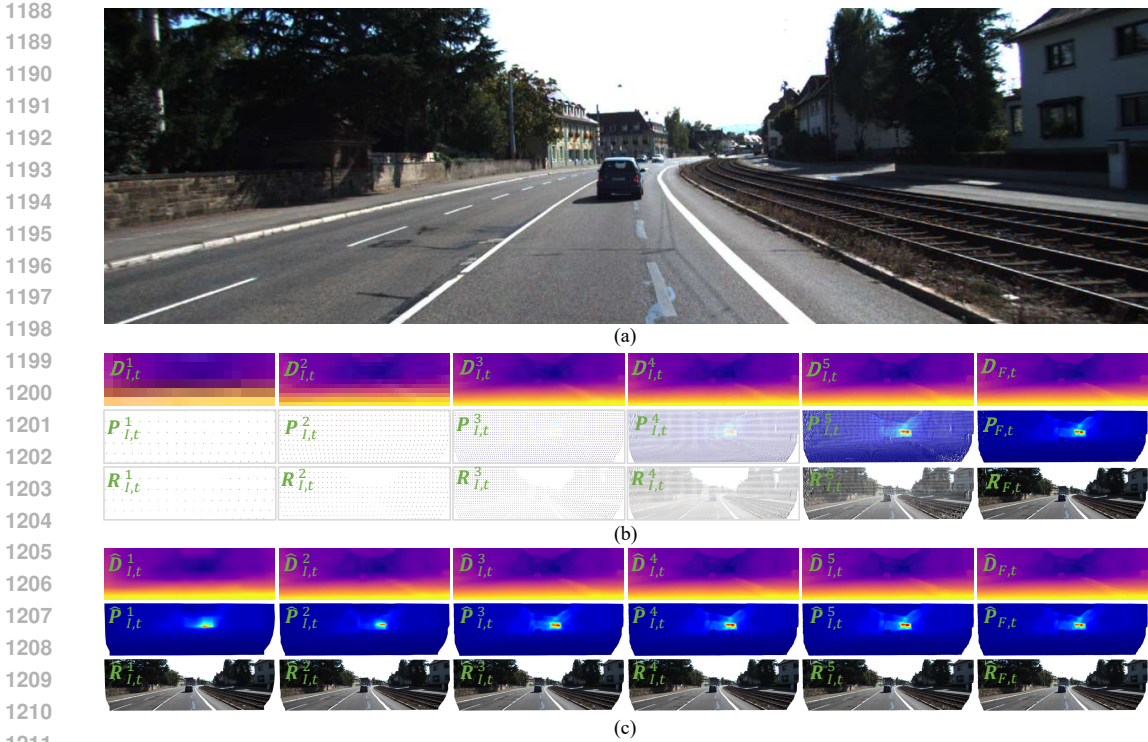


Figure 10: **Qualitative of the depth maps, point clouds without and with color values of the target image.** (a) The target image with a resolution of 375×1242 . (b) Qualitative results at different scales when target image are input at training resolution 192×640 , including predicted depth maps ($\mathbf{D}_{I,t}^s \in \mathbb{R}^{192/2^s \times 640/2^s}$ and $\mathbf{D}_{F,t}^s \in \mathbb{R}^{192 \times 640}$) (2nd column), projected point clouds ($\mathbf{P}_{I,t}^s \in \mathbb{R}^{192/2^s \times 640/2^s \times 3}$, $\mathbf{P}_{F,t}^s \in \mathbb{R}^{192 \times 640 \times 3}$, $\mathbf{R}_{I,t}^s \in \mathbb{R}^{192/2^s \times 640/2^s \times 3}$ and $\mathbf{R}_{F,t}^s \in \mathbb{R}^{192 \times 640 \times 3}$) without (3rd column) and with (4rd column) color values of the target image (c) Qualitative results after the predicted depth is upsampled to the real image resolution.

perspectives in Fig.12. The more qualitative results are reported in Fig. 13. We observed that under the three viewing perspectives, the overall structure of point clouds corresponding to different scales is roughly consistent.

