# Monocular Road Planar Parallax Estimation

Haobo Yuan, Teng Chen, Wei Sui, Jiafeng Xie, Lefei Zhang, *Senior Member, IEEE*
Yuan Li, and Qian Zhang

*Abstract*—Estimating the 3D structure of the drivable surface and surrounding environment is a crucial task for assisted and autonomous driving. It is commonly solved either by using 3D sensors such as LiDAR or directly predicting the depth of points via deep learning. However, the former is expensive, and the latter lacks the use of geometry information for the scene. In this paper, instead of following existing methodologies, we propose Road Planar Parallax Attention Network (RPANet), a new deep neural network for 3D sensing from monocular image sequences based on planar parallax, which takes full advantage of the omnipresent road plane geometry in driving scenes. RPANet takes a pair of images aligned by the homography of the road plane as input and outputs a $\gamma$ map (the ratio of height to depth) for 3D reconstruction. The $\gamma$ map has the potential to construct a two-dimensional transformation between two consecutive frames. It implies planar parallax and can be combined with the road plane serving as a reference to estimate the 3D structure by warping the consecutive frames. Furthermore, we introduce a novel cross-attention module to make the network better perceive the displacements caused by planar parallax. To verify the effectiveness of our method, we sample data from the Waymo Open Dataset and construct annotations related to planar parallax. Comprehensive experiments are conducted on the sampled dataset to demonstrate the 3D reconstruction accuracy of our approach in challenging scenarios.

*Index Terms*—Planar Parallax Estimation, 3D Computer Vision, Deep Learning.

## I. INTRODUCTION

**3**D reconstruction of road environment [1]–[3] has received increasing attention in the assisted and autonomous driving field [4], [5] as it is crucial for obstacle detection [6], distance measurement [7], [8], and road condition recognition [9], etc. Existing methods commonly exploit 3D sensors such as Li-DAR or adopt vision-based 3D reconstruction algorithms. 3D sensors [10] can provide reasonably accurate 3D information, but the high price, the sparse nature, along with the concerns about reliability limit their potential in mass production. In contrast, vision-based methods such as Structure from Motion (SfM) [11] are low-cost yet suffer from various conditions

such as weakly-textured regions, lighting variation, and rapid movement. Besides, these methods usually require laborious manual parameter tuning to guarantee good performance. Recently, deep learning-based methods have been applied to 3D road reconstruction [12]–[14] and have shown promising performance. In these methods, the learning-based depth estimation tasks can be summarized as directly regressing depth numerical values from image pixels using convolutional neural networks (CNNs) trained in a supervised [15] or unsupervised manner [16]. However, simply applying deep neural networks to depth estimation and considering it as a per-pixel regression task means a lack of use of scene geometry information.

In this paper, we would like to bring attention to a family of methods utilizing planar parallax geometry for 3D reconstruction [17]–[20]. Planar parallax was first proposed in the 1990s used for planar motion modeling. The core idea behind it is that the 3D structure is strongly related to the residual image displacements caused by homography warping between two views. Although planar parallax-based methods require a "plane" as a reference, which may be hard to find in some scenes, 3D road reconstruction is naturally a good application of planar parallax as the ground plane serves perfectly as the reference plane. However, they are susceptible to image noise and only suitable for rigid scenes, which prevents them from being widely adopted [21]. In the real-world datasets on the way, such as Waymo Open Dataset [22], the traditional planar parallax method may fail according to our experiments.

To overcome the drawbacks of traditional methods, we design a deep neural network named RPANet to estimate dense planar parallax. RPANet utilizes a novel cross-attention mechanism, takes a pair of images aligned by road plane homography as input, and outputs a pixel-wise $\gamma$ map that represents the pixel-wise ratio of height to depth. A photometric loss could then be applied to train the network since the planar parallax is a *de facto* bridge between the homography-aligned images. Note that the photometric loss relies on the residual flow, which is in the realm of traditional planar parallax geometry. Leveraging both traditional geometry and deep learning with the derived geometry formulas, our method benefits from the robustness of deep learning and the interpretability of geometry-based algorithms.

As there is no publicly available dataset for planar parallax estimation on the road, we build a Road Planar Parallax Dataset (RP2-Waymo) based on the Waymo Open Dataset [22] to train and evaluate our method. We sample data from the Waymo Open Dataset for its diversity in scenes, seasons, and time of day as well as the excellent synchronization between LiDAR and cameras. Beyond the original samples,

we estimate the road plane from the LiDAR points, with which both the homography matrix and the ground truth $\gamma$ can be computed. Sparse ground truth of depth and height is generated by projecting LiDAR points to images. Except for the depth and height, we construct the high-precision road homography, which warps the input image pairs. Thanks to the homography matrix, we can get the image pairs with the aligned road, which is the input of our proposed RPANet. The main contributions of our work are summarized as follows:

- Inspired by the traditional planar parallax geometry, we propose to take the omnipresent road as the reference for 3D structure estimation with deep learning. The predicted planar parallax can be used to reconstruct depth and height of each pixel as well as boost the learning.
- Motivated by the attention mechanism in the stereo-matching, we propose a novel deep neural network called RPANet, which leverages a novel cross-attention module. The cross-attention module can be utilized to find the matching relationship between two images easily and is conducive to predicting the planar parallax.
- To validate our proposed method, we build the RP2-Waymo dataset based on the Waymo Open Dataset [22] for planar parallax estimation. Extensive experiments are conducted on this dataset, and the results demonstrate that our method can recover accurate 3D road surface structures.

## II. RELATED WORK

*a) Planar Parallax:* The planar parallax model is first proposed in [20], [23] to derive a 3D structure relative to a planar surface. They demonstrate that by leveraging the planar parallax model to remove the camera rotation, the reconstruction becomes more accurate and stable. Inspired by [20], [23]. Kumar *et al.* [19] propose a method applying the planar parallax model to estimate the height in aerial images. Observing that the depth of all points in aerial images is nearly the same, they eliminate the depth factor in the parallax equation by calculating the normal of the ground plane. Irani *et al.* [21] further extend the two frames algorithms of [20], [23] to multiple frames and improve the robustness of 3D structure reconstruction. However, most traditional planar parallax algorithms need to obtain accurate correspondence in advance, therefore susceptible to noise. Furthermore, traditional (non-learning-based) planar parallax algorithms also suffer from low speed. haney *et al.* [24] train a deep neural network constrained by planar parallax for the the event-based camera. A concurrent work [25] also uses planar parallax geometry but focusing more on depth estimation. Our aim in this paper is to utilize deep learning to estimate planar parallax, which can effectively determine the depth and height of each pixel, to estimate the structure of a driving scene.

*b) Learning-based 3D Structure Estimation:* One of the earliest works trying to estimate 3D structure from a 2D image by the neural network is proposed by Eigen *et al.* [15]. They directly regress the depth information from a single image. Except for their attempt to predict the depth map to estimate 3D structure, other works try to estimate the 3D structure

by predicting the point cloud [26], voxel [27], mesh [28], or implicit function [29]. However, although these works [26]–[29] have the potential to predict the whole 3D models, even including the backside, they mainly focus on the single object 3D reconstruction rather than scene 3D reconstruction. At the same time, especially for the autonomous driving use case, geometric constraints are introduced to boost the prediction of dense depth. E.g., [12], [16], [30], [31] use 3D geometry constraints to train the depth estimation networks by re-projecting the depth map, while [32]–[35] construct constraints through stereo geometry. The application of different 3D geometric constraints above helps reduce the difficulty of training and further improves the accuracy of 3D structure estimation. Distinct from the previous works, we apply a novel geometric constraint based on planar parallax.

