WHEN GLASS DISAPPEARS AT NIGHT: A NOVEL NIR-RGB MULTI-MODAL SOLUTION

Anonymous authors

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

033

037

040

041

042

043

044

045

047

048

051

052

Paper under double-blind review

ABSTRACT

Glass surface detection (GSD) has recently been attracting research interests. However, existing GSD methods focus on modeling glass surface properties for daytime scenes only, and can easily fail in nighttime scenes due to significant lighting discrepancies. We observe that, due to the spectral differences between Near-Infrared (NIR) light sources and common LED lights, NIR and RGB cameras capture complementary visual patterns (e.g., light reflections, shadows, and edges) of glass surfaces, and cross-comparing their lighting and reflectance properties can provide reliable cues for nighttime GSD. Inspired by this observation, we propose a novel approach for nighttime GSD based on the multi-modal NIR and RGB image pairs. We first construct a nighttime GSD dataset, which contains 6,192 RGB-NIR image pairs captured in diverse real-world nighttime scenes, with corresponding carefully-annotated glass surface masks. We then propose a novel network for the nighttime GSD task with two novel modules: (1) a RGB-NIR Guidance Enhancement (RNGE) module for extracting and enriching the NIR reflectance features with the guidance of RGB reflectance features, and (2) a RGB-NIR Fusion and Localization (RNFL) module for fusing RGB and NIR reflectance features into glass features conditioned on the multi-modal illumination discrepancy-aware features. Extensive experiments demonstrate that our method outperforms state-of-the-art methods in nighttime scenes while generalizing well to daytime scenes. We will release our dataset and codes.

1 Introduction

Glass surfaces, such as glass doors, walls, and windows, are ubiquitous in our daily lives. Their lack of intrinsic visual texture patterns can easily conceal glass surfaces within the background scene, causing significant detection difficulties. Failing to detect glass surfaces may cause the downstream vision applications, including 3D scene reconstruction and robotic navigation, to fail as well. Hence, glass surface detection (GSD) is a challenging, but fundamental task.

A few deep learning-based methods are proposed to detect glass surfaces based on RGB features (Mei et al., 2020; He et al., 2021; Lin et al., 2021; Fan et al., 2023; Liu et al., 2024; Lin et al., 2022; Yan et al., 2025; Qi et al., 2024), or RGB-X multi-modal features (Lin et al., 2025; Huo et al., 2023; Yan et al., 2024). These methods typically focus on modeling different priors for detecting glass surfaces, including contrasted RGB (Mei et al., 2020) or RGB-Thermal (Huo et al., 2023) features, boundaries (He et al., 2021; Fan et al., 2023), reflections (Lin et al., 2021; Liu et al., 2024; Yan et al., 2024), perceived noisy depth (Lin et al., 2025), semantic correlations (Lin et al., 2022), and ghosting effects (Yan et al., 2025). However, these GSD priors are specifically developed for daytime scenes and can be drowned in the low-light or complex artificial lighting of nighttime scenes. For example, as shown in the top two rows of Fig. 1, compared to NIR imaging, depth or thermal imaging provides limited contextual information in nighttime scenes for existing RGB-depth/thermal-based GSD methods (Lin et al., 2025; Huo et al., 2023) to detect glass regions. Meanwhile, modeling intrinsic cues such as boundary (Fan et al., 2023), reflections (Lin et al., 2021), and ghosting effects (Yan et al., 2025) for glass surface localization in the RGB domain (Lin et al., 2021; Fan et al., 2023; Yan et al., 2025), or comparing reflections between the RGB and NIR modalities (Yan et al., 2024) is unreliable in nighttime scenes, as demonstrated in the bottom two rows of Fig. 1.

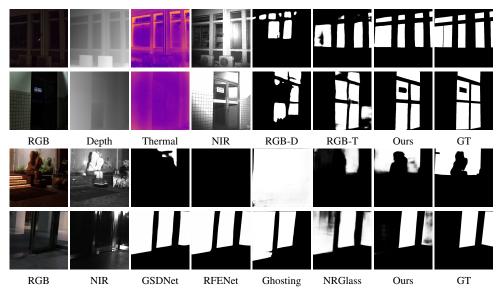


Figure 1: Upper two rows: Depth (Lin et al., 2025) and thermal (Huo et al., 2023) cues provide very limited contextual information for localizing glass surfaces in nighttime scenes, compared to NIR imaging. Bottom two rows: Intrinsic cues such as reflections (in either RGB (Lin et al., 2021) or RGB-NIR (Yan et al., 2024) domain), and boundary (Fan et al., 2023) and ghosting effects (Yan et al., 2025) (in the RGB domain), can be easily buried in nighttime scenes. We propose to cross-compare the lighting and reflectance information between RGB and NIR modalities for accurate glass surface detection in nighttime scenes.