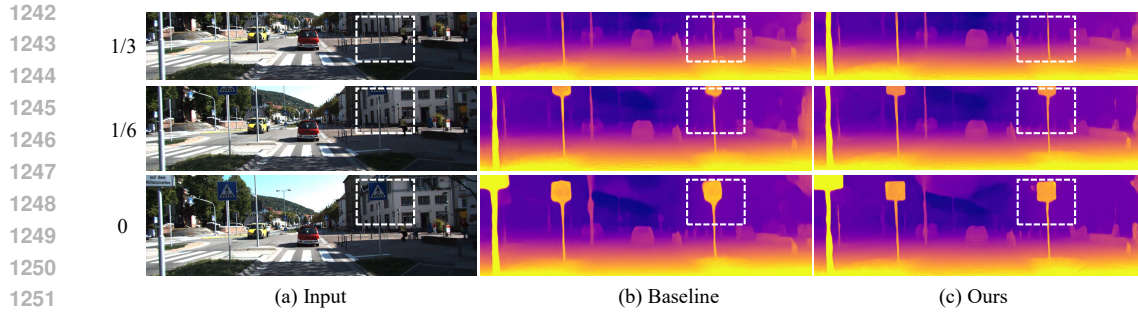
corresponding point clouds without (3rd column) and with (4rd column) color values of the target image

K QUALITATIVE RESULTS OF CROPPED IMAGE

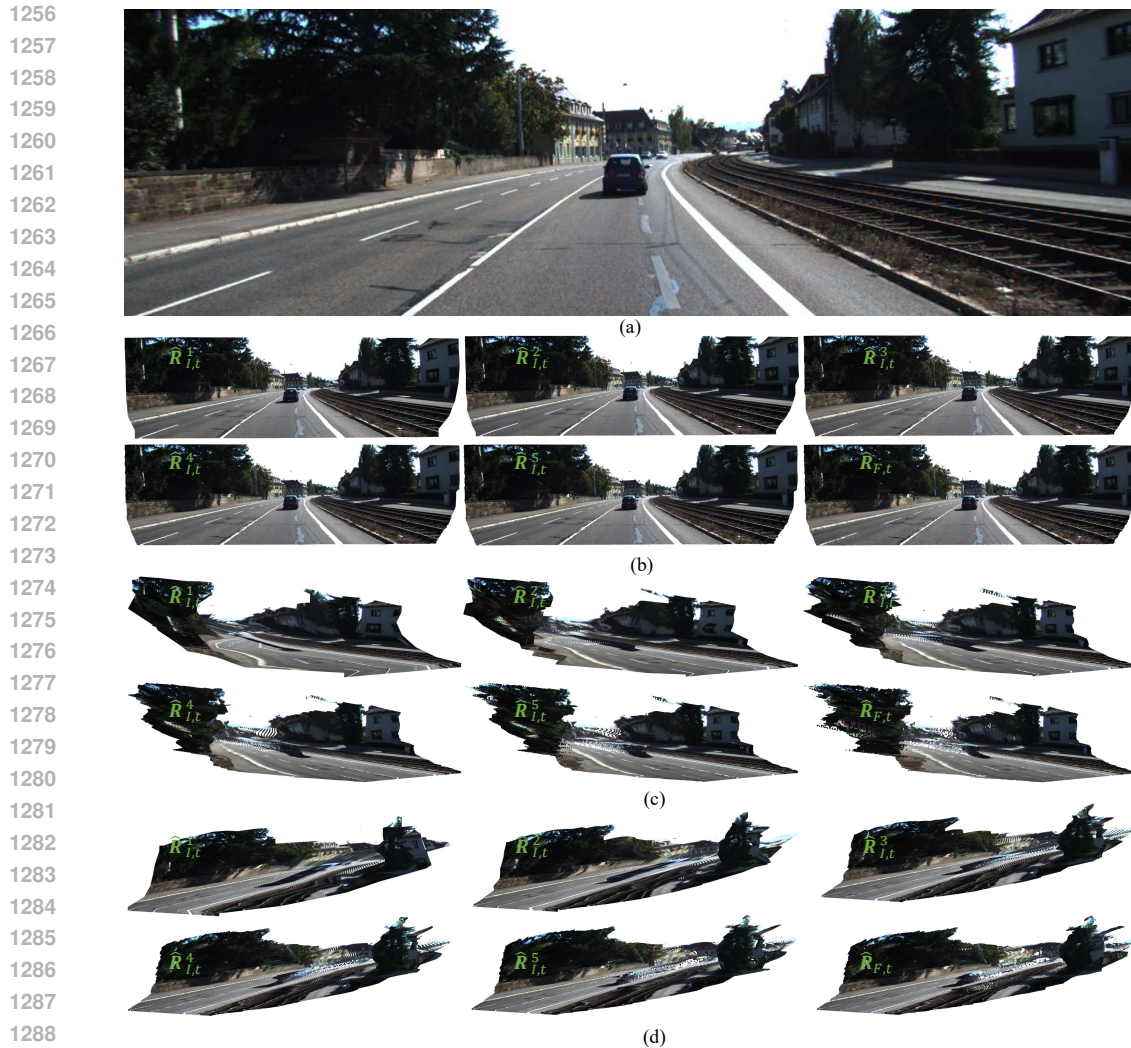
As shown in Fig. 11, we present the visualization results of the depth maps corresponding to different cropping ratios. We observe that the Baseline produces discontinuous and inaccurate depth results when the spatial layout of the image changes. In contrast, our method generates consistent results across all three cropping ratios, demonstrating its applicability to challenging cases such as cropped images.

