# MonoGlass3D: Monocular 3D Glass Detection with Plane Regression and Adaptive Feature Fusion

Kai Zhang, Guoyang Zhao, Jianxing Shi, Bonan Liu, Weiqing Qi, and Jun Ma, *Senior Member, IEEE*

*Abstract*—Detecting and localizing glass in 3D environments poses significant challenges for visual perception systems, as the optical properties of glass often hinder conventional sensors from accurately distinguishing glass surfaces. The lack of real-world datasets focused on glass objects further impedes progress in this field. To address this issue, we introduce a new dataset featuring a wide range of glass configurations with precise 3D annotations, collected from distinct real-world scenarios. On the basis of this dataset, we propose MonoGlass3D, a novel approach tailored for monocular 3D glass detection across diverse environments. To overcome the challenges posed by the ambiguous appearance and context diversity of glass, we propose an adaptive feature fusion module that empowers the network to effectively capture contextual information in varying conditions. Additionally, to exploit the distinct planar geometry of glass surfaces, we present a plane regression pipeline, which enables seamless integration of geometric properties within our framework. Extensive experiments demonstrate that our method outperforms state-of-the-art approaches in both glass segmentation and monocular glass depth estimation. Our results highlight the advantages of combining geometric and contextual cues for transparent surface understanding. The dataset and code will be released at https://github.com/Kai0139/MonoGlass3D.

*Index Terms*—Glass detection, plane regression, deep learning, robotics perception.
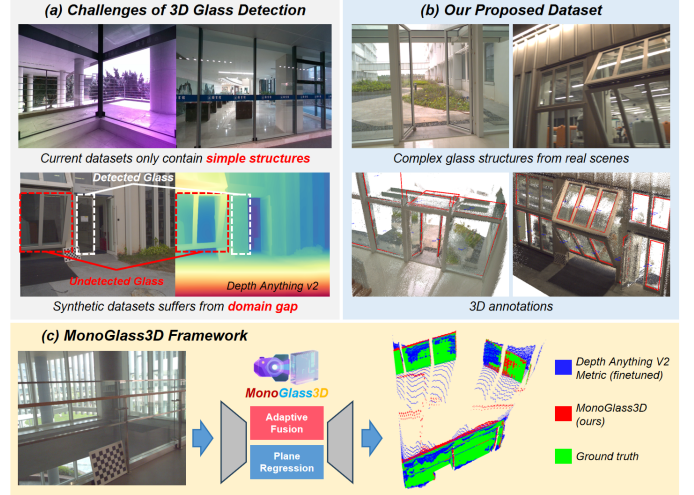


Fig. 1. **Challenges and our proposed method.** (a) Existing 3D glass dataset [11] only contains simple glass structures with one or two planes. Synthetic datasets suffers from domain gap problems, which makes the models hard to generalize in real-world scenarios, as depicted by depth map from Depth Anything V2 Metric [12] (pretrained on Hypersim [13]). (b) Our dataset is collected from complex glass structures in real-world. (c) Our MonoGlass3D network achieve superior performance in monocular 3D glass detection when compared to the baseline.

## I. INTRODUCTION

The detection and localization of glass surfaces represent a significant and enduring challenge in the field of visual perception [1], primarily attributable to the material's non-distinctive visual appearance and unique transmissive optical properties [2]. These inherent characteristics create a critical gap in environmental awareness for robotic systems. Given the ubiquitous presence of glass structures in modern civil environments, ranging from windows and doors to glass walls and facades, the ability to accurately perceive and localize glass surfaces has become increasingly crucial for the safe and reliable operation of autonomous platforms [3]. The inability to detect glass surfaces not only poses substantial safety risks, but also severely limits the operational capabilities of robots in both structured and unstructured environments, making this challenge a pressing issue in the advancement of robotics and autonomous systems.

In some of recent literatures, glass detection is approached as an image segmentation problem [2], [4], [5], they are capable of recognizing the presence of glass in images, but lacks crucial 3D information. Practical robotic applications such as path planning [6] and object manipulation [7], [8] demand a comprehensive 3D understanding of the environment. The primary challenge stems from glass's inherent transparency [9], [10]. Robots predominantly rely on cameras and LiDAR sensors, both of which allow visible light and laser beams to pass through glass, these surfaces often appear indistinguishable from the background or even 'invisible' to perception systems. As a result, detecting and localizing glass in 3D space is particularly difficult. Unlike 2D segmentation, which only requires identifying the shape of glass in images, 3D glass detection must also infer spatial position and depth, making the problem substantially more challenging.

Some research efforts consider 3D glass detection as an image based depth estimation problem. GWDepth [11] collects RGB-D data for training a network to estimate the depth of glass walls in civil settings. More recently, Depth Anything V2 [12] introduces a vision transformer based model that is

Kai Zhang, Guoyang Zhao, Jianxing Shi, and Weiqing Qi are with the Robotics and Autonomous Systems Thrust, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou 511453, China (e-mail: kzhang740@connect.hkust-gz.edu.cn; gzhao492@connect.hkust-gz.edu.cn; jshi627@connect.hkust-gz.edu.cn; wqiad@connect.hkust-gz.edu.cn).

Bonan Liu is with the Computational Media and Arts Thrust, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou 511453, China (e-mail: bliu404@connect.hkust-gz.edu.cn).

Jun Ma is with the Robotics and Autonomous Systems Thrust, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou 511453, China, and also with the Division of Emerging Interdisciplinary Areas, The Hong Kong University of Science and Technology, Hong Kong SAR, China (e-mail: jun.ma@ust.hk).

capable of estimating the depth of glass surfaces. However, the aforementioned approaches have several limitations. First, in 2D glass segmentation, optical features such as blurriness [14], reflections [11], and boundaries [4], [5] are often utilized to extract contextual information; however, optical features often rely on lighting and may not always be present, while boundary contexts appear as thin lines and are inherently difficult for neural networks to learn. Second, the prior knowledge that glass surfaces are generally planar is not leveraged, overlooking a key geometric constraint that can regulate 3D glass detection. Third, existing 3D glass datasets have notable shortcomings: GWDepth [11] only includes scenes with one or two glass planes (which are either coplanar or perpendicular, as illustrated in Fig. 1(a)), and Depth Anything V2 is trained on synthetic data [13] only, which introduces domain gap issues and hinders generalization to real-world scenes (Fig. 1(a)).

While our work shares a similar objective with previous studies [11], [12], we take a fundamentally different approach. We introduce a new real-world glass dataset and Mono-Glass3D, a novel 3D glass detection network. To overcome the limitations of existing datasets and to facilitate efficient collection and annotation of 3D glass surfaces, we build a new real-world glass dataset (Fig. 1(b)) along with a novel data acquisition pipeline. Motivated by the need for more diverse and complex glass configurations, our pipeline leverages a LiDAR-Visual mapping system to produce 3D point clouds containing various glass structures. Glass surfaces are jointly labeled in both 2D images and 3D point clouds. For glass detection, rather than designating specific contextual features, we introduce an adaptive feature fusion module inspired by the centerness concept from [15], to enhance the network's adaptability to the diverse shapes present in regions surrounding glass surfaces. Moreover, instead of estimating depth values, our method focuses on regressing plane parameters for glass regions. By reformulating depth estimation as plane parameter regression, we are able to leverage the inherent geometric properties of planes to regulate 3D estimations, resulting in significantly improved 3D glass detection performance, as shown in Fig. 1(c). Extensive experiments on both glass segmentation and glass depth estimation tasks show that our approach achieves superior performance compared to state-of-the-art methods.

In summary, the contribution of this paper is fourfold:

- We present a new 3D glass dataset capturing diverse glass structures and lighting conditions in real-world scenarios, annotated with segmentation masks, 3D plane parameters, and depth maps.
- We propose an adaptive feature fusion module that leverage centerness maps to enhance the network's adaptability to varying glass contexts.
- We develop a plane regression pipeline that formulates 3D glass detection into a 3D plane parameter prediction task, enabling the network to benefit from underlying geometric properties.
- Our method demonstrates noteworthy performance in both glass segmentation and depth estimation, outperforming existing methods with a lightweight framework.