We observe that NIR and RGB cameras can capture complementary visual cues (*e.g.*, reflection/transmission discrepancies and illumination discrepancies) for nighttime GSD. By projecting their own lights, active NIR cameras ensure consistent visibility in low-light/uneven lighting conditions, producing geometry/reflectance patterns on glass surfaces that complement those (*e.g.*, colors, textures, and semantics) in the RGB modality. Based on this observation, we propose in this paper a novel RGB-NIR-based approach, which considers the complementary patterns of glass surfaces between NIR and RGB images, for nighttime GSD.

As there are no available datasets for this task, we first construct a large-scale RGB-NIR glass surface detection dataset, with a hybrid imaging system consisting of a DSLR camera and a NIR camera accompanied by an active NIR light source. Our dataset contains 6, 192 RGB-NIR image pairs captured from diverse real-world nighttime scenes, with the corresponding manually annotated glass surface masks. We then propose a novel neural network to model the complementary patterns on glass surfaces between NIR and RGB images for glass surface detection.

To extract the complementary patterns between RGB and NIR images as cues, our method first performs a learning-based image decomposition to decompose the input RGB and NIR images into two pairs of reflectance and illumination components, and then uses two separate encoders to extract the semantics-aware and material-aware contextual features from the reflectance components of the two modalities. We further introduce two novel modules. First, we propose a novel RGB-NIR Guidance Enhancement (RNGE) module, which leverages the semantic features extracted from the reflectance component of the input RGB image for feature extraction-and-enhancement of the reflectance component of the NIR image. In other words, the RNGE module aims to explore the semantics from the RGB image to assist glass feature extraction, especially boundary prediction, from the NIR image. Meanwhile, we model the illumination discrepancies between RGB and NIR images as gating matrices based on their derived Illumination components, and propose a novel RGB-NIR Fusion and Localization (RNFL) module for decoding the multi-modal reflectance features into glass features conditioned on the derived gating matrices. As shown in Fig. 1, our method can produce more accurate detection results under challenging lighting conditions of night-time scenes. Our method can be deployed on popular surveillance cameras that switch between RGB and active NIR modes. The main contributions of this work can be summarized as follows:

¹Nighttime surveillance systems always use active NIR cameras accompanied by active NIR light sources.

- We propose the first approach for nighttime glass surfaces detection, by modeling the complementary patterns of glass surface regions between RGB and NIR image pairs.
- Our network includes two novel modules: a RNGE module for enriching NIR reflectance features with RGB reflectance features, and a RNFL module for GSD based on multi-modal reflectance feature aggregation guided by the multi-modal illumination differences.
- We construct the first large-scale nighttime GSD dataset, which contains 6K RGB-NIR glass image pairs (with corresponding masks) captured from diverse nighttime scenes.
- Extensive evaluations show that our proposed method outperforms SOTA methods in nighttime scenes and generalizes well to daytime scenes.

2 Related Work

Glass Surface Detection (GSD) has recently gained significant research attention with several deep learning-based methods proposed. A line of GSD methods are based on input RGB frames, which model the contrasted contextual features (Mei et al., 2020) and incorporate reflection priors (Lin et al., 2021; Liu et al., 2024), boundary detection (He et al., 2021; Lin et al., 2021; Fan et al., 2023), semantic correlations (Lin et al., 2022), blurry effects (Qi et al., 2024), and ghosting cues (Yan et al., 2025). Another line of methods explore RGB-X multi-modal imaging for GSD. Kalra et al. (2020) leverages the polarization information to segment transparent objects (*e.g.*, wine glass and glass balls), which may not generalize well to glass surfaces with irregular shapes. Huo et al. (2023) models the contrasted glass features between RGB and Thermal modalities. Yan et al. (2024) model the reflection differences between RGB and NIR images for daytime GSD. They use an NIR filter attached to the DSLR camera lens and rely on the ambient light to capture NIR images. Lin et al. (2025) propose to model the noise differences between RGB and depth images. Most recently, Zhang et al. (2025) propose the MonoGlass3D method, performing 3D glass segmentation and 3D plane regression simultaneously.