*c) Visual Transformer:* Transformer [36] is first used in the Natural Language Processing (NLP) tasks and origins from the self-attention mechanism. It shows extraordinary performance on many NLP tasks [37]. Inspired by the self-attention mechanism, many researchers explore a similar mechanism to solve computer vision tasks, such as classification [38], object detection [39], [40] and image generation [41]. In general, the attention mechanism can enhance the global perception of neural networks beyond CNNs [38], [42]. However, the self-attention mechanism needs lots of computational resources. They often use patch-based [38] or axial-based [43] strategies to reduce the computational overhead. In the field of 3D structure estimation, thanks to the epipolar constraint, some stereo match schemes [44], [45] also apply attention-based mechanisms to build their neural networks without extremely high computational costs. In our proposed method, we also apply a proposed cross-attention mechanism to model the displacement caused by the homography transformation.
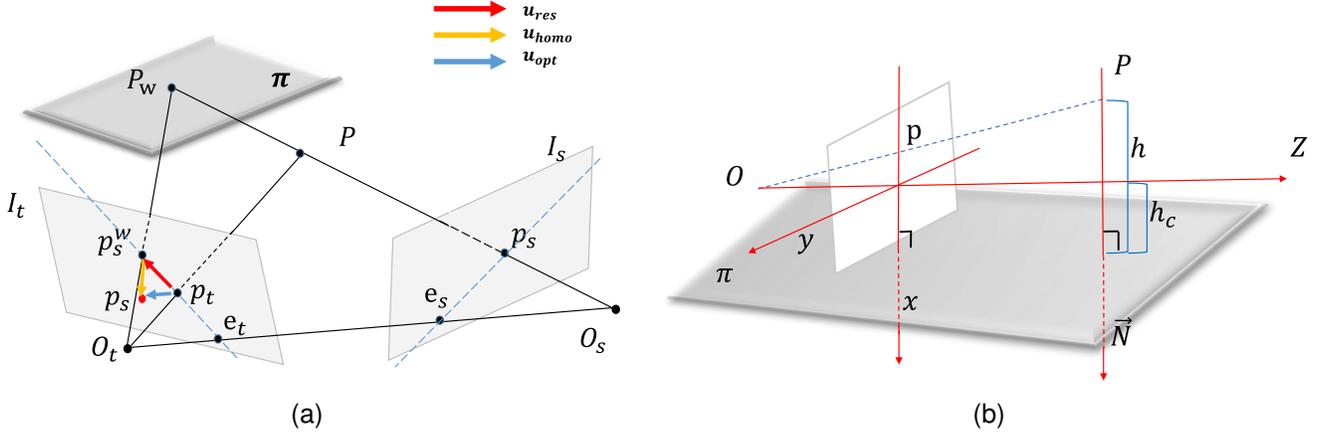
## III. OUR APPROACH

Instead of predicting depth directly through deep neural network such as in [12], [16], [31], [46], our proposed method estimates a $\gamma$ map, which is denoted as

$$\gamma = \frac{h}{d}. \tag{1}$$

The $\gamma$ map represents the ratio of height $h$ to depth $d$ for each pixel from two consecutive images $I_s$ and $I_t$. The pair of input images would be aligned twice, first by the road plane homography and then by the residual flow generated from the $\gamma$ map. After the final alignment, the static region between two images should be well aligned, hence photometric error can be calculated.

### A. Planar Parallax Geometry

*a) Planar Parallax.:* A reference plane is denoted as $\pi$ in 3D space and $\mathbf{P}$ is a point off $\pi$. The point $\mathbf{P}$ is observed by two camera views whose optical centers are represented as $\mathbf{O}_t$ and $\mathbf{O}_s$. $\mathbf{p_t}$ and $\mathbf{p}_s$ are the re-projections on image $I_t$ and $I_s$ respectively. Supposing $\mathbf{P}^w$ is the intersection of the ray $\overrightarrow{\mathbf{O}_s\mathbf{P}}$ and $\pi$, we can obtain its re-projection in camera $\mathbf{O}_t$ denoted as $\mathbf{p}_s^w$ by the homography $\mathbf{H}_{s \rightarrow t}$ (See Fig. 1a for

Fig. 1. (a) The planar parallax geometry. (b) The illustration of 3D reconstruction from $\gamma$.

details). Based on this, the relationship of $\mathbf{p}_t$, $\mathbf{p}_s$, and $\mathbf{p}_t^w$ can be represented as follows:

$$
\begin{aligned}
\mathbf{u}_{homo} &= \mathbf{p}_s - \mathbf{p}_s^w, \\
\mathbf{u}_{res} &= \mathbf{p}_s^w - \mathbf{p}_t, \\
\mathbf{u}_{opt} &= \mathbf{p}_s - \mathbf{p}_t = \mathbf{u}_{homo} + \mathbf{u}_{res},
\end{aligned}
\tag{2}
$$

where $\mathbf{u}_{homo}$ is the displacement caused by homography while $\mathbf{u}_{res}$ is the residual flow. The residual flow represents the displacement of corresponding pixels between a pair of images already aligned by the road homography. We hold that $\mathbf{u}_{opt}$ is a de facto bridge between traditional 3D geometry ($\mathbf{u}_{homo}$) and deep neural network ($\mathbf{u}_{res}$ derived from $\gamma$).

In Fig. 1a, following [16] each 3D point $\mathbf{P}'$ in the source camera coordinate system denoted by $\mathbf{O}_s$ could be transformed to the target camera coordinate system denoted by $\mathbf{O}_t$ by a rigid transformation

$$
\mathbf{P} = \mathbf{R}\mathbf{P}' + \mathbf{T}, \tag{3}
$$

where $\mathbf{T} = (\mathbf{T}_x, \mathbf{T}_y, \mathbf{T}_z)$ is the translation vector in $\mathbf{O}_t$, and $\mathbf{R}$ is the rotation matrix from $\mathbf{O}_s$ to $\mathbf{O}_t$.

Except for the depth, we need to introduce the height of each 3D point in our proposed method. Given an arbitrary $\mathbf{P} = (X, Y, Z)$ in camera coordinate system, where $Z$ is the depth, its height could be expressed as

$$
h = h_c - \vec{\mathbf{N}}^T \mathbf{P}. \tag{4}
$$

where $\vec{\mathbf{N}}$ is the normal of $\pi$, and $h_c$ is the height of camera to plane $\pi$. Then, we have

$$
\frac{h + \vec{\mathbf{N}}^T \mathbf{P}}{h_c} = 1. \tag{5}
$$

Multiplying $\mathbf{T}$ by Eqn. 5 we could obtain

$$
\begin{aligned}
\mathbf{P} &= \mathbf{R}\mathbf{P}' + \mathbf{T}\frac{h - \vec{\mathbf{N}}^T \mathbf{P}'}{h_c} \\
&= (\mathbf{R} + \frac{\mathbf{T}\vec{\mathbf{N}}^T}{h_c})\mathbf{P}' + \frac{h}{h_c}\mathbf{T}.
\end{aligned}
\tag{6}
$$

With $\mathbf{t} = \mathbf{KT} = (t_x, t_y, t_z)^T$, $\mathbf{p} = \frac{\mathbf{K}}{Z}\mathbf{P} = (x, y, 1)^T$, $\mathbf{p}' = \frac{\mathbf{K}'}{Z'}\mathbf{P}' = (x', y', 1)^T$, we obtain

$$
Z\mathbf{K}^{-1}\mathbf{p} = (\mathbf{R} + \frac{\mathbf{T}\vec{\mathbf{N}}^T}{h_c})Z'\mathbf{K}^{-1}\mathbf{p}' + \frac{h}{h_c}\mathbf{T}. \tag{7}
$$

We multiply Eqn. 7 by $\frac{\mathbf{K}}{Z'}$ on both sides and obtain

$$
\begin{aligned}
\frac{Z}{Z'}\mathbf{p} &= \mathbf{K}(\mathbf{R} + \frac{\mathbf{T}\vec{\mathbf{N}}^T}{h_c})\mathbf{K}^{-1}\mathbf{p}' + \frac{h}{h_c Z'}\mathbf{t} \\
&= \mathbf{H}\mathbf{p}' + \frac{h}{h_c Z'}\mathbf{t},
\end{aligned}
\tag{8}
$$

where $\mathbf{H} = \mathbf{K}(\mathbf{R} + \frac{\mathbf{T}\vec{\mathbf{N}}^T}{h_c})\mathbf{K}^{-1}$ represents the homography matrix between the two images. Eqn. 8 could be reformulated as

$$
\mathbf{p} = \frac{\mathbf{H}\mathbf{p}' + \frac{h}{h_c Z'}\mathbf{t}}{\frac{Z}{Z'}}. \tag{9}
$$