## II. RELATED WORKS

### A. 2D Glass Detection

Image-based glass detection raises unique challenges compared to general object detection [16] due to the transparent nature of glass elements. Mei et al. [2] pioneers this field by establishing the first benchmark for the glass segmentation problem through introducing GDD, the first glass detection dataset, and GDNet, a glass detection network that identifies glass surfaces by extracting multi-scale contextual features. Although glass surfaces exhibit non-static and ambiguous appearances, certain distinctive features can still be leveraged to detect their presence. For instance, some works [4], [5] explicitly use glass boundaries to supervise networks in learning boundary context, while others [14], [17], [18] extract glass reflections to enhance the network's capability of recognizing glass. In addition, semantic labels have been utilized to guide networks in learning glass context features associated with semantic segmentation [12], [19]. To better distinguish glass from the background, different sensor modalities have been incorporated in previous works. Polarization imaging sensors, which capture the linear polarization states of pixels, have been employed in [20], [21], and thermal imaging sensors have also been explored [22]. However, these alternative sensors are not as accessible as RGB cameras, which hinders their application in realistic scenarios.

### B. 3D Glass Detection

Unlike glass segmentation in images, it is much harder to accurately measure the position of glass in 3D space. Common 3D sensors, such as depth cameras and LiDARs, both rely on measuring the time of flight of transmitted infrared (IR) signals. However, similar to visible light, IR signals often penetrate the glass surface and reflected by the object behind glass, the reflected signal is also unreliable because of light refraction. To avoid these difficulties, synthetic data becomes a practical and efficient approach to generate annotated 3D data for glass. Synthetic dataset is employed in [10], [23] to train depth estimation network for small transparent objects, Depth Anything V2 [12] proposes a synthetic-to-real knowledge transfer pipeline with a teacher model pre-trained on Hypersim [13], a large scale synthetic dataset featuring hundreds of indoors scenes with glass surfaces.

The shape distribution of glass surface in real world is highly long tailed, with the majority being flat surfaces. For flat glass, dense depth maps can be interpolated from a few points. GWDepth [11] uses a depth camera to capture the depth values at glass boundaries and and interpolate pixel-wise glass depths. However, this requires manual annotation of reliable pixels in each frame. Additionally, while the interpolability of flat glass surfaces is a valuable geometric constraint for glass detection, their approach does not fully leverage this property.

### C. 3D Perception with Planar Constraints

Among geometric primitives, planes are especially useful as they regulate unconstrained 3D points to lie on a 2D surface. Recent advances in deep learning have enabled significant progress in 3D perception leveraging planar constraints.
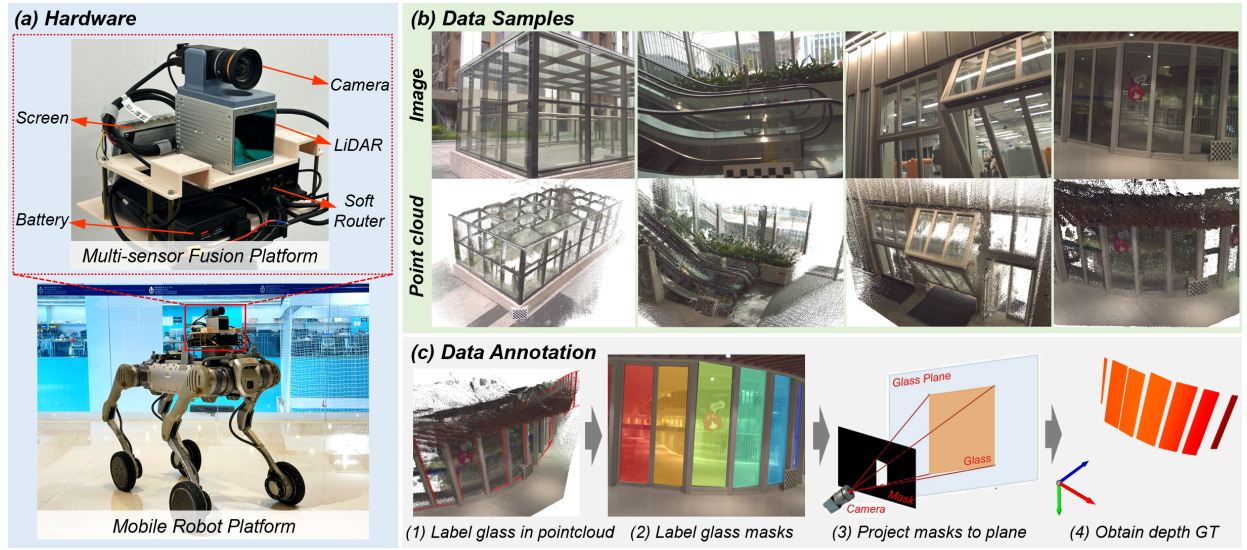
Fig. 2. **Proposed dataset construction procedures.** (a) Data collection device, we integrated a LiDAR and a camera with a wheel-legged robot. (b) Collected scenes with complex glass structures. (c) 3D glass data annotation procedures.

PlaneNet [24] introduces the first end-to-end neural network for piecewise planar reconstruction from a single RGB image, they estimate plane parameters and segmentation masks to yield structured depth maps with high accuracy. Similarly, PlaneRecover [25] proposes a convolutional neural network framework that exploits a novel plane structure-induced loss, allowing the network to simultaneously predict planar segmentation without requiring explicit plane annotations. Building upon these ideas, PlaneRCNN [26] adapts a Mask R-CNN architecture for plane detection and jointly refines segmentation masks by enforcing multi-view consistency during training. Moving beyond explicit plane detection, P3Depth [27] introduces a piecewise planarity prior to monocular depth estimation, learning to aggregate information from coplanar pixel groups through a two-headed network structure that adaptively fuses predictions, resulting in depth maps with sharp geometric boundaries and leading performance on standard benchmarks. These advances collectively demonstrate the effectiveness of integrating planar priors into deep networks for improved 3D perception from single RGB images. Most of the related works treat plane parameter estimation as an auxiliary task to improve object detection or depth estimation. In contrast, our approach is designed to predict the plane parameters of glass surfaces, which tackles 3D glass detection problem from a new perspective.

## III. MONOCULAR 3D GLASS DETECTION DATASET

### A. Hardware Setup

Compared to depth cameras [10], [11], LiDAR systems generally provide more accurate distance measurements. However, conventional LiDAR produces inherently sparse data and does not deliver depth values for every pixel within the Field of View (FoV), limiting its utility for dense 3D reconstruction of transparent surfaces such as glass. To address this, we select the Livox-Avia, a compact, non-repetitive scanning LiDAR. Unlike repetitive scanning LiDARs, which produce a fixed point cloud pattern, non-repetitive scanning allows point cloud density to accumulate over time as the sensor moves slowly during data acquisition, resulting in more comprehensive 3D coverage of glass surfaces. We configure the LiDAR to operate in single-first return mode to ensure stable and reliable mapping, as multi-return modes are usually incompatible with LiDAR-Visual SLAM systems.

For the RGB camera, we select the MV-CU013-80GC global shutter camera from HIKROBOT, which captures images at resolution of $1280 \times 1024$. A $6\,\mathrm{mm}$ lens is deliberately chosen to ensure the camera's FoV is slightly smaller than that of the LiDAR. This design guarantees that all glass surfaces captured in the images can be reliably mapped to the corresponding 3D LiDAR point clouds, thereby ensuring consistent cross-modal annotation.

To optimize image quality under varying lighting conditions, we tailor the camera settings to different environments. In indoor environments, automatic exposure is enabled to maintain uniform brightness across frames. In outdoor environments, where lighting can fluctuate rapidly, we disable automatic exposure and instead dynamically adjust the exposure time in real-time according to brightness changes of adjacent frames, effectively preventing overexposure and ensuring high-quality image capture regardless of ambient brightness.

For sensor integration, we adopt the configuration in FAST-LIVO2 [28], as illustrated in Fig. 2(a). Unlike FAST-LIVO2 that synchronizes sensors using external PPS signals, we use a soft router with multiple Ethernet interfaces that support the IEEE1588 Precision Time Protocol (PTP). This configuration enables precise timestamp synchronization between image frames and LiDAR scans through a PTP master device, without requiring any external signal sources. Extrinsic calibration between the LiDAR and camera is performed using a chessboard target, wherein the transformation matrix is computed by matching chessboard corners detected from LiDAR intensity images and RGB images, ensuring accurate spatial correspondence between modalities.

To further enrich the diversity and realism of our dataset,

we mount the synchronized sensor suite on the wheel-legged Unitree B2-W robot, which offers high mobility and stability over challenging terrains, such as stairs and curbs. This hardware integration enables data collection in both handheld and robot-carried modes, allowing for comprehensive coverage of complex environments and facilitating the construction of a realistic and versatile 3D glass dataset.