All these existing methods, however, are designed for daytime scenes, and their proposed cues may not be effective under low-light or complex artificial lighting conditions of nighttime scenes. In this work, we propose to detect glass surfaces at nighttime scenes by utilizing dual RGB-active NIR cameras. In contrast to passive NIR filters Yan et al. (2024), where both RGB and NIR spectra may not be illuminated in nighttime scenes, our imaging setup creates an induced photometric discrepancy (from active NIR illumination and ambient lighting), and our network learns to cross-compare this discrepancy for the detection.

Mirror Detection (MD) aims to detect mirror regions. Existing methods focus on learning the correlations between the reflected and real surrounding contents, by modeling the contextual contrasted RGB (Yang et al., 2019; Xu et al., 2024) or RGB-Depth (Mei et al., 2021b) features, pearance correspondences (Lin et al., 2020; 2023; Huang et al., 2023), visual chirality (Tan et al., 2023), spatial/frequency-based specular textures (Xie et al., 2024), and inconsistent motions (Warren et al., 2024). Although these methods have achieved impressive progress for mirror detection, glass surfaces have a very different property from mirrors. While mirrors only contain reflection, glass surfaces contain both reflection as well as transmission, making glass surfaces more challenging to detect. In this paper, we explore the use of RGB-NIR images for glass surface detection.

3 PROPOSED DATASET

We construct the first large-scale night-time glass surface detection dataset for training and evaluation.

Hybrid Imaging System. We design a hybrid imaging system consisting of a DSLR camera (Canon 70D), and an NIR camera (HIKVISION MVCH250-90GN) with an active NIR light source (40W). Both cameras are synchronized to capture RGB and NIR image pairs from the target scene simultaneously. Although the reflectivity of most glass surfaces decreases with increasing wavelength of the incident ray within the range of $780 \sim 1100$ nm (Huo et al., 2023; Planinsic, 2011), excessively long wavelengths will diminish night vision effectiveness (Ariff et al., 2015). Hence, we set the wavelength of the active near-infrared light to $850 \sim 940$ nm (Ariff et al., 2015) to balance the suppression of near-infrared reflection and night vision capability, allowing for clearer NIR images

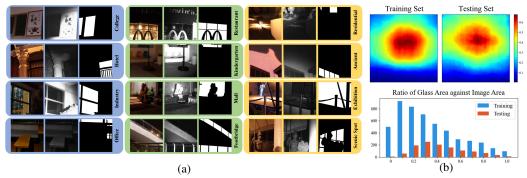


Figure 2: Examples and some statistics of our dataset: (a) Example RGB/NIR/Glass mask triplets from our dataset; (b) Glass location distributions (top) and glass/image ratios (bottom).

with less glass surface reflection. We use the binocular image alignment method (Shen et al., 2020) to align each RGB-NIR image pair.

Dataset Statistics. Our dataset consists of 6, 192 triplets of RGB and NIR images and the corresponding manually labeled glass surface masks. Our dataset covers 12 types of daily life scenes, such as campus, hotel, and shopping mall, as shown in Fig. 2a. We randomly split our dataset into 5,000 and 1,192 triplets for training and evaluation, respectively. Fig. 2b shows the statistics of glass location distribution and the ratio of glass area over the image. The glass regions cluster around the upper part of the image, as they are more likely placed around the eye level. Our dataset contains varying glass/image ratios, posing a significant challenge to detect.

4 PROPOSED METHOD

Our core idea is to cross-compare and fuse the lighting and material information captured in the RGB and NIR modalities for night-time glass surface detection. Fig. 3 shows an overview of our method, which includes the Retinex Decomposition, the Encoder, and the Decoder stages.

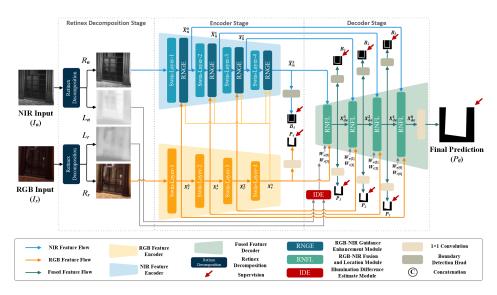


Figure 3: Method Overview. The NIR (I_n) and RGB (I_r) images captured at nighttime are first decomposed into corresponding reflectance $(R_n \text{ and } R_r)$ and illumination $(L_n \text{ and } L_r)$ layers. We use two encoders to process the two reflectance layers with a novel RNGE module to guide material-aware features extraction from the NIR reflectance R_n . We cross-compare the illumination layers $(L_n \text{ and } L_r)$ via the IDE module to control the feature flow in the decoder, where we propose the RNFL module to selectively fuse the dual-modality features for glass surface detection.