The z-axis of $\mathbf{p}$ and $\mathbf{p}'$ is 1. Only the third row of $\mathbf{H}$ and $\mathbf{t}$ contains the information of z-axis. To obtain $\mathbf{p} = (x, y, 1)$, we apply $\frac{Z}{Z'}$ to normalize the z-axis of $\mathbf{H}\mathbf{p}' + \frac{h}{h_c Z'}\mathbf{t}$. Notice that we could get the derivation

$$
\frac{Z}{Z'} = \mathbf{H}_3\mathbf{p}' + \frac{h\mathbf{T}_z}{h_c Z'}, \tag{10}
$$

where $\mathbf{H}_3$ and $\mathbf{T}_z$ are third component of $\mathbf{H}$ and $\mathbf{T}$ respectively. By substituting $\frac{Z}{Z'}$, we obtain

$$
\begin{aligned}
\mathbf{p} &= \frac{\mathbf{H}\mathbf{p}' + \frac{h}{h_c Z'}\mathbf{t}}{\mathbf{H}_3\mathbf{p}' + \frac{h\mathbf{T}_z}{h_c Z'}} \\
&= \frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} - \frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} + \frac{\mathbf{H}\mathbf{p}' + \frac{h}{h_c Z'}\mathbf{t}}{\mathbf{H}_3\mathbf{p}' + \frac{h\mathbf{T}_z}{h_c Z'}} \\
&= \frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} - \frac{\frac{h\mathbf{T}_z}{h_c Z'}}{\mathbf{H}_3\mathbf{p}' + \frac{h\mathbf{T}_z}{h_c Z'}}\frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} + \frac{\frac{h}{h_c Z'}}{\mathbf{H}_3\mathbf{p}' + \frac{h\mathbf{T}_z}{h_c Z'}}\mathbf{t} \\
&= \frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} - \frac{h\mathbf{T}_z}{Z h_c}\frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} + \frac{h}{Z h_c}\mathbf{t}.
\end{aligned}
\tag{11}
$$

When $\mathbf{T}_z = 0$, Eqn. 11 becomes

$$
\mathbf{p} = \frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} + \frac{h}{Z h_c}\mathbf{t}. \tag{12}
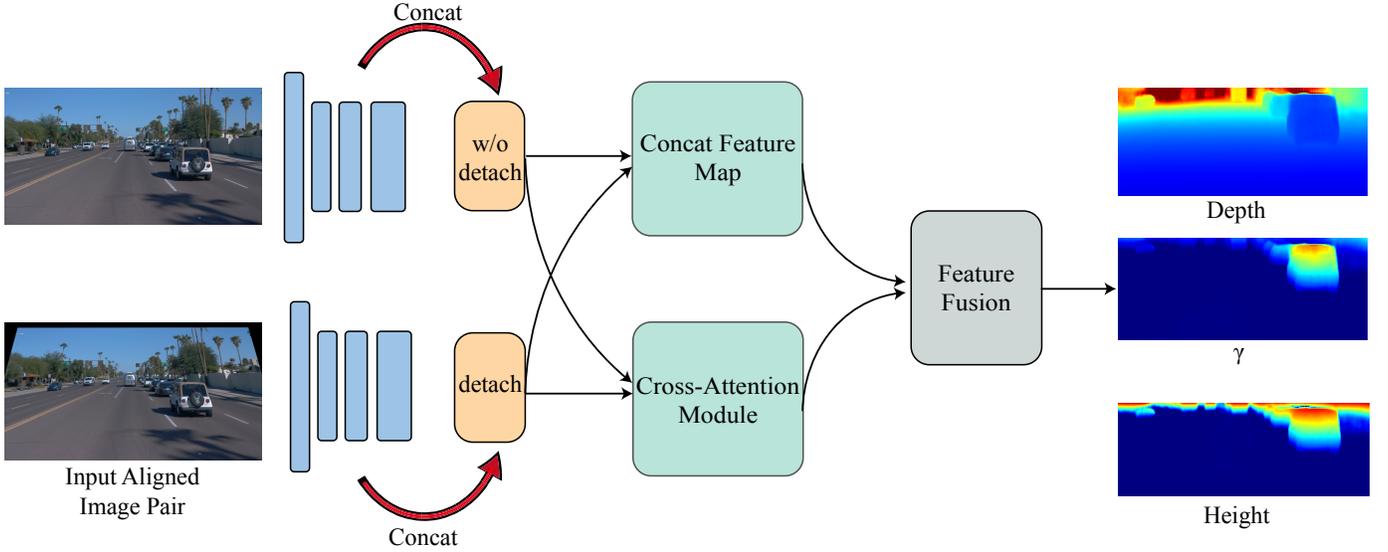$$

Fig. 2. The proposed Road Planar Parallax Attention Neural Network.

When $\mathbf{T}_z \neq 0$, we obtain

$$\mathbf{p} = \frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} - \frac{h\mathbf{T}_z}{Zh_c}\left(\frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} - \mathbf{e}\right), \tag{13}$$

where $\mathbf{e} = \frac{1}{\mathbf{T}_z}\mathbf{t}$. Given $\frac{h}{Z} = \gamma$ and $\frac{\mathbf{H}\mathbf{p}'}{\mathbf{H}_3\mathbf{p}'} = \mathbf{p}^w$, Eqn. 13 could be finally converted to

$$\mathbf{p} = \left(1 - \frac{h\mathbf{T}_z}{Zh_c}\right)\mathbf{p}^w + \frac{h\mathbf{T}_z}{Zh_c}\mathbf{e}, \tag{14}$$

$$\left(1 - \frac{h\mathbf{T}_z}{Zh_c}\right)\mathbf{p} = \left(1 - \frac{h\mathbf{T}_z}{Zh_c}\right)\mathbf{p}^w + \frac{h\mathbf{T}_z}{Zh_c}\mathbf{e} - \frac{h\mathbf{T}_z}{Zh_c}\mathbf{p}, \tag{15}$$

$$\mathbf{p} - \mathbf{p}^w = \frac{-\frac{h\mathbf{T}_z}{Zh_c}}{1 - \frac{h\mathbf{T}_z}{Zh_c}}(\mathbf{p} - \mathbf{e}), \tag{16}$$

$$\mathbf{p} - \mathbf{p}^w = \frac{-\gamma\frac{\mathbf{T}_z}{h_c}}{1 - \gamma\frac{\mathbf{T}_z}{h_c}}(\mathbf{p} - \mathbf{e}). \tag{17}$$

Applying the definition in Eqn. 2, we can calculate $\mathbf{u}_{res}$ from $\gamma$ by

$$\mathbf{u}_{res} = \frac{-\gamma\frac{\mathbf{T}_z}{h_c}}{1 - \gamma\frac{\mathbf{T}_z}{h_c}}(\mathbf{p} - \mathbf{e}). \tag{18}$$

From Eqn. 18, we know that the planar parallax can be easily gotten when the height of camera, the translation along z-axis, and $\gamma$ are available. In our framework, the $\gamma$ is estimated from the neural network, while others are from sensors or calibration. The benefit of the planar parallax is twofold. First, the homography perfectly depicts the underlying geometry of the road plane. After warping the source image $I_s$ via the road homography, pixels in road region would be aligned strictly with the target image $I_t$. In comparison, pixels in non-road regions would be affected by distortion of various degrees which is related to the height to the road plane. The distortion can provide vital cues for 3D reconstruction of the scene. Second, the homography can also remove the effect of rotation, which makes our method more robust to small baseline motion.

*b) Road 3D Geometry Recovery:* Although our proposed method only uses 2D flow-based transformations to build the training loss, we can also perform 3D reconstruction from $\gamma$ (see Fig. 1b for details). Similar to Eqn. 4, we have

$$h = h_c - \vec{\mathbf{N}}^T\mathbf{P}. \tag{19}$$

Supposing $\mathbf{K}$ is the camera's intrinsic matrix, $\mathbf{p}$ is the projection of $\mathbf{P}$ on image plane, whose homogeneous coordinate denoted by $\mathbf{p} = (x, y, 1)$, $\mathbf{P}$ could be calculated by an inverse projection given as

$$\mathbf{P} = Z\mathbf{K}^{-1}\mathbf{p}. \tag{20}$$

Substituting Eqn. 20 into Eqn. 19 we could get

$$h = h_c - \vec{\mathbf{N}}^T(\mathbf{K}^{-1}Z\mathbf{p}). \tag{21}$$

Dividing both sides of Eqn. 21 by $Z$ gives

$$\frac{h}{Z} = \frac{h_c}{Z} - \vec{\mathbf{N}}^T(\mathbf{K}^{-1}\mathbf{p}). \tag{22}$$

Defining $\gamma = \frac{h}{Z}$, Eqn. 22 could be finally reorganized as

$$Z = \frac{h_c}{\gamma + \vec{\mathbf{N}}^T(\mathbf{K}^{-1}\mathbf{p})}. \tag{23}$$