### B. Data Collection

In this work, we use FAST-LIVO2 [28] to produce 3D scans of glass and surrounding environments. FAST-LIVO2 is a state-of-the-art open sourced multi-modal SLAM system, which employs sensor fusion techniques to combine LiDAR, RGB camera, and IMU data to estimate the camera poses while generating a colored point cloud map. To ensure a wide diversity of glass geometries, we capture data in a variety of locations, including corridors, cafes, libraries, offices, glass doors, glass walls, ceilings, escalators, etc. In total, we scan 50 distinct scenes featuring a broad range of glass configurations, as exemplified in Fig. 2(b).

### C. Data Annotation

The data annotation mainly consists of 4 steps, as shown in Fig. 2(c). First, for each scene, after we obtain the point cloud map, we annotate the 3D positions of glass surfaces in the point cloud using thin bounding boxes. The vertices of bounding boxes are then transformed to the camera coordinate frame by applying spatial transformations:

$$P_c = T_{cw}P_w \tag{1}$$

where $P_c$ and $P_w$ denote the bounding box vertex coordinates in camera and world frames, respectively, $T_{cw}$ is the transformation matrix mapping world coordinates to camera frame.

Then, we fit the 3D bounding boxes to planes. A plane in 3D can be described by 4 variables characterized by:

$$\mathbf{n} \cdot \mathbf{v} + d = 0 \tag{2}$$

where $\mathbf{n}$ is the 3D plane normal vector (contains 3 variables), and $\mathbf{v}$ is arbitrary point on the plane, $d$ is the interception, with magnitude equals to the distance from origin to the plane. Given 3D positions of bounding box vertices $P_c$, we can fit the vertices to a set of plane parameters using least squares approximation:

$$\hat{\mathbf{n}} = (P_c^T P_c)^{-1} P_c^T \left[-1\right]^{3 \times 1}$$
$$\mathbf{n} = \frac{\hat{\mathbf{n}}}{\|\hat{\mathbf{n}}\|}, \quad d = \frac{1}{\|\hat{\mathbf{n}}\|} \tag{3}$$

After glass instances are registered as 3D planes, we annotate glass in images. Segmentation masks for glass surfaces are annotated on selected image frames from each scan.

Finally, after we obtain 2D masks and their corresponding 3D plane parameters, the depth value for each pixel within the mask is computed by projecting each 2D mask to its associated 3D plane:

$$\text{depth} = d \left(\mathbf{n}^T K^{-1} \begin{bmatrix} u_x & u_y & 1 \end{bmatrix}^T\right)^{-1} \tag{4}$$

where $K$ is the camera intrinsic matrix, $u_x$ and $u_y$ are the pixel coordinate. This projection process is depicted as the last two steps of Fig. 2(c).
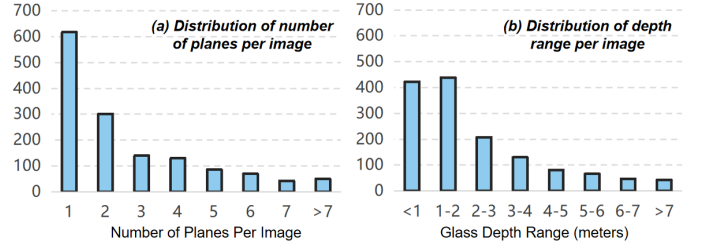
### D. Dataset Information



Fig. 3. **Statistics of proposed dataset.** (a) Distribution of number of planes in individual images. (b) Distribution of the range of glass depths in individual images.

The proposed labeling pipeline enables efficient generation of precise depth ground truth for a large number of images with minimal manual effort. Building on this capability, our dataset exhibits a substantially greater diversity in glass planar configurations. Eventually, we obtain 1437 annotated frames, with 1,070 images allocated to the training set and 367 images to the validation set. As illustrated in Fig. 3(a), the number of glass planes per image in our dataset ranges from 1 up to 10, providing a broader spectrum of geometric complexity. Furthermore, as shown in Fig. 3(b), the range of glass depths within individual images spans from $0.07\,\text{m}$ to $17.76\,\text{m}$, increasing the challenge of 3D estimation. Additionally, unlike previous approaches, 21 out of our 50 scenes are collected at night, introducing a wide variety of illumination conditions and further challenging the robustness of glass detection methods.

## IV. METHODOLOGY

### A. Framework Overview

Our proposed MonoGlass3D is designed to perform both glass segmentation and plane regression simultaneously with a simple overall architecture. As shown in Fig. 4. Specifically, the outputs from the last three layers of the DINOv2 encoder, each with dimensions $C \times \frac{H}{14} \times \frac{W}{14}$, are processed by the adaptive fusion modules. The adaptive fusion module first reduces the feature channel dimension to 256. Then the prediction results from the previous level, namely centerness $c^{i-1}$, plane parameters $p^{i-1}$, and segmentation $s^{i-1}$, are concatenated with the features and further refined using convolution layers. Centerness $c^i$ is predicted first; it is then fused with the segmentation and plane features by element-wise multiplication, enhancing the features according to the geometric properties captured by centerness. The adapted features, $F_s$ for segmentation and $F_p$ for plane regression, are fed into their respective prediction heads to generate segmentation masks and plane parameter predictions at $\frac{H}{14} \times \frac{W}{14}$ resolution for each encoder layer. Features from the last layer, $F_p^3$ and $F_s^3$ are concatenated, augmented with position encoding, then passed to self-attention layers for further refinement, then the final glass segmentation mask is upsampled via transpose
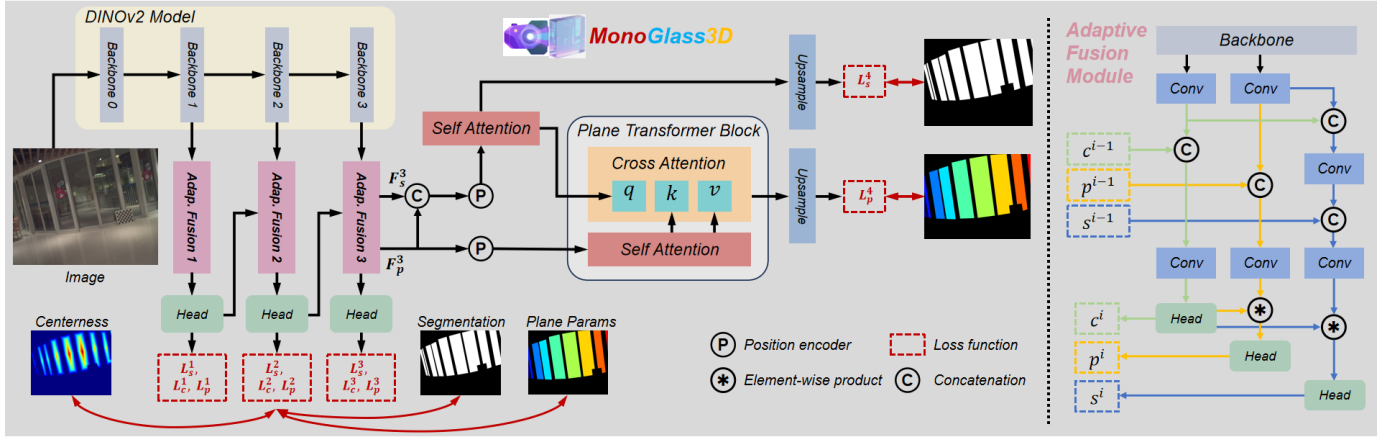
Fig. 4. **Proposed MonoGlass3D framework.** Our network consists DINOv2 backbone, centerness adaptive fusion module, and attention layers. Our network jointly performs glass segmentation and plane regression, which are supervised by segmentation loss $L_s$, centerness loss $L_c$, and plane loss $L_p$.

convolution layers. The output plane features from the last adaptive fusion layer, $F_p^3$, is combined with position encoding, and further refined using self-attention layers with transformer encoders, followed by cross-attention with the segmentation features using transformer decoder layers. Finally, the refined features are up sampled to predict the plane parameters via transpose convolutions.

### B. Plane Parameter Regression

As mentioned, a 3D plane can be defined by four variables: a 3D normal vector and interception. However, directly regressing four unconstrained variables is not ideal. The 3D space has 6 Degrees of Freedom (DoF), three for translation and three for rotation, but the intrinsic dimensionality of a plane is only 3 DoF. This is because two translation axes and one rotation axis (along the plane direction) are inherently undefined for a plane. Consequently, representing a plane using four free variables leads to over-parameterization. Usually, the plane normal vector is normalized to fix the vector length DoF. To avoid the need for post hoc normalization, we re-map the plane normal to a polar coordinate system with two angles, $[\theta_1, \theta_2]$, as shown in Fig. 5. The algebraic conversion is defined as:

$$
\theta_1 = \begin{cases} \arccos\left(\frac{r_{xz}}{r}\right), & \text{if } y > 0, \\ -\arccos\left(\frac{r_{xz}}{r}\right), & \text{otherwise} \end{cases}
$$
$$
\theta_2 = \begin{cases} \pi - \arccos\left(\frac{z}{r_{xz}}\right), & \text{if } x > 0, \\ -\arccos\left(\frac{z}{r_{xz}}\right), & \text{otherwise} \end{cases} \tag{5}
$$

where $r_{xz}$ is the vector magnitude in $x$ and $z$ axes, $r_{xz} = \sqrt{r_x^2 + r_z^2}$, and $r = 1$ for unit vectors. When the plane is perpendicular to $z$ axis, $\theta_1 = \theta_2 = 0$, which corresponds to normal vector $[0, 0, -1]$.