L LIMITATION

Although our proposed method has better overall performance in both quantitative and qualitative aspects compared to existing self-supervised methods, there are still limitations. First, insufficient smoothness in glass-related reflective areas, such as the fifth sample in Fig. 4, and the last two samples in Fig. 9. Possible solutions to address this limitation is to incorporate semantic information into the framework, enabling the model to better differentiate reflective regions and enforce smoothness constraints. The second limitation is that predictions may be inaccurate for individual targets that are similar to the overall scene, such as the cabinet in the last sample of Fig. 7. This limitation



1252 Figure 11: **Qualitative comparison of cropped images on the KITTI dataset.** We present depth
1253 estimation results on images with different crop ratios (1/3, 1/6, and 0). Highlighted regions (dot-
1254 ted boxes) demonstrate that our method achieves more accurate depth predictions compared to the
1255 Baseline under varying crop conditions.



1289 Figure 12: **Qualitative results of reconstructed point clouds from multiple perspectives.** These
1290 visualization correspond to Fig. 10.

1291
1292 is likely due to the scarcity of such samples in the training datasets, as these datasets were primarily
1293 used for generalization testing. To mitigate this, future work could incorporate semantic information
1294 and leverage plane priors to improve predictions in these challenging scenarios.

1295

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

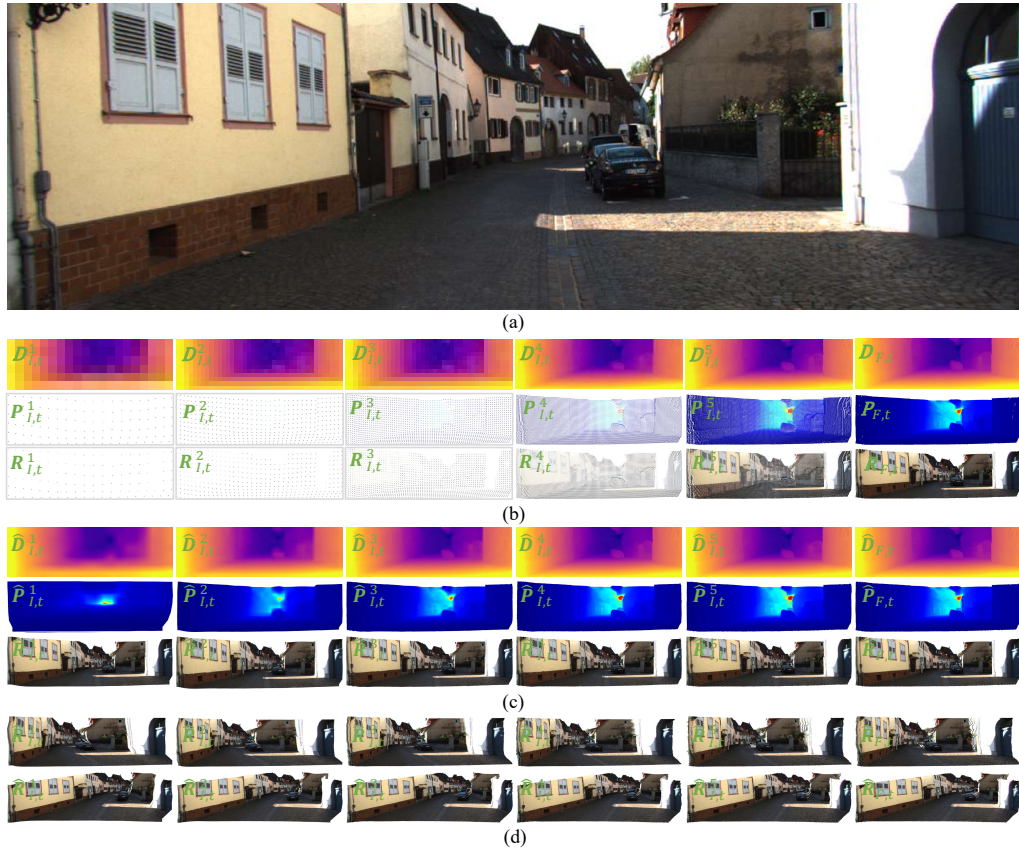


Figure 13: **More qualitative results of the depth maps, point clouds without and with color values of the target image.** (a) The target image with a resolution of 375×1242 . (b) Visualization results at different scales when target image are input at training resolution 192×640 , including predicted depth maps ($D_{I,t}^s \in \mathbb{R}^{192/2^s \times 640/2^s}$ and $D_{F,t}^s \in \mathbb{R}^{192 \times 640}$) (2nd column), projected point clouds ($P_{I,t}^s \in \mathbb{R}^{192/2^s \times 640/2^s \times 3}$, $P_{F,t}^s \in \mathbb{R}^{192 \times 640 \times 3}$, $R_{I,t}^s \in \mathbb{R}^{192/2^s \times 640/2^s \times 3}$ and $R_{F,t}^s \in \mathbb{R}^{192 \times 640 \times 3}$) without (3rd column) and with (4rd column) color values of the target image (c) Qualitative results after the predicted depth is upsampled to the real image resolution. (d) Visualization results of reconstructed point clouds from multiple perspectives.

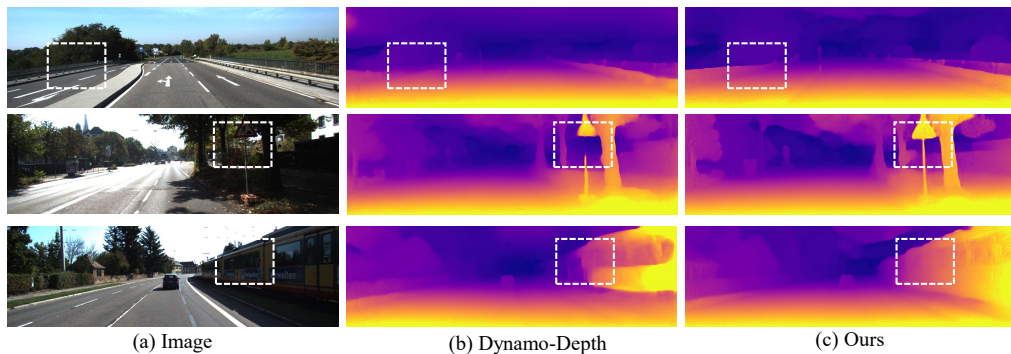


Figure 14: **Qualitative comparison with recent SOTA work Dynamo-Depth (Sun et al., 2024).**