Note that the height of a pixel can be easily calculated by

$$h_p = \gamma Z. \tag{24}$$

The above formulas theoretically proves that the planar parallax estimated by the deep neural network can be directly converted into height and depth of each pixel. In the experiment part, we will convert the planar parallax estimated by deep neural network into depth and height to verify our proposed method.
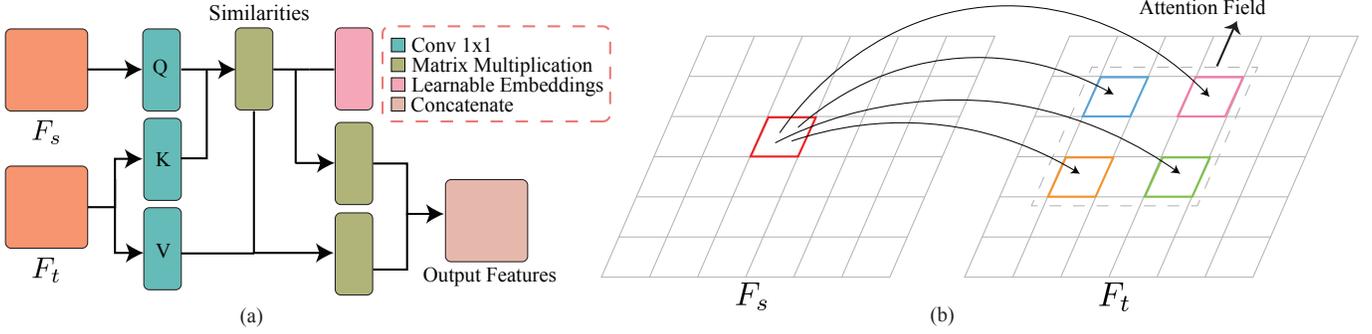
Fig. 3. (a) Proposed cross-attention module. (b) Illustration of cross-attention.

## B. Network Architecture

Planar parallax estimation requires perceiving the image displacements between two images. Based on this observation, in this section, we introduce RPANet, which consists of three modules: a CNN-based feature extraction layer, a cross-attention module, and an output layer. As shown in Figure. 2, the input of RPANet includes two consecutive images, one of which is warped by road homography, and the output of RPANet is the $\gamma$ map. The feature extraction module extracts features from both input images with shared weights. Similar to [47], to avoid collapsing, we apply a stop-gradient operation on the feature extracted from warped image $I_s^w$. Then, the extracted features are fed into the cross-attention module, through which the key clues between the features can be modeled more effectively. Finally, a feature fusion module is applied to predict $\gamma$.

As illustrated in Fig. 3(b), we adopt a cross-attention module. The cross-attention module tries to extract key clues for $\gamma$ estimation by performing neighborhood matching on two feature maps. We call the neighborhood area the attention field, which represents the field where one pixel of the feature map can be matched in the other feature map. Inspired by [44], we add a series of learnable parameters in order to make the network learn the implicit matching relationship efficiently.

As depicted in Fig. 3(a), in our proposed cross-attention module, 1x1 convolution is firstly performed to extract $Q$, $K$, and $V$, which were described in [36]. The feature maps $F_s$ extracted from the source image are fed to a $1 \times 1$ convolution to produce $Q$. $K$ along with $V$ are extracted from $F_t$ which are feature maps extracted from the target image through another two $1 \times 1$ convolution modules. Then, partial matrix multiplication is conducted between $Q$ and $K$ where one pixel in $Q$ has a set of corresponding pixels in a fixed attention field of $K$. A set of learnable parameters is utilized to find the matching relationship between the pixels from $F_s$ and $F_t$. Those parameters would be multiplied to the similarities generated from the result of partial matrix multiplication between $Q$ and $K$. A similar process of partial matrix multiplication would be conducted between the similarities and $V$. For every pixel of output,

$$y_o = \sum_{p \in A_{k \times k}} softmax(q_o^T * k_p)(v_p + r_{o,p}), \quad (25)$$

where $y_o$ is one pixel of the output $O$, and $A_{k \times k}$ refers to the attention field. $q_i$, $k_i$, and $v_i$ are vectors of a specific pixel from $Q$, $K$, and $V$ respectively, their shapes depend on the dimensions on channel. The $r_{a,b}$ is the learnable parameters represent a flow from $a$ to $b$ in two feature maps respectively. The partial matrix multiplication method is implemented with einsum and can be found in our implementation for details. Different from previous works [45], [48] which build attention module based on the epipolar constraint, our cross-attention module is based on the local geometry of planar parallax.

By applying a dilated attention on the 1/2 downsampled feature maps and $19 \times 19$ attention field, we avoid setting the attention field dense and global areas to reduce the number of parameters, which is similar to the dilated convolution [49].

## C. Loss Function

As the output $\gamma$ can be used to reconstruct a warped target image $I_t'$ from $I_s$, the widely used photometric loss can naturally be applied as supervisory signals. Besides, we use sparse ground truth $\gamma^*$ to build a sparse loss. Considering the photometric loss is not informative when applied on low-texture or homogeneous regions, we introduce additional smoothness loss to regularize our output. The total loss is given by:

$$E_{total} = \lambda_s E_s + \lambda_p E_p + \lambda_{sm} E_{sm}, \quad (26)$$

where $E_s$ is sparse loss for $\gamma$ map, $E_p$ is photometric loss function, $E_{sm}$ is smoothness loss function. $\lambda_s$, $\lambda_p$ and $\lambda_{sm}$ are the loss weights on the respective loss term.

*a) Photometric Loss Function:* After obtaining $\gamma$ map, we can calculate the correspondence of each pixel between the source image $I_s$ and the target image $I_t$. Given the $\mathbf{u}_{res}$ of a pixel $\mathbf{p}_s$ in $I_s$, we can get its corresponding pixel in $I_t$ as

$$\mathbf{p}_t' = \mathbf{u_p} + \mathbf{H}_{s \to t} * \mathbf{p}_s$$
$$= \mathbf{u_p} + \mathbf{p}_s^w. \quad (27)$$

Based on the above equation, frame $I_{t'}$ can be reconstructed as

$$I_{t'}[\mathbf{p}_t] = I_s \langle \mathbf{p}_s \rangle, \quad (28)$$

where $I_{t'}[\mathbf{p}_t]$ are pixel intensities at position $\mathbf{p}_t$, and $\langle \rangle$ is a bilinear sampling operator. Accordingly, the photometric loss function can be constructed to measure the difference between

Fig. 4. The Waymo Open Dataset used to construct RP2-Waymo dataset. The points in the figures are the LiDAR used for constructing the dataset. The blue points are the the points labeled with road. The road plane is extracted based on these points. Among all of the images, despite time, weather, scene vary, the road plane can be considered as a reference for planar parallax estimation.

$I_t$ and $I'_t$. We use the robust photometric error combining SSIM [50] and L1 norm between two images which is given by

$$E_p(I_t, I_{t'}) = \alpha \frac{1 - SSIM(I_t, I_{t'})}{2} + (1 - \alpha)||I_t - I_{t'}||, \quad (29)$$

where $\alpha$ is a hyper-parameter.

*b) Sparse Loss Function:* Since the ground truth $\gamma^*$ can be generated from sparse LiDAR points, we directly use it to train the network. The sparse loss can be defined as

$$E_s = \sum_{\mathbf{p} \in \Omega^l} |\gamma_{\mathbf{p}} - \gamma_{\mathbf{p}}^*|, \quad (30)$$

where $\gamma_{\mathbf{p}}$ and $\gamma_{\mathbf{p}}^*$ are the predicted and ground truth $\gamma$ value for pixel $\mathbf{p}$ respectively. $\Omega$ is the union of pixels that have ground truth.

*c) Smoothness Loss Function:* Edge-aware smoothness loss [31], [51] is widely used by existing methods to enhance the depth's local consistency. Different from these methods, we apply the second order smoothness constraint [52], [53] on the residual flow. The smoothness loss function is defined as

$$E_{sm} = \sum_{\mathbf{d}} \sum_{\mathbf{p}} \left( \left|\nabla^{\mathbf{d}}\mathbf{u}_{res}(\mathbf{p})\right|^2 e^{-\beta|\nabla^{\mathbf{d}}I_t(\mathbf{p})|} \right), \quad (31)$$

where $\nabla^{\mathbf{d}}$ stands for gradient calculated along the direction $\mathbf{d}$, $\beta$ is the weight for the gradient of image $I_t$ and $e$ is the natural base. In our method, we calculate the gradient of $\mathbf{u}_{res}$ and $I_t$ in both horizontal and vertical direction. By leveraging the smoothness loss, the collinearity of neighboring flows is improved, and hence the final depth and height are regularized.