The plane normal vector can be regulated in the $-z$ hemisphere, this is because according to (2), when the normal vector $\mathbf{n}$ points to the $+z$ direction, we can flip the sign of both $\mathbf{n}$ and $d$ to yield the equivalent plane representation. The polar representation effectively reduces the parameter space, and narrows down the range of plane parameters to $\theta_1, \theta_2 \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ and $d \in \mathbb{R}$. Compared with conventional approaches that estimate four parameters (three for the normal
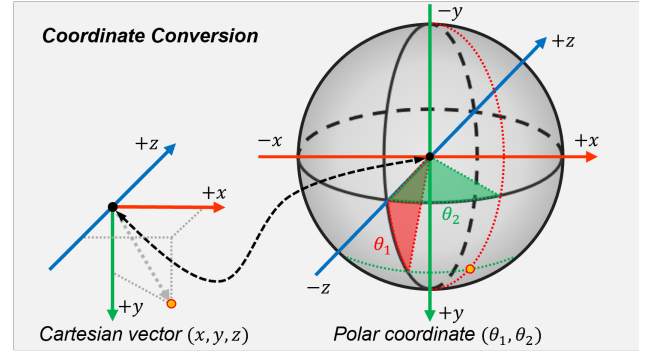


Fig. 5. **Coordinate conversion of plane normal vector.** Cartesian vector $(x, y, z)$ is mapped to its corresponding polar coordinate $(\theta_1, \theta_2)$.

vector and one for the intercept), our angular representation reduces plane estimation to three parameters, $[\theta_1 \quad \theta_2 \quad d]$, and eliminates the need for additional normalization. In our network, the prediction heads estimate the angles $\theta_1$ and $\theta_2$ using $\tanh$ activation, up-scaled by $\frac{\pi}{2}$, while the interception $d$ is up-scaled by a factor of 5.

For planes with a same normal direction, the value and sign of interception $d$ vary with the plane's position within the image, making positional encoding critical for accurate inference. Since the DINOv2 backbone utilizes transformer encoder layers with positional encoding, we directly apply adaptive fusion module to the outputs of the last three encoder layers to predict the plane parameters. Additionally, we enhance both the segmentation features $F_s^3$ and the plane features $F_p^3$ with extra positional encoding. We then refine $F_p^3$ by applying self-attention through transformer encoder layers, providing global context, and followed by cross-attention using transformer decoder layers, where segmentation features serves as the query embedding to localize glass across the image. The final plane parameter estimations are up sampled and projected to $3 \times H \times W$, with 3 channels corresponding to $[\theta_1 \quad \theta_2 \quad d]$.

### C. Adaptive Feature Fusion with Centerness

In real-world scenarios, glass typically exhibit regular shapes such as rectangles, trapezoids, and ellipses. Geometric
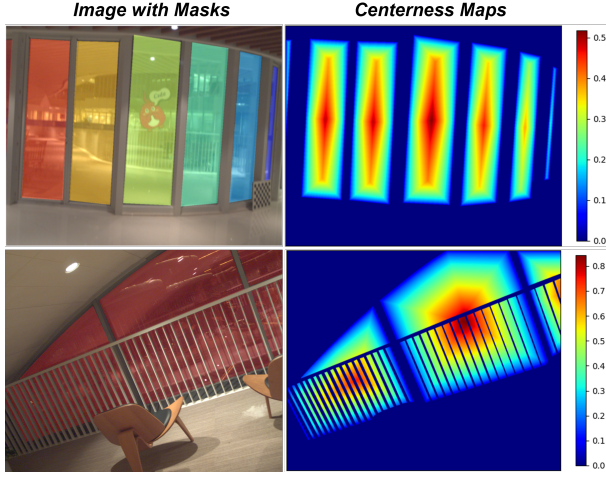
Fig. 6. **Centerness map examples.** Image with masks (left) and corresponding centerness maps (right).

context is therefore a crucial property for both plane estimation and glass segmentation. Previous works have leveraged boundary information [4], [5] to enhance contextual learning; however, boundaries are represented as thin lines in images, making them challenging to be learned efficiently.

In our approach, we adopt the concept of centerness introduced by [29] and [15]. Centerness was originally designed for instance segmentation, it is quantified by the normalized distance from a pixel to the boundary of an object instance. Specifically, for an pixel within an instance, centerness is defined by:

$$\mathcal{C} = \sqrt{\frac{d_{min}}{d_{max}}} \qquad (6)$$

where $d_{min}$ and $d_{max}$ are the shortest and longest distance to the instance contour from that pixel, respectively.

As shown in Fig. 6, the centerness distribution closely reflects the geometric shape of the mask. Typically, the centerness forms a smooth, pyramid-like distribution, instances with aspect ratios closer to 1 exhibit higher peak values at their centers. Compared to boundary context, centerness provides a much smoother signal that is easier for neural networks to learn and generalize from, while still capturing rich geometric context at the instance level.

In our network, predicted centerness is utilized to enhance both segmentation and plane features in adaptive fusion modules, but we multiply with $1 - \mathcal{C}$, because the features extracted from the glass surface is less important and possibly distracting, we would like the network to focus more on the context near the boundaries, which are more useful for extracting the geometric properties of glass. Centerness is integrated by feature fusion through:

$$F_{fused} = F * (1 - \mathcal{C}) \qquad (7)$$

where $*$ denotes element-wise multiplication.

By integrating centerness in this manner, our model adaptively emphasizes features according to their geometric context, thereby improving both segmentation accuracy and plane parameter estimation.

## D. Loss Function

The training of our network is supervised by three types of ground-truth: centerness, segmentation masks, and plane parameters:

$$L = L_c + L_s + l_p \qquad (8)$$

where $L$ is the total loss, $L_c$ denotes the centerness loss, it is computed using binary cross-entropy (BCE) loss. For segmentation, we employ a composite loss $L_s$, that combines BCE loss and Intersection over Union (IoU) loss:

$$L_s = 0.5L_{bce} + L_{iou} \qquad (9)$$

The plane loss (denoted as $L_p$) is designed to more effectively supervise the learning of plane parameters. Previous works such as [24], [27] use depth loss to supervise plane parameter estimation. However, this is not appropriate, because depth is the projection of 3D points onto the $z$-axis from the camera's perspective, it is distributed non-uniformly for different plane orientations.

To further elaborate this issue, as shown in Fig. 7(a), consider an image pixel at $x$, for a fixed angular error $\delta\theta$ between predicted planes $p_1, p_2$, and ground-truth planes $p_1^*, p_2^*$, the depth errors $d_1$ and $d_2$ can vary significantly. As depicted in Fig. 7(b), the depth error increases nonlinearly with the plane angle (when the plane interception fixed at 1), leading to highly uneven error distributions in the plane angle domain. Therefore, directly employing depth error as the supervision signal for plane parameters can result in suboptimal training.

To address this issue, we propose the plane distance loss. Given a plane estimation at a particular pixel location $(u_x, u_y)$, we first determine its corresponding 3D position $p$ by projecting the pixel onto the estimated plane. Next, we select 4 points on this plane, each at an equal distance from $p$ and oriented along orthogonal directions, as illustrated in Fig. 7(c). The plane distance loss is defined as the sum of the distances from these four points to the ground-truth plane (represented by the white dotted lines). Unlike depth error, these distances represent the true geometric discrepancy between the two planes and are invariant with respect to the camera viewpoint and plane orientation.