## IV. DATASETS

Since there is no existing datasets dedicated to the planar parallax estimation task, we build a dataset named RP2-Waymo by carefully selecting data from the Waymo Open Dataset [22] and calculating the homography matrix. The RP2-Waymo dataset contains 13,030 training samples and 1,287

validation samples, which is challenging as it contains various scenes such as city, highway, suburb, and different weather conditions. To ensure fairness, we sample data uniformly from different sequences. The ground plane is extracted from the point cloud via robust algorithms such as RANSAC [54]. Combined with odometry provided by the Waymo Open Dataset, the homography matrix needed by RPANet can be easily calculated.

**Training Set.** In the training set, we take full advantage of the LiDAR points to calculate the homography matrix, the road plane, and the $\gamma$ numbers of the pixels that are available. Each sample consists of two consecutive images, one is aligned by homography matrix, which we call the source image, and the other is called target image.

**Validation Set.** In the validation set, we create two modes to evaluate the proposed method thoroughly in order to measure the performance of different use cases in the real world.

a). The road plane is available (PA). In this setting, we use LiDAR to construct not only the ground truth of $\gamma$, depth, and height, but also the road plane and homography matrix. We evaluate our methods on this setting because this setting can help us evaluate the deep neural network excluding the errors caused by the homography matrix and road plane calibration as much as possible. In the real world, the homography matrix and road plane can be calibrated with sensors inside the car. We leave the analysis of errors caused by these sensors for future work, since the sensors are not available in the Waymo Open Dataset now.

b). The road plane is not avaliable (PNA). In this setting, we only use LiDAR to construct the ground truth of $\gamma$, depth, and height for evaluation. In other words, the 3D reconstruction does not require LiDAR at all. In order to generate the road plane and homography matrix, we apply a homographynet with ResNet-18 [54] backbone and two head each contains 3 convolution layers. The homographynet takes raw images as input and output the road plane and pose, then the homography can be computed from them. When training the homographynet, we apply cosin ssimilarity loss for ground

TABLE I

**THE MEAN ABSOLUTE ERROR OF HEIGHT AND DEPTH.** DEPTH AND HEIGHT ARE OBTAINED FROM $\gamma$ ACCORDING TO EQN. 1. "GEOMETRY" REFERS TO USING THE TRADITIONAL GEOMETRY BASED METHOD. "DEPTH-BASELINE-R18" REPRESENTS A DEPTH ESTIMATION BASELINE WITH RESNET-18 BACKBONE (HEIGHT RESULTS ARE CALCULATED BY EQN. 24 WITH ROAD PLANE GROUND TRUTH). "GAMMANET-R18" IS A SIMPLE BASELINE THAT HAS THE SAME NETWORK STRUCTURE BUT PREDICT GAMMA SUPERVISED BY GROUND TRUTH. "W/O U" REFERS TO OUR RPANET WITHOUT THE UNSUPERVISED LOSS (EQN. 29). FOR THE VALIDATION SET, "PA" REFERS TO THAT ROAD PLANE IS AVAILABLE AND "PNA" INDICATES THAT ROAD PLANE IS NOT AVAILABLE. THE BEST RESULTS ARE SHOWN IN BOLD.

| Method | Validation Set | Absolute Height Error(m) | | | | Absolute Depth Error(m) | | |
|---|---|---|---|---|---|---|---|---|
| | | $h < 0.1m$ | $h < 0.3m$ | $h < 0.5m$ | $h < 1m$ | $d < 30m$ | $d < 50m$ | $d < 80m$ |
| **Geometry** | **PA** | 0.290 | 0.420 | 0.513 | 0.603 | 2.82 | 10.80 | 14.93 |
| **Depth-Baseline-R18** | **PA** | 0.057 | 0.066 | 0.069 | 0.077 | 0.545 | 0.947 | 1.352 |
| **GammaNet-R18** | **PA** | 0.021 | 0.034 | 0.041 | 0.056 | 0.388 | 0.775 | 1.200 |
| **RPANet w/o "U"** | **PA** | 0.022 | 0.034 | 0.041 | 0.053 | 0.358 | 0.703 | 1.162 |
| **RPANet** | **PA** | **0.019** | **0.031** | **0.038** | **0.051** | **0.337** | **0.702** | **1.140** |
| **RPANet** | **PNA** | 0.023 | 0.036 | 0.043 | 0.057 | 0.362 | 0.772 | 1.252 |

norm and photometric loss for homography matrix. We are surprised to find that a very simple network can get very effective results, although it will cause acceptable errors. More details can be found in Sec.V.

**Metrics.** We use the Mean Absolute Error (MAE) to evaluate the proposed methods on height and depth.

$$MAE = \frac{\sum_i^n |\hat{y}_i - y_i|}{n}. \tag{32}$$

In Eqn. 32, only pixels with ground truth are calculated in. The height and depth are also evaluated under different depth and height intervals to report the range where errors happen. Following [12], we also apply the widely used metrics in depth estimation evaluation to compare with methods for comparison. It is worth noting that height and $\gamma$ cannot be measured by relative errors, because the ground truth may have zero or negative values. For this reason, we do not apply relative metrics (e.g. the absolute relative error) on $\gamma$ and height evaluation.

## V. EXPERIMENTS

In this section, we evaluate the proposed method by comparing it with the depth estimation methods on our proposed datasets. After that, we give an analysis of the performance of these methods.

### A. Implementation Details

Our framework is implemented using Pytorch [55]. We adopt Adam optimizer [56] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ to update the parameters of RPANet. All the experiments are performed on a stand-alone server with 4 NVIDIA TITAN Xp GPUs. With the setting of corresponding default hyperparameters described above, each training step costs about 0.7 seconds. In the default settings, we use 3 GPUs with a batch size of 6 and each GPU needs about 8GB RAM though it has about 12GB RAM. Our model is trained for 20 epochs on the training set with the learning rate reduced by a factor of 10 after the tenth epoch where the initial learning rate is 0.0001. All hyper-parameters are tuned based

on a 1000-size development set, and we make no further adjustments except for the number of epochs for training on the complete training set. The input images are first cropped (from $1920 \times 1280$ to $1920 \times 1024$) and then resized to $960 \times 512$ using the Lanczos interpolation algorithm. We adopt a ResNet-18 backbone following [12].

**Baselines.** Since our method mainly focuses on the planar parallax geometry and the cross-attention module for finding the relationship between consecutive images, it can be adapted with different feature extractors. We adopt a ResNet-18 backbone following [12] to build RPANet. As there lacks planar parallax methods to be compared with, to verify the effectiveness of RPANet, we also build a baseline depth estimation with ResNet-18 backbone (denoted as "**Depth-Baseline-R18**") for fair comparison. Note that the Depth-Baseline-R18 is with full supervision rather than self supervision in [12]. To compare our method with recent transformer-based depth estimation method, we also build a "**Depth-Baseline-DPT**" beyond the DPT-Hybrid [57] feature extractor. To compare our method with a naive baseline that predicts gamma map with supervision, we use a "**GammaNet-R18**" with the same network structure as "**Depth-Baseline-R18**" but predicts gamma by supervision. We also provide a geometry-only method "**Geometry**" to compare RPANet with traditional methods. Specifically, we estimate the optical flow with OpenCV and calculate the depth and height with planar parallax geometry as described in Sec. III-A, which is similar to the practice in the traditional planar parallax methods [19], [23].