In addition to the plane distance loss, we add an L1 loss directly on the plane parameters. Thus, the plane loss for each pixel is the sum of the plane distance loss and the plane parameter loss:

$$L_p = L_{param} + L_{dist} \qquad (10)$$

Since different glass instances occupy varying areas in the image, we normalize the plane loss at the instance level to mitigate the effect of mask size. For an image with $N$ instances, where the $i$-th instance contains $M_i$ pixels, the plane loss is averaged across all pixels and all instances, which can be defined as:

$$L_p = \frac{1}{N} \sum_{i=0}^{N} \frac{1}{M_i} \sum_{j=0}^{M_i} L_{p,(i,j)} \qquad (11)$$

Segmentation and plane losses are computed for all output stages, while centerness loss is only evaluated at cascade

**(a) Depth error changes with plane angle** | **(b) Depth error on plane angle domain** | **(c) Plane distance loss (white dotted lines)**
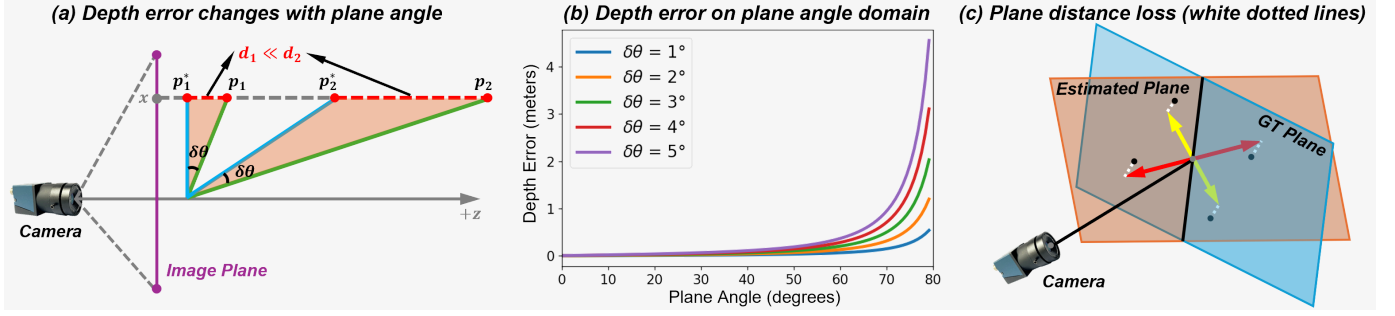
Fig. 7. **Drawbacks of depth loss and proposed plane distance loss.** (a) Depth error at same angular error varies largely depending on plane angles. (b) Plot of depth error on the plane angle domain. (c) Plane distance loss, measures the distance between 3D planes (white dotted lines).

layers. To balance the contributions from different stages, we apply stage-specific weights to each loss, defined as:

$$L_p = 0.1L_p^1 + 0.1L_p^2 + 0.2L_p^3 + 0.6L_p^4$$
$$L_s = 0.1L_s^1 + 0.1L_s^2 + 0.2L_s^3 + 0.6L_s^4 \quad (12)$$
$$L_c = 0.2L_c^1 + 0.3L_c^2 + 0.5L_c^3$$

## V. EXPERIMENTS

### A. Training Setup

Our network is implemented with PyTorch [30]. For the backbone, we employ DINOv2 ViT-S, initialized with pre-trained weights from Depth Anything v2 (ViT-S). Our network is trained and evaluated at resolution of $504 \times 630$. We use AdamW optimizer, the backbone learning rate is set to $5 \times 10^{-6}$, all other layers are set to $5 \times 10^{-5}$. We adopt a learning rate scheduler which decays the learning rate by a factor of 0.95 if the loss does not decrease for 8 epochs. Our model is trained for 230 epochs on 4 Nvidia RTX 4090 GPUs.

### B. Evaluation Metrics

We demonstrate the performance of our approach by comparing with recent works on two metrics: glass segmentation and depth estimation.

For segmentation performance, we use the following metrics: IoU, defined as $\frac{TP}{TP+FP+FN}$, it measures the overlap between prediction and ground-truth. Mean Absolute Error (MAE), defined as $\frac{1}{|\mathcal{I}|}\sum_{i \in \mathcal{I}}|y_i^* - y_i|$, where $y_i$ is the predicted probability at pixel $i$ in image $\mathcal{I}$, $y_i^*$ is the ground-truth. F1 score, defined as $\frac{2(\text{Precision} \times \text{Recall})}{\text{Precision}+\text{Recall}}$, it is a harmonic mean of Precision and Recall, where Precision is $\frac{TP}{TP+FP}$ and Recall is $\frac{TP}{TP+FN}$. Balanced Error Rate (BER), defined as $1 - 0.5(\frac{TP}{TP+FP} + \frac{TN}{TN+FN})$, it is the average of errors on both True and False predictions.

Although our network outputs plane parameters, we report results using depth estimation metrics for more intuitive comparison with previous works. We use the following depth evaluation metrics: Mean absolute error (MAE), similar to segmentation, except here $y$ denotes depth instead of segmentation probability. Root Mean Squared Error (RMSE), defined as $\sqrt{\frac{1}{|\mathcal{I}|}\sum_{i \in \mathcal{I}}(y_i^* - y_i)^2}$, which is a quadratic mean of the errors. Absolute Relative Error (ARE), defined as $\frac{1}{|\mathcal{I}|}\sum_{i \in \mathcal{I}}|y_i^* - y_i|/y_i^*$, it measures the error normalized by the true value. Accuracy with Threshold, defined as $S_t = \{i| \max(\frac{y}{y^*}, \frac{y^*}{y}) < t\}/|\mathcal{I}|\}$, It measures the proportion of



**Day Scenes** | **Night Scenes**

*(a) Coplanar Structures*

*(b) Multi-Angle Surfaces*

*(c) Multi-Occluded Layers*

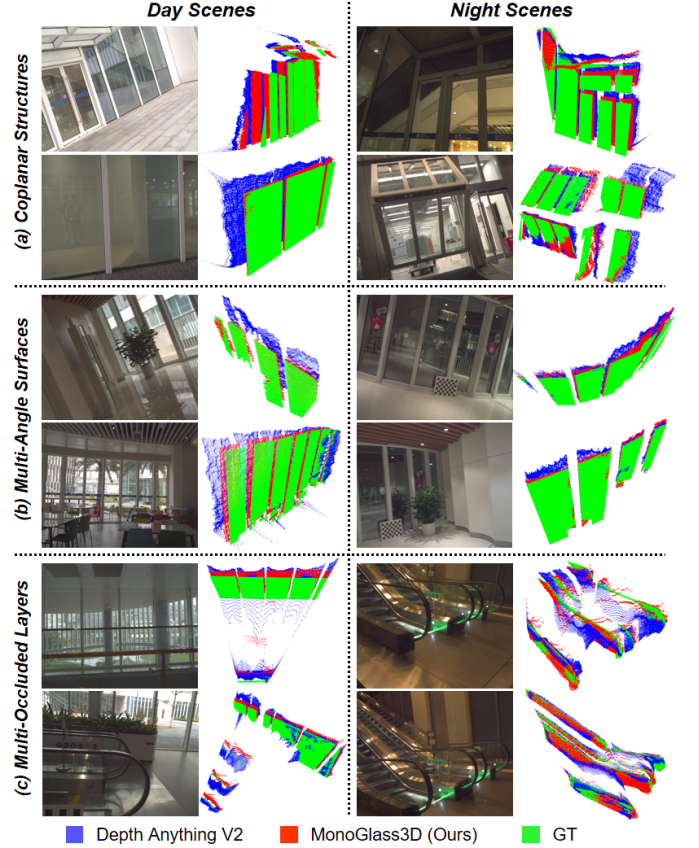Depth Anything V2 ■ | MonoGlass3D (Ours) ■ | GT ■

Fig. 8. **Comparisons of glass detection results.** (a) Coplanar structures contain glass surfaces lying on same plane. (b) Multi-angle surfaces contain scenes with glass instances with relatively small angular difference. (c) Multi-occluded layers contain scenes with overlapping glass surfaces.

pixels where the ratio of estimated to ground-truth depth is below threshold $t$. Following previous works [11], [12], we use thresholds of $1.25$, $1.25^2$, and $1.25^3$. These metrics are denoted as $\sigma_1 = S_{1.25}$, $\sigma_2 = S_{1.25^2}$, and $\sigma_3 = S_{1.25^3}$, respectively.