### B. Quantitative Results

**Ablation Study.** As shown in Table.I and Table.II, we conduct ablation study of our method on the RP2-Waymo dataset. The RPANet is trained on the training set in supervised or semi-supervised manner, and evaluated on the two validation sets including **PA** (road plane is available through LiDAR) and **PNA** (road plane is not available, which means LiDAR is not used during inference). As expected, our RPANet containing the proposed cross-attention module as well as trained with all the above loss functions, achieves the best results. With the

TABLE II
ABLATION STUDY ON THE PROPOSED RPANET. RPANET IS OUR PROPOSED METHOD. "W/O RE" MEANS THAT WITHOUT THE RELATIVE EMBEDDING IN OUR PROPOSED METHOD. "W/O DE" MEANS THAT WITHOUT THE DETACHMENT AFTER THE FEATURE OF THE WARPED IMAGE IN THE NETWORK. ALL OF THE RESULTS IN THIS TABLE ARE GOTTEN UNDER **PNA**.

| Method | Absolute Height Error(m) | | | | Absolute Depth Error(m) | | |
|---|---|---|---|---|---|---|---|
| | $h < 0.1m$ | $h < 0.3m$ | $h < 0.5m$ | $h < 1m$ | $d < 30m$ | $d < 50m$ | $d < 80m$ |
| **RPANet** | 0.019 | 0.031 | 0.038 | 0.051 | 0.337 | 0.702 | 1.140 |
| **RPANet w/o RE** | 0.039 | 0.050 | 0.056 | 0.067 | 0.475 | 0.856 | 1.297 |
| **RPANet w/o DE** | 0.030 | 0.041 | 0.047 | 0.059 | 0.405 | 0.788 | 1.228 |

TABLE III
THE MEAN ABSOLUTE ERROR OF HEIGHT AND DEPTH FOR **RPANET** AND **DEPTH-BASELINE-R18** IN DIFFERENT DEPTH AND HEIGHT INTERVALS. $h_0/d_0$ MEANS THAT IN THE SPECIFIC HEIGHT AND DEPTH RANGE, $h_0$ $m$ AND $d_0$ $m$ ARE THE MEAN ABSOLUTE ERROR OF HEIGHT AND DEPTH RESPECTIVELY.

| Depth \ Height | Method | $h < 0.1m$ | $h < 0.3m$ | $h < 0.5m$ | $h < 1m$ |
|---|---|---|---|---|---|
| $d < 30m$ | **RPANet** | 0.014/0.13 | 0.020/0.19 | 0.022/0.22 | 0.029/0.34 |
| | **Depth-Baseline-R18** | 0.049/0.39 | 0.053/0.45 | 0.054/0.47 | 0.058/0.55 |
| $d < 50m$ | **RPANet** | 0.017/0.20 | 0.026/0.36 | 0.031/0.45 | 0.041/0.70 |
| | **Depth-Baseline-R18** | 0.055/0.55 | 0.061/0.69 | 0.063/0.76 | 0.069/0.95 |
| $d < 80m$ | **RPANet** | 0.019/0.27 | 0.031/0.52 | 0.038/0.69 | 0.051/1.14 |
| | **Depth-Baseline-R18** | 0.057/0.64 | 0.066/0.88 | 0.069/1.02 | 0.077/1.35 |

TABLE IV
COMPARISON RESULTS OF OUR METHOD AND THE REPRESENTATIVE DEPTH ESTIMATION METHOD. THE DEFINITIONS OF METRICS ARE SAME AS [12]. "**RPANET + DPT**" REFERS TO A DPT-HYBRID [57] BACKBONE ADAPTED INTO OUR RPANET. THE METRICS WITH ORANGE BACKGROUND MEANS "**LOWER** IS BETTER". THE METRICS WITH BLUE BACKGROUND MEANS "**HIGHER** IS BETTER"

| method | Abs Rel | Sq Rel | RMSE | RMSE log | $\sigma < 1.25$ | $\sigma < 1.25^2$ | $\sigma < 1.25^3$ |
|---|---|---|---|---|---|---|---|
| **Depth-Baseline-R18** | 0.0474 | 0.3302 | 3.505 | 0.0876 | 0.970 | 0.992 | 0.997 |
| **RPANet** | 0.0378 | 0.4491 | 3.934 | 0.0896 | 0.964 | 0.989 | 0.996 |
| **Depth-Baseline-DPT** [57] | 0.0247 | 0.2573 | 2.110 | 0.0732 | 0.982 | 0.997 | 0.999 |
| **RPANet + DPT** | 0.0201 | 0.2660 | 2.209 | 0.0755 | 0.986 | 0.997 | 0.999 |

depth in the range of $0 - 80$ meters, our proposed RPANet has achieved a mean absolute error of **1.140** which is the best results among all of the methods. Note that even without road plane (the PNA setting), RPANet still outperform **Depth-Baseline-R18** with a large margin.

The results of different networks are reported in in Table.I. We can see that the traditional geometry based method is almost failed in our dataset and the accuracy of "**Depth-Baseline-R18**" is significantly lower than the others predicting $\gamma$ especially in height estimation, which validates the effectiveness of our proposed method. We can also notice that the results of network "**GammaNet-R18**" is comparative to that of network "**RPANet w/o U**" with respect to depth metric while worse in the height metric. The results indicate that the cross-attention module may provide more useful information for depth evaluation. Comparing the results of "**RPANet w/o U**" and "**RPANet**" in Table.I, we can conclude that the photometric loss reduces the MAE in all intervals regardless of height and depth. This is because that the photometric loss helps RPANet learn more accurate correspondences and supply complementary supervision for areas lack of ground truth.

In Table.II, we do the ablation studies on some strategies in our proposed RPANet. The results show that both the relative embedding in the cross-attention module and the detachment after the feature of the homography warped image are very useful for the network. The former is because that the proposed network highly relies on finding the displacement between homography aligned images to estimate the $\gamma$ map. The latter is to prevent the whole network collapses.

As shown in Table.III, we also report the detailed mean absolute error of height and depth in different intervals of depth and height for providing more information of our final setting. From Table. III, we can see that the results of both networks become poor for objects that are farther and higher. This is because the ground truth in the distance is more sparse, and the distance target is too small to estimate the matching relationship.

**Comparison between PA and PNA.** During training, we can use LiDAR to obtain the accuracy road plane and homography matrix, but during inference we may only have image sequences. To validate the utility of our method, we further compare the results of **PA** and **PNA** mentioned in

TABLE V
COMPARISON OF DIFFERENT SCENES IN ABSOLUTE HEIGHT ERROR AND ABSOLUTE DEPTH ERROR.

| Scene | Absolute Height Error(m) | | | | Absolute Depth Error(m) | | |
|---|---|---|---|---|---|---|---|
| | $h < 0.1m$ | $h < 0.3m$ | $h < 0.5m$ | $h < 1m$ | $d < 30m$ | $d < 50m$ | $d < 80m$ |
| city | 0.022 | 0.032 | 0.038 | 0.049 | 0.361 | 0.700 | 0.993 |
| suburb | 0.019 | 0.030 | 0.035 | 0.043 | 0.224 | 0.537 | 1.000 |
| highway | 0.012 | 0.018 | 0.024 | 0.041 | 0.299 | 0.571 | 1.005 |
| night | 0.027 | 0.056 | 0.074 | 0.122 | 1.003 | 1.873 | 2.469 |

TABLE VI
THE NUMBER OF EACH SCENE IN THE VALIDATION SET.

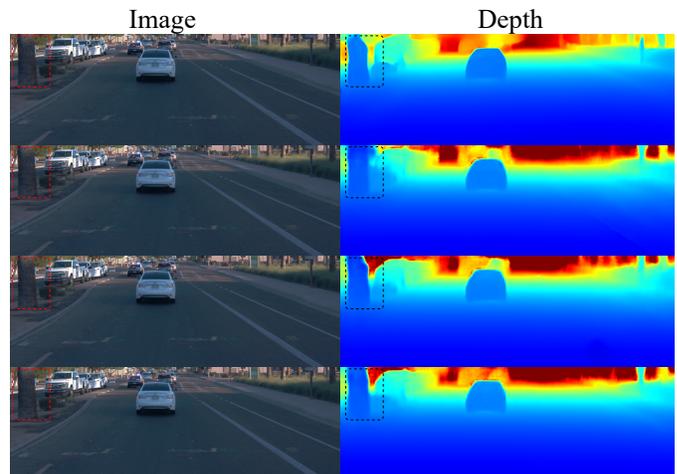| Scenes | city | suburb | highway | night |
|---|---|---|---|---|
| Number of Samples | 249 | 718 | 197 | 123 |



Fig. 5. Comparison on depth of different methods. Coloring from blue to red represents depth value from small to large. Rectangles depict regions for comparison. From top to bottom: (1).“**Depth-Baseline-R18**”; (2).“**GammaNet-R18**”; (3).“**RPANet w/o U**”; (4).“**RPANet**”.

Sec. IV. Comparing the results of height in **PNA** and **PA** of our method, **PNA** is worse than **PA** of all setting, but the gap is within an acceptable range. When comparing the results of depth, the **PNA** of “**RPANet**” is better than the result of “**Depth-Baseline-R18**”, this indicates that unsupervised loss and attention are more efficient for depth. Further more, even the plane and homography information supplied from a simple network, the result of **PNA** is better than “**Depth-Baseline-R18**” with a large margin.

**Comparison with Depth Estimation.** To fully evaluate the effectiveness of our method, comparative experiments with the depth estimation network are conducted and the results are reported in Table. IV. The ground truth of depth is generated from LiDAR points and used as sparse supervision. As reported in Table. IV, the proposed RPANet outperforms the depth estimation networks both with CNN (ResNet-18 [58]) and transformer (DPT-Hybrid [57]) backbone in absolute relative error but lags behind on square-based errors. We speculate that this may be due to the complexity of the scenes. In the complex scenes, our RPANet is affected by error of both plane estimation and gamma estimation at relatively far pixels. Considering the depth map is got from Eqn. 23, only a small error in the gamma and road plane will lead to a large error in depth map. For the far pixels, we would like to leave the more accurate depth estimation as future work. For example, one possible solution is combining the advantages of depth estimation and planar parallax and have better performance at both shorter distances and longer distances. But we still want to emphasize that depth estimation of pixels with a relatively close distance (e.g., less than 30m) may be more useful for driving scenes, and our RPANet has a pretty good performance in those pixels. This is due to the benefit brought by the geometric constraint of planar parallax, which makes the neural network learn to predict much easier.

**Comparison of Different Scenes.** To evaluate the performance of our method on different scenes, the validation set of RP2-Waymo dataset are clustered into 4 categories : **city**, **suburb**, **highway**, and **night**. The distribution of the validation set in the four scenarios and the original dataset

is kept consistent without deliberate adjustments. The number of samples in each scene can be found in Table VI. We test our proposed RPANet under all 4 scenes and the results are reported in Table. V. It can be seen from the data that in the night scene, the performance of our proposed RPANet has degraded. this is possibly because that the image quality at night is poor and there are fewer clues to recover the three-dimensional information from the image.

### C. Qualitative Results

The output of RPANet are represented in Fig. 6. The $\gamma$ is the direct result of the RPANet, while the depth and height are converted from $\gamma$ according to Eqn. 1. Colors from blue to red represents values from small to large. It can be seen from the $\gamma$ map that the $\gamma$ value of the nearby vehicle is relatively large, while the distant vehicle is relatively small, which is consistent with the definition of $\gamma$. Compared with depth information, $\gamma$ and height information can better distinguish obstacles above the road such as sidewalks, which to a certain extent is essential for autonomous or assisted driving. From the results, we can see that the drivable surface can be easily identified with the help of height and depth extracted from $\gamma$. **Comparison of Different Methods.** As shown in Fig. 5, we compare the different methods with the visualization results. The results show that RPANet could predict accurate depth,
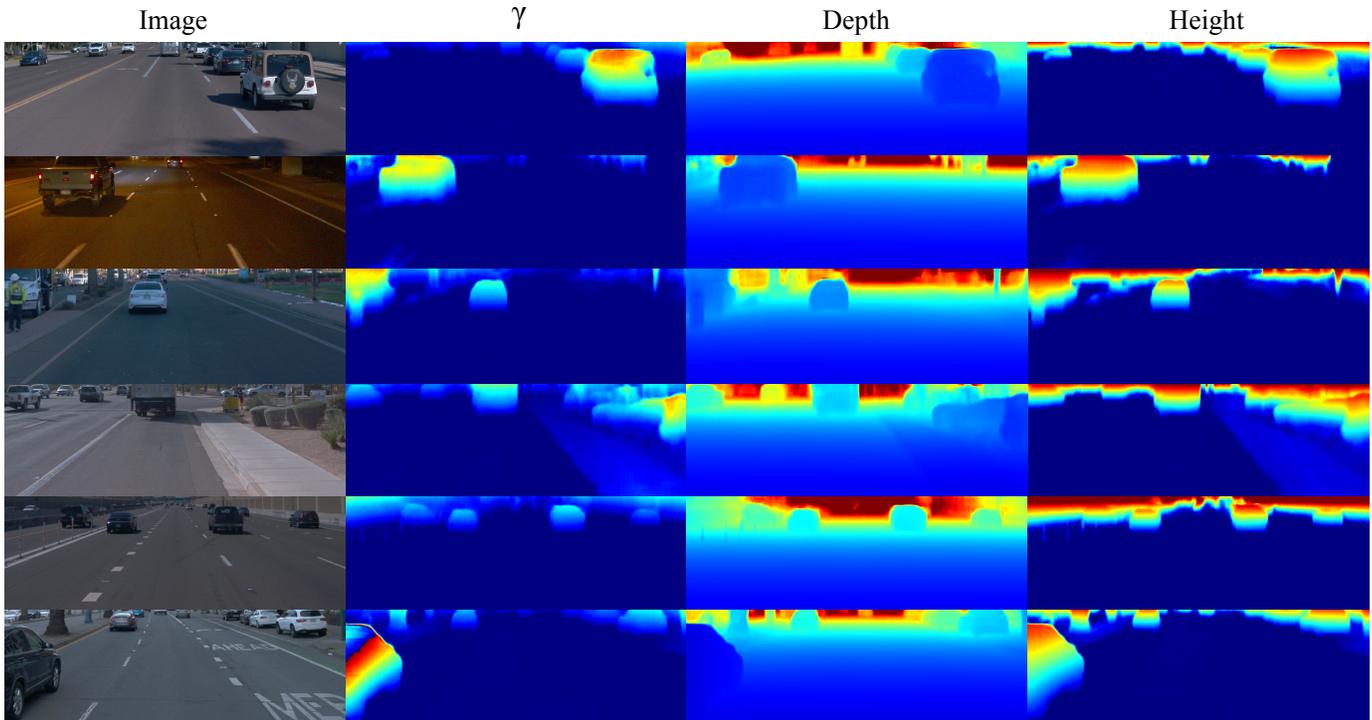
| Image | γ | Depth | Height |
|-------|---|-------|--------|



Fig. 6. The visualization results of γ, depth, and height of different inputs.



| (a) | (b) | (c) | (d) |
|-----|-----|-----|-----|

Fig. 7. (a) Source image $I_s$. (b) Homography warped source image $I_s^w$. (c) Reconstructed image $I_t'$. (d) Target image $I_t$. (a) to (b) refers to the homography aligning process; (b) to (c) refers to the planar parallax reconstruction process; (c) to (d) illustrates the difference between reconstructed image and the target image.

especially in discontinuous regions (see the boundaries of trees in rectangles). At the edge of the tree, the illumination of it is relatively low, making depth prediction difficult. Nonetheless, after applying the geometric constraint of planar parallax, the depth information can be well calculated leveraging γ. We take the tree besides the road as an example. Although the depth estimation network directly outputs the depth and thus avoiding the error caused by the road surface calibration or magnified by Eqn. 1, the depth estimation network does not identify the boundary of the tree successfully. "**GammaNet-R18**" does not distinguish between the tree and the road surface next to it, while RPANet distinguishes the tree from the pavement better. It shows that the cross-attention module and the photometric loss help the neural network to learn γ more easily.

**Image Reconstruction via Residual Flow.** As illustrated in Fig. 7, the road plane of source image and target image are aligned after homography warping, while other static areas can be further aligned by the residual flow warping. The

visualization shows that the two warping steps build a bridge between the source image and the target image, due to that we can easily obtain the reconstructed image and leverage photometric loss to train our RPANet. In Fig. 7, it is easy to find whether the proposed framework has successfully reconstructed the 3D scene with planar parallax. From (a) to (b), the road is well aligned with homography matrix while things above the road are with a displacement. From (b) to (c), things above the road are corrected with planar parallax estimated by neural network. The difference between (c) and (d) indicates that errors still exist, especially for the pixels with a large height or depth in the image. The well-aligned pixels from (c) to (d) indicate that their 3D structures are well estimated.

## VI. CONCLUSION

In this paper, we propose a planar parallax estimation method, which combines neural networks and planar parallax geometry. The input of our method is aligned image pairs via

road homography. The output $\gamma$ map is utilized to recover the 3D structure (depth and height). We also devise the cross-attention module to learn planar parallax more easily. Since no public dataset provides aligned images via road homography, we collect data from the Waymo Open Dataset and build the RP2-Waymo dataset. Comprehensive experiments conducted on the datasets valid the effectiveness of our method.

## REFERENCES

[1] T. Asai, K. Yamaguchi, Y. Kojima, T. Naito, and Y. Ninomiya, "3D line reconstruction of a road environment using an in-vehicle camera," in *Proc. ISVC*, 2008, pp. 897–904.

[2] D. Chen and X. He, "Fast automatic three-dimensional road model reconstruction based on mobile laser scanning system," *Optik*, vol. 126, no. 7-8, pp. 725–730, 2015.

[3] H. Gao, L. Liu, Y. Tian, and S. Lu, "3d reconstruction for road scene with obstacle detection feedback," *IJPRAI*, p. 1855021, 2018.