### C. Performance Evaluation

*1) Depth Estimation:* For the depth estimation task, we finetuned ZoeDepth [31] and Depth Anything V2 Metric [12] on our dataset. Quantitative results across various metrics are summarized in Table I, our method achieves the best performance under all evaluation metrics, and maintains stable performance in both day and night time data, despite having only $52.4$ M trainable parameters, which is only $54\%$ of Depth

TABLE I
DEPTH ESTIMATION METRICS COMPARISON

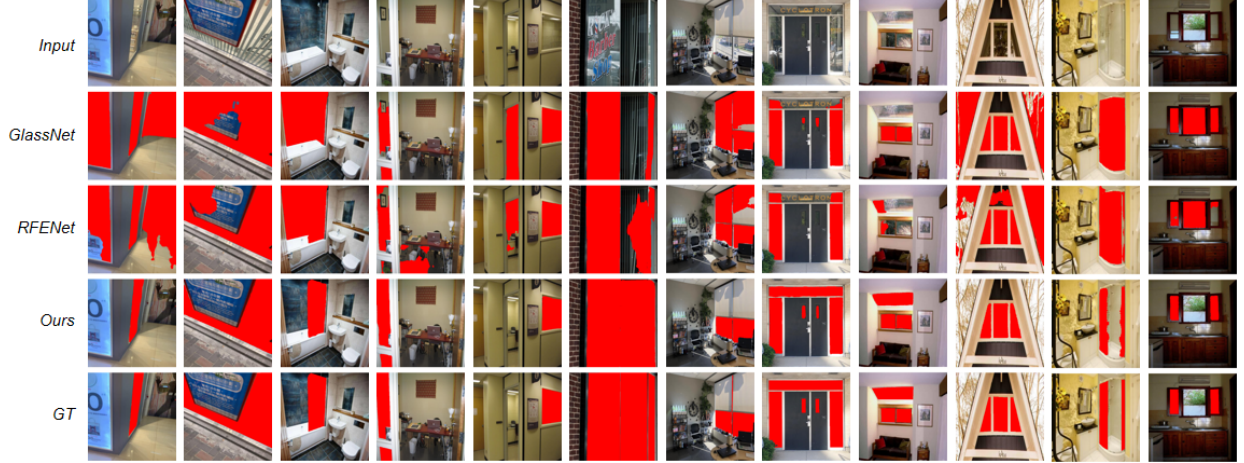| Test Data | Methods | Param. ↓ | Abs. Rel. ↓ | MAE ↓ | RMSE ↓ | $\sigma_1$ ↑ | $\sigma_2$ ↑ | $\sigma_3$ ↑ |
|---|---|---|---|---|---|---|---|---|
| | ZoeDepth (BEiT-B) [31] | 112.0M | 0.156 | 0.720 | 0.806 | 0.757 | 0.937 | 0.989 |
| All | Depth Anything V2 (ViT-S) [12] | **24.8M** | 0.090 | 0.387 | 0.419 | 0.911 | 0.993 | 0.997 |
| | Depth Anything V2 (ViT-B) [12] | 97.5M | 0.067 | 0.288 | 0.318 | 0.970 | 0.993 | 0.997 |
| | MonoGlass3D | 52.4M | **0.063** | **0.274** | **0.298** | **0.982** | **0.994** | **0.997** |



Fig. 9. **Comparisons of glass segmentation results.** Compared to previous approaches, our method is better at accurately detecting glass in images, even when the surrounding context is ambiguous.

Anything V2 (ViT-B). To facilitate a more intuitive comparison, we also visualize 3D projections of the estimated plane parameters, as shown in Fig. 8. For analysis, we categorize our dataset into three representative scene types.

*Coplanar Structures* (Fig. 8(a)): This subset consists of images containing glass surfaces with regular geometric structures, often featuring multiple instances sharing identical plane parameters. In such cases, our approach demonstrates robust performance owing to the explicit plane parameter estimation, which ensures consistent and accurate surface fitting. Compared to direct depth regression, the depth maps projected from plane parameters are notably flatter and better aligned with the ground truth surfaces.

*Multi-Angle Surfaces* (Fig. 8(b)): These examples contain multiple glass instances with minor angular differences between planes. Our method maintains stable and accurate fitting across surfaces with subtle variations, a task that is challenging for conventional depth estimation approaches [12], [31] that lack explicit geometric supervision (e.g., normal vectors or plane parameters).

*Multi-Occluded Layers* (Fig. 8(c)): This subset includes scenes with multiple overlapping glass surfaces, presenting the most challenging cases in our dataset. Our approach outperforms Depth Anything V2 in handling significant depth discontinuities. In the plane parameter space, parallel planes share the same normal but differ in intercepts, making such features more distinguishable and learnable by the network.

*2) Glass Segmentation:* For a fair comparison with other methods, we disable the layers required for plane regression in our model, then trained and evaluated on two commonly used glass segmentation datasets: GDD [2] and GSD [17]. Quantitative results are reported in Table II and Table III. Our model achieves state-of-the-art performance on both datasets, while

TABLE II
GLASS SEGMENTATION COMPARISON RESULTS ON THE GDD DATASET

| Methods | Param. ↓ | IoU ↑ | F1 ↑ | MAE ↓ | BER ↓ |
|---|---|---|---|---|---|
| GDNet [2] | 183.2M | 0.876 | 0.937 | 0.063 | 5.62 |
| GlassNet [17] | 83.7M | 0.881 | 0.932 | 0.059 | 5.71 |
| EBLNet [4] | 46.2M | 0.887 | 0.940 | 0.055 | 5.36 |
| GlassSemNet [19] | 240.2M | 0.908 | 0.950 | 0.045 | 4.34 |
| VBNet [14] | - | 0.907 | 0.948 | 0.048 | 4.70 |
| MonoGlass3D | **34.0M** | **0.920** | **0.951** | **0.038** | **3.76** |

TABLE III
GLASS SEGMENTATION COMPARISONS RESULTS ON THE GSD DATASET

| Methods | Param. ↓ | IoU ↑ | F1 ↑ | MAE ↓ | BER ↓ |
|---|---|---|---|---|---|
| GDNet [2] | 183.2M | 0.790 | 0.869 | 0.069 | 7.72 |
| GlassNet [17] | 83.7M | 0.836 | 0.903 | 0.055 | 6.12 |
| GlassSemNet [19] | 240.2M | 0.856 | 0.920 | 0.044 | 5.60 |
| RFENet [5] | 152.6M | 0.865 | **0.931** | 0.048 | 6.23 |
| VBNet [14] | - | 0.861 | 0.921 | 0.043 | 5.51 |
| MonoGlass3D | **34.0M** | **0.872** | 0.917 | **0.040** | **5.16** |

using significantly fewer parameters compared to competing methods. Qualitative examples are shown in Fig. 9. Compared to previous works that focus on specialized contextual cues, such as boundaries [4], [5] or optical features [11], [14], our approach leverages centerness maps for richer instance information and improved geometric context, enabling more accurate glass region identification, especially when boundaries are unclear or optical cues are weak.

### D. Ablation Studies

To evaluate the effectiveness of individual components in our proposed model, we conducted a series of ablation studies. The comparative results for both glass segmentation and depth estimation are presented in Table IV. We begin with a 'Base'

TABLE IV
ABLATION STUDIES ON DIFFERENT COMPONENTS

| Base | AF | SA | CA | $L_{depth}$ | $L_{dist}$ | Cascades | Param. ↓ | IoU ↑ | Abs. Rel. ↓ | MAE ↓ | RMSE ↓ | $\sigma_1$ ↑ | $\sigma_2$ ↑ | $\sigma_3$ ↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ✓ | | | | | ✓ | ✓ | **31.15M** | 0.925 | 0.082 | 0.340 | 0.374 | 0.950 | 0.987 | 0.995 |
| ✓ | | ✓ | | | ✓ | ✓ | 38.97M | 0.924 | 0.072 | 0.317 | 0.346 | 0.972 | 0.993 | 0.997 |
| ✓ | | ✓ | ✓ | | ✓ | ✓ | 48.56M | 0.923 | 0.071 | 0.306 | 0.340 | 0.977 | 0.993 | 0.997 |
| ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | 52.43M | 0.940 | 0.073 | 0.313 | 0.336 | 0.975 | 0.993 | 0.997 |
| ✓ | ✓ | ✓ | ✓ | | ✓ | | 45.87M | 0.939 | 0.064 | 0.278 | 0.302 | 0.978 | 0.994 | 0.997 |
| ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | 52.43M | **0.942** | **0.063** | **0.274** | **0.298** | **0.982** | **0.994** | **0.997** |

'Base' represents backbone and prediction heads, 'AF' represents adaptive fusion modules, 'SA' represents self-attention layers, 'CA' represents cross-attention layers.
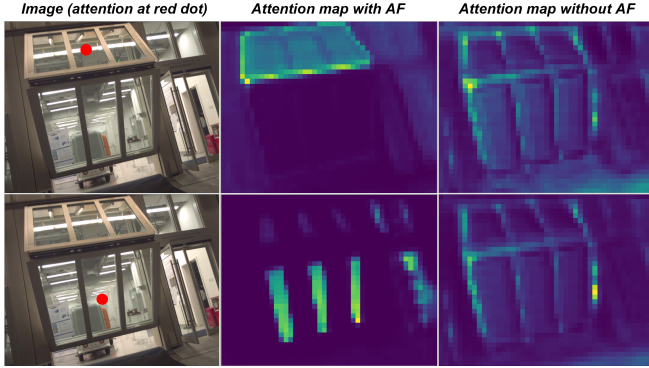


Fig. 10. **Attention maps with and without AF.** The attention maps at the red dots show that AF modules enable the network to recognize different planes.

model comprising only the backbone and prediction heads. Subsequently, we incrementally add the cross-attention (CA) module, self-attention (SA) module, and the adaptive fusion (AF) module to assess their respective contributions.

The results indicate that the transformer-based attention modules, CA and SA have limited impact on segmentation performance in IoU, but yield noticeable improvements in plane regression accuracy. The cascade structure is implemented to improve the learning efficiency of backbone layers, it also demonstrated noticeable improvements in overall performance. The AF module contributes the most significant performance gains among all components, which demonstrates the effectiveness of our centerness based adaptive feature fusion. To further illustrate its impact, attention maps generated with and without the AF module are compared in Fig. 10. The model enhanced with the AF module is able to more effectively capture the coplanar relations of the three glass windows at the top and the three at the bottom of the image.

Conventionally, segmentation and depth supervision do not provide explicit instance-level guidance, which is generally acceptable when object distinction is not required. In contrast, our approach seeks to distinguish glass surfaces residing on different planes, necessitating instance-level information. Centerness maps supply this instance-level knowledge without the computational overhead of explicit instance-level inference, which would unnecessarily complicate the problem.

We also evaluate the impact of different plane regression supervision strategies. As discussed in Section IV-D, we argue that directly supervising plane parameter estimation using depth error is problematic due to its non-uniformity with respect to plane orientation. To address this, we adopt a

TABLE V
DEPTH ESTIMATION METRICS COMPARISON ON REAL SCENARIOS

| Methods | Abs. Rel. ↓ | MAE ↓ | RMSE ↓ | $\sigma_1$ ↑ |
|---|---|---|---|---|
| ZoeDepth (BEiT-B) [31] | 0.215 | 1.056 | 1.159 | 0.600 |
| Dep. Any. V2 (ViT-S) [12] | 0.123 | 0.547 | 0.577 | 0.841 |
| Dep. Any. V2 (ViT-B) [12] | 0.088 | 0.386 | 0.409 | 0.952 |
| MonoGlass3D | **0.083** | **0.383** | **0.405** | **0.976** |

loss based on the geometric distance between predicted and ground-truth planes. After retrained our network with depth loss $L_{depth}$, in place of plane distance loss $L_{dist}$, as shown in Table IV, the results are consistent with our earlier analysis. The segmentation performance remains comparable, but depth estimation metrics deteriorate significantly. This confirms that the network is unable to achieve similar performance levels with depth supervision alone, further validating the effectiveness of our proposed plane distance loss.

### E. Real Scene Experiment

The ultimate goal of this research is to provide a practical monocular 3D glass detection solution. To further assess real-world performance, from our dataset, we deliberately select six scenes in which all frames are excluded from the training set, including indoor and outdoor scenes with different illumination conditions, resulting in a total of 180 images. This experimental setup allows us to rigorously evaluate the generalization capability of our model by explicitly testing it on entirely unseen environments. The results on depth metrics are presented in Table V, with segmentation performance of 0.940 IoU. Some result samples are shown in Fig. 11. Although our model's performance decreases compared to evaluations where most scenes are included in training, it still consistently outperforms competing methods on these challenging samples. Our approach demonstrated stable and reliable performance, particularly when detecting large glass surfaces at close range (< 10 meters). However, false negatives are more likely to occur for glass surfaces located farther from the camera. We attribute this limitation to the predominance of short-range data in our training set, which constrains the model's ability to effectively detect distant glass surfaces.

### F. Discussion

*1) Limitations:* Our method works on the basis of assuming all glass surfaces can be approximated as 3D planes, which

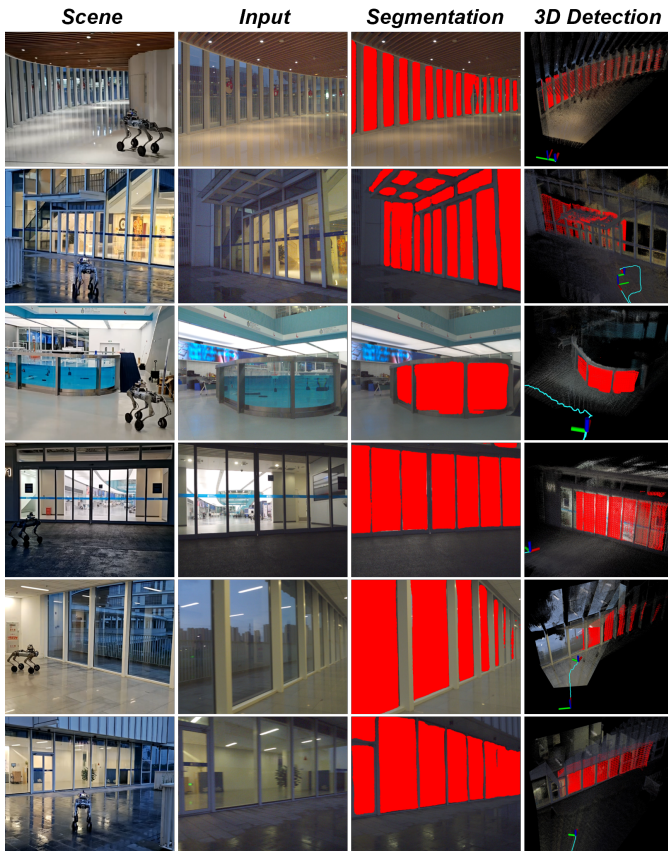| Scene | Input | Segmentation | 3D Detection |
|---|---|---|---|



Fig. 11. **3D glass detection results in real scenes.** We evaluate our approach in scenes that are absent in the training set, and it maintains strong performance.

## VI. CONCLUSION

In this work, we address the challenging problem of 3D glass detection with monocular input. To facilitate this research, we build a novel dataset featuring complex glass structures with precise annotations. We propose MonoGlass3D, a unified network performs glass segmentation and plane regression simultaneously. Central to our design is the Adaptive Fusion module, which leverages centerness maps to provide geometric context guidance, significantly enhancing both segmentation and plane parameter estimation. We formulate 3D glass detection as a plane regression problem, which allows us to exploit the inherent geometric properties of 3D planes for more effective learning. Extensive experiments demonstrate that our approach achieves state-of-the-art performance in both glass segmentation and glass depth estimation.

does not always hold in real-world scenarios. In practice, glass surface shapes follow a long-tailed distribution: while the majority can be reasonably approximated as planes, there remain numerous cases with curved or irregular geometries. This intrinsic variability presents a significant challenge for designing a 3D glass detection method that generalizes to arbitrary glass shapes.

Additionally, despite our efforts to streamline data collection and annotation, constructing our dataset remains a labor-intensive process. This raises challenges for scaling up the dataset to cover a wider variety of glass scenarios. The primary difficulty arises from the need to annotate glass surfaces both in 3D point clouds and in 2D RGB images, and then accurately associate these labels across modalities. Since the 3D positions of glass are not directly observable, we have to rely on supplementary data modalities and manual annotation to achieve precise 3D labeling.

*2) Future Work:* While end-to-end monocular 3D glass detection is conceptually simple and elegant, monocular RGB cameras inherently suffer from limited 3D perception capabilities, primarily due to scale ambiguity and susceptibility to varying lighting conditions. For increased robustness and accuracy in 3D perception, it is desirable to incorporate data from additional sensor modalities. In future work, we plan to investigate the integration of multi-modal inputs—such as depth sensors, LiDAR, or stereo vision—to further enhance the performance and reliability of 3D glass detection systems.

## REFERENCES

[1] J. Jiang, G. Cao, J. Deng, T.-T. Do, and S. Luo, "Robotic perception of transparent objects: A review," *IEEE Transactions on Artificial Intelligence*, vol. 5, no. 6, pp. 2547–2567, 2023.

[2] H. Mei, X. Yang, Y. Wang, Y. Liu, S. He, Q. Zhang, X. Wei, and R. W. Lau, "Don't hit me! glass detection in real-world scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3687–3696.