[4] C. Badue, R. Guidolini, R. V. Carneiro, P. Azevedo, V. B. Cardoso, A. Forechi, L. Jesus, R. Berriel, T. M. Paixao, F. Mutz *et al.*, "Self-driving cars: A survey," *ESWA*, p. 113816, 2020.

[5] B. Schoettle and M. Sivak, "A survey of public opinion about autonomous and self-driving vehicles in the us, the uk, and australia," University of Michigan, Ann Arbor, Transportation Research Institute, Tech. Rep., 2014.

[6] B. Lu, B. Tam, and N. Kottege, "Autonomous obstacle legipulation with a hexapod robot," *arXiv preprint arXiv:2011.06227*, 2020.

[7] H. Xu and J. Zhang, "Aanet: Adaptive aggregation network for efficient stereo matching," in *Proc. CVPR*, 2020, pp. 1959–1968.

[8] F. Zhang, V. Prisacariu, R. Yang, and P. H. Torr, "Ga-net: Guided aggregation net for end-to-end stereo matching," in *Proc. CVPR*, 2019, pp. 185–194.

[9] H. Fujita, M. Itagaki, K. Ichikawa, Y. K. Hooi, K. Kawano, and R. Yamamoto, "Fine-tuned pre-trained mask r-cnn models for surface object detection," *arXiv preprint arXiv:2010.11464*, 2020.

[10] J. Zhang and S. Singh, "Loam: Lidar odometry and mapping in real-time," in *Proc. RSS*, 2014, pp. 1–9.

[11] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proc. CVPR*, 2016, pp. 4104–4113.

[12] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, "Digging into self-supervised monocular depth estimation," in *Proc. ICCV*, 2019, pp. 3828–3838.

[13] Z. Yin and J. Shi, "Geonet: Unsupervised learning of dense depth, optical flow and camera pose," in *Proc. CVPR*, 2018, pp. 1983–1992.

[14] I. Tishchenko, S. Lombardi, M. R. Oswald, and M. Pollefeys, "Self-supervised learning of non-rigid residual flow and ego-motion," in *Proc. 3DV*, 2020, pp. 150–159.

[15] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Proc. NIPS*, 2014, pp. 2366–2374.

[16] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, "Unsupervised learning of depth and ego-motion from video," in *Proc. CVPR*, 2017, pp. 1851–1858.

[17] M. Irani and P. Anandan, "Parallax geometry of pairs of points for 3d scene analysis," in *Proc. ECCV*, 1996, pp. 17–30.

[18] M. Irani, B. Rousso, and S. Peleg, "Recovery of ego-motion using region alignment," *IEEE TPAMI*, vol. 19, no. 3, pp. 268–272, 1997.

[19] R. Kumar, P. Anandan, and K. Hanna, "Direct recovery of shape from multiple views: A parallax based approach," in *Proc. ICPR*, 1994, pp. 685–688.

[20] A. Shashua and N. Navab, "Relative affine structure: Theory and application to 3d reconstruction from perspective views," in *Proc. CVPR*, 1994, pp. 483–489.

[21] M. Irani, P. Anandan, and M. Cohen, "Direct recovery of planar-parallax from multiple frames," *IEEE TPAMI*, vol. 24, no. 11, pp. 1528–1534, 2002.

[22] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, Y. Zhang, J. Shlens, Z. Chen, and D. Anguelov, "Scalability in perception for autonomous driving: Waymo open dataset," in *Proc. CVPR*, 2020, pp. 2446–2454.

[23] H. S. Sawhney, "3d geometry from planar parallax," in *Proc. CVPR*, 1994, pp. 929–934.

[24] K. Chaney, A. Z. Zhu, and K. Daniilidis, "Learning event-based height from plane and parallax," in *Proc. IROS*, 2019, p. 3690–3696.

[25] H. Xing, Y. Cao, M. Biber, M. Zhou, and D. Burschka, "Joint prediction of monocular depth and structure using planar and parallax geometry," *PR*, vol. 130, p. 108806, 2022.

[26] H. Fan, H. Su, and L. J. Guibas, "A point set generation network for 3d object reconstruction from a single image," in *Proc. CVPR*, 2017, pp. 605–613.

[27] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese, "3d-r2n2: A unified approach for single and multi-view 3d object reconstruction," in *Proc. ECCV*, 2016, pp. 628–644.

[28] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang, "Pixel2mesh: Generating 3d mesh models from single rgb images," in *Proc. ECCV*, 2018, pp. 52–67.

[29] S. Saito, Z. Huang, R. Natsume, S. Morishima, A. Kanazawa, and H. Li, "Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization," in *Proc. ICCV*, 2019, pp. 2304–2314.

[30] R. Garg, V. K. Bg, G. Carneiro, and I. Reid, "Unsupervised cnn for single view depth estimation: Geometry to the rescue," in *Proc. ECCV*, 2016, pp. 740–756.

[31] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proc. CVPR*, 2017, pp. 270–279.

[32] J. Zbontar and Y. LeCun, "Computing the stereo matching cost with a convolutional neural network," in *Proc. CVPR*, 2015, pp. 1592–1599.

[33] A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, R. Kennedy, A. Bachrach, and A. Bry, "End-to-end learning of geometry and context for deep stereo regression," in *Proc. ICCV*, 2017, pp. 66–75.

[34] J.-R. Chang and Y.-S. Chen, "Pyramid stereo matching network," in *Proc. CVPR*, 2018, pp. 5410–5418.

[35] X. Guo, K. Yang, W. Yang, X. Wang, and H. Li, "Group-wise correlation stereo network," in *Proc. CVPR*, 2019, pp. 3273–3282.

[36] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. NIPS*, 2017, pp. 5998–6008.

[37] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu *et al.*, "A survey on visual transformer," *IEEE TPMAI*, vol. 45, no. 1, pp. 87–110, 2023.

[38] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. ICLR*, 2021, pp. 1–9.

[39] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Proc. ECCV*, 2020, pp. 213–229.

[40] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," in *Proc. ICLR*, 2021, pp. 1–9.

[41] N. Parmar, A. Vaswani, J. Uszkoreit, L. Kaiser, N. Shazeer, A. Ku, and D. Tran, "Image transformer," in *Proc. ICML*, 2018, pp. 4055–4064.

[42] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proc. ICCV*, 2021, pp. 10012–10022.

[43] J. Ho, N. Kalchbrenner, D. Weissenborn, and T. Salimans, "Axial attention in multidimensional transformers," *arXiv preprint arXiv:1912.12180*, 2019.

[44] L. Wang, Y. Guo, Y. Wang, Z. Liang, Z. Lin, J. Yang, and W. An, "Parallax attention for unsupervised stereo correspondence learning," *IEEE TPAMI*, vol. 44, no. 4, pp. 2108–2125, 2020.

[45] Z. Li, X. Liu, N. Drenkow, A. Ding, F. X. Creighton, R. H. Taylor, and M. Unberath, "Revisiting stereo depth estimation from a sequence-to-sequence perspective with transformers," in *Proc. ICCV*, 2021, pp. 6197–6206.

[46] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3D packing for self-supervised monocular depth estimation," in *Proc. CVPR*, 2020, pp. 2485–2494.

[47] X. Chen and K. He, "Exploring simple siamese representation learning," in *Proc. CVPR*, 2021, pp. 15750–15758.

[48] L. Wang, Y. Wang, Z. Liang, Z. Lin, J. Yang, W. An, and Y. Guo, "Learning parallax attention for stereo image super-resolution," in *Proc. CVPR*, 2019, pp. 12250–12259.

[49] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *Proc. ICLR*, 2016, pp. 1–9.

[50] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE TIP*, vol. 13, no. 4, pp. 600–612, 2004.

[51] J.-W. Bian, Z. Li, N. Wang, H. Zhan, C. Shen, M.-M. Cheng, and I. Reid, "Unsupervised scale-consistent depth and ego-motion learning from monocular video," in *Proc. NeurIPS*, 2019, pp. 35–45.

[52] S. Meister, J. Hur, and S. Roth, "Unflow: Unsupervised learning of optical flow with a bidirectional census loss," in *Proc. AAAI*, 2018, pp. 7251–7259.

[53] R. Jonschkowski, A. Stone, J. T. Barron, A. Gordon, K. Konolige, and A. Angelova, "What matters in unsupervised optical flow," in *Proc. ECCV*, 2020, pp. 557–572.

[54] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[55] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *Proc. NeurIPS*, 2019, pp. 8026–8037.

[56] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015, pp. 1–9.

[57] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," in *Proc. CVPR*, 2021, pp. 12 179–12 188.

[58] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016, pp. 770–778.