[3] B. Liu, G. Zhao, J. Jiao, G. Cai, C. Li, H. Yin, Y. Wang, M. Liu, and P. Hui, "Omnicolor: a global camera pose optimization approach of lidar-360camera fusion for colorizing point clouds," in *2024 IEEE International Conference on Robotics and Automation*, 2024, pp. 6396–6402.

[4] H. He, X. Li, G. Cheng, J. Shi, Y. Tong, G. Meng, V. Prinet, and L. Weng, "Enhanced boundary learning for glass-like object segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 859–15 868.

[5] K. Fan, C. Wang, Y. Wang, C. Wang, R. Yi, and L. Ma, "Rfenet: Towards reciprocal feature evolution for glass segmentation," *arXiv preprint arXiv:2307.06099*, 2023.

[6] X. Tao, N. Lang, H. Li, and D. Xu, "Path planning in uncertain environment with moving obstacles using warm start cross entropy," *IEEE/ASME Transactions on Mechatronics*, vol. 27, no. 2, pp. 800–810, 2021.

[7] J. Jiang, G. Cao, A. Butterworth, T.-T. Do, and S. Luo, "Where shall i touch? vision-guided tactile poking for transparent object grasping," *IEEE/ASME Transactions on Mechatronics*, vol. 28, no. 1, pp. 233–244, 2022.

[8] L. Jiang, S. Liu, Y. Cui, and H. Jiang, "Path planning for robotic manipulator in complex multi-obstacle environment based on improved_rrt," *IEEE/ASME Transactions on Mechatronics*, vol. 27, no. 6, pp. 4774–4785, 2022.

[9] J. Xu, N. Xi, C. Zhang, Q. Shi, and J. Gregory, "A robot-assisted back-imaging measurement system for transparent glass," *IEEE/ASME Transactions on Mechatronics*, vol. 17, no. 4, pp. 779–788, 2011.

[10] L. Zhu, A. Mousavian, Y. Xiang, H. Mazhar, J. van Eenbergen, S. Debnath, and D. Fox, "RGB-D local implicit function for depth completion of transparent objects," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4649–4658.

[11] Y. Liang, B. Deng, W. Liu, J. Qin, and S. He, "Monocular depth estimation for glass walls with context: a new dataset and method," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 12, pp. 15 081–15 097, 2023.

[12] L. Yang, B. Kang, Z. Huang, Z. Zhao, X. Xu, J. Feng, and H. Zhao, "Depth anything v2," *Advances in Neural Information Processing Systems*, vol. 37, pp. 21 875–21 911, 2024.

[13] M. Roberts, J. Ramapuram, A. Ranjan, A. Kumar, M. A. Bautista, N. Paczan, R. Webb, and J. M. Susskind, "Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 912–10 922.

[14] F. Qi, X. Tan, Z. Zhang, M. Chen, Y. Xie, and L. Ma, "Glass makes blurs: Learning the visual blurriness for glass surface detection," *IEEE Transactions on Industrial Informatics*, vol. 20, no. 4, pp. 6631–6641, 2024.

[15] E. Xie, P. Sun, X. Song, W. Wang, X. Liu, D. Liang, C. Shen, and P. Luo, "Polarmask: Single shot instance segmentation with polar representation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 193–12 202.

[16] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson *et al.*, "SAM 2: Segment anything in images and videos," *arXiv preprint arXiv:2408.00714*, 2024.

[17] J. Lin, Z. He, and R. W. Lau, "Rich context aggregation with reflection prior for glass surface detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13 415–13 424.

[18] F. Liu, Y. Liu, J. Lin, K. Xu, and R. W. Lau, "Multi-view dynamic reflection prior for video glass surface detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 4, 2024, pp. 3594–3602.

[19] J. Lin, Y.-H. Yeung, and R. Lau, "Exploiting semantic relations for glass surface detection," *Advances in Neural Information Processing Systems*, vol. 35, pp. 22 490–22 504, 2022.

[20] A. Kalra, V. Taamazyan, S. K. Rao, K. Venkataraman, R. Raskar, and A. Kadambi, "Deep polarization cues for transparent object segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8602–8611.

[21] H. Mei, B. Dong, W. Dong, J. Yang, S.-H. Baek, F. Heide, P. Peers, X. Wei, and X. Yang, "Glass segmentation using intensity and spectral polarization cues," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 622–12 631.

[22] D. Huo, J. Wang, Y. Qian, and Y.-H. Yang, "Glass segmentation with RGB-Thermal image pairs," *IEEE Transactions on Image Processing*, vol. 32, pp. 1911–1926, 2023.

[23] S. Sajjan, M. Moore, M. Pan, G. Nagaraja, J. Lee, A. Zeng, and S. Song, "Clear grasp: 3D shape estimation of transparent objects for manipulation," in *2020 IEEE International Conference on Robotics and Automation*, 2020, pp. 3634–3642.

[24] C. Liu, J. Yang, D. Ceylan, E. Yumer, and Y. Furukawa, "PlaneNet: Piece-wise planar reconstruction from a single RGB image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2579–2588.

[25] F. Yang and Z. Zhou, "Recovering 3D planes from a single image via convolutional neural networks," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 85–100.

[26] C. Liu, K. Kim, J. Gu, Y. Furukawa, and J. Kautz, "PlaneRCNN: 3D plane detection and reconstruction from a single image," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4450–4459.

[27] V. Patil, C. Sakaridis, A. Liniger, and L. Van Gool, "P3Depth: Monocular depth estimation with a piecewise planarity prior," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1610–1621.

[28] C. Zheng, W. Xu, Z. Zou, T. Hua, C. Yuan, D. He, B. Zhou, Z. Liu, J. Lin, F. Zhu *et al.*, "FAST-LIVO2: Fast, direct lidar-inertial-visual odometry," *IEEE Transactions on Robotics*, 2024.

[29] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: Fully convolutional one-stage object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9627–9636.

[30] 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," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[31] S. F. Bhat, R. Birkl, D. Wofk, P. Wonka, and M. Müller, "Zoedepth: Zero-shot transfer by combining relative and metric depth," *arXiv preprint arXiv:2302.12288*, 2023.

**Kai Zhang** (Student Member, IEEE) Recieved the B.Eng. degree in mechatronics engineering from the University of Sydney, Sydney, Australia, in 2020 and the M.Phil. degree in robotics and autonomous systems in the The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China, in 2024, where he is currently working toward the Ph.D. degree in robotics and autonomous systems at the Intelligent Autonomous Driving Center. His research interests include robotics perception, SLAM, and embodied intelligence.

**Guoyang Zhao** (Student Member, IEEE) received the B.Eng. degree in logistics engineering from Northeast Agricultural University, Harbin, China, in 2022, and the M.Phil. degree in robotics and autonomous systems from The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China, in 2024. He is currently pursuing the Ph.D. degree at the Intelligent Autonomous Driving Center, Robotics and Autonomous Systems Thrust, The Hong Kong University of Science and Technology, Guangzhou, China. His research interests include computer vision, robotics navigation, and deep learning.

**Jianxing Shi** (Student Member, IEEE) received his B.Eng. degree in Mechanical Engineering from Shandong University in 2022 and his M.Phil. degree in Robotics and Autonomous Systems from The Hong Kong University of Science and Technology (Guangzhou) in 2025. His research interests include sensor fusion, LiDAR SLAM, and robotic perception.

**Bonan Liu** (Student Member, IEEE) received the bachelor's degree in engineering mechanic from Central South University ChangSha, China, in 2020, the master's degree in Data Science from the City university of Hong Kong in 2022. He is currently a Ph.D. candidate student in the information hub in HKUST-GZ.

**Weiqing Qi** received the B.S. degree in Computer Science from University of California, Santa Barbara, CA, USA, in 2021, and the M.Phil. degree in robotics and autonomous systems from The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China, in 2024. His current research interests include lane detection, drivable area segmentation, and semantics segmentation, etc.

**Jun Ma** (Senior Member, IEEE) received the B.Eng. degree with First Class Honours in electrical and electronic engineering from Nanyang Technological University, Singapore, in 2014, and the Ph.D. degree in electrical and computer engineering from the National University of Singapore, Singapore, in 2018. From 2018 to 2021, he held several positions at the National University of Singapore; University College London, London, U.K.; University of California, Berkeley, Berkeley, CA, USA; and Harvard University, Cambridge, MA, USA. He is currently an Assistant Professor with the Robotics and Autonomous Systems Thrust, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China, and also with the Division of Emerging Interdisciplinary Areas, The Hong Kong University of Science and Technology, Hong Kong SAR, China. He is also the Director of Intelligent Autonomous Driving Center, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China. His research interests include motion planning and control for robotics and autonomous driving.