The Retinex Decomposition stage aims to extract lighting and reflectance information from the input multi-modal images, since the imaging of both RGB and NIR images is dependent on an external light source. We fine-tune the SOTA Content-Transfer Decomposition Network (Jiang et al., 2024) to decompose the input RGB image I_r into Reflectance component R_r and Illumination component L_r , and the input NIR image I_n into Reflectance component R_n and Illumination component L_n .

The Encoder stage aims to extract the semantics-aware features from the RGB reflectance component, which guides and enhances the material-aware feature extraction from the NIR reflectance component. We use the Swin Transformer V2 (Liu et al., 2022) to extract multi-scale features (denoted as X_r^i and X_n^i , where $i \in \{0,1,2,3\}$) from the two Reflectance components (R_r and R_n), respectively. At each scale, we propose the RNGE module (Sec. 4.1) to leverage the high-level semantic features X_r^3 and the low-level features X_r^i to guide and enhance the extraction of features X_n^i from R_n . Meanwhile, since the NIR reflectance component R_n exhibits clearer and more complete boundaries of glass surfaces than those in R_r , we perform the boundary detection on the deepest features \bar{X}_n^3 . Notably, for handling extreme low-light conditions, our method is designed to rely on the active NIR modality.

The Decoder stage aims to transform the multi-modal features into glass surface features for detection. We first model the illumination differences between the NIR and RGB images via the Illumination Difference Estimation (IDE) module (Sec. 4.2), which takes the illumination components of two modalities (L_n and L_r) as input and predicts two weight matrices W_r and W_n for controlling the feature flows in the decoder. We then propose the RGB-NIR Fusion and Localization (RNFL) module (Sec. 4.3), which works at multi-scales, integrating the extracted multi-modal features R_r and R_n from the Encoder stage with guidance from the IDE module for predicting the glass surface masks.

4.1 RGB-NIR GUIDANCE ENHANCEMENT (RNGE) MODULE

The proposed RNGE module aims to extract and enhance multi-scale NIR features from R_n conditioned on the RGB features. As shown in Fig. 4, we first concatenate X_n^i and X_r^i , and then use a convolution layer to produce the fused features X_f^i . Meanwhile, we compute the difference between X_n^i and X_r^i as X_d^i through element-wise subtraction. Since X_d^i is expected to capture cross-modal differences, which serve as cues for potential glass surface locations, we apply a convolution layer and a sigmoid (\cdot) activation to normalize X_d^i , yielding the activation

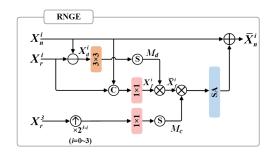


Figure 4: The proposed RNGE module.

map M_d . We then multiply X_f^i by M_d to obtain the enriched glass features \bar{X}_f^i . In addition, we apply supervision to X_r^3 to capture abundant semantic features of glass surfaces and use them as guidance to the multi-modal feature fusion. Specifically, X_r^3 is first upsampled to the size of X_n^i and normalized by a sigmoid (\cdot) function to produce the activation map M_c , which is then applied to \bar{X}_f^i for further enrichment. Finally, we use a self-attention (SA) block (Vaswani, 2017) to enhance the glass features and transform them back to the initial X_n^i through a residual connection to produce the output features \bar{X}_n^i . The whole process can be formulated as:

$$\begin{split} X_f^i &= \operatorname{Conv}(\operatorname{Concat}(X_n^i, X_r^i)), \\ \bar{X}_f^i &= X_f^i \otimes \operatorname{sigmoid}(\operatorname{Conv}(X_n^i - X_r^i)), \\ \bar{X}_n^i &= \operatorname{SA}(\bar{X}_f^i \otimes \operatorname{sigmoid}(\operatorname{Conv}(\operatorname{Up}(X_r^3)))) + X_n^i. \end{split} \tag{1}$$

Fig. 5 shows the input features (X_n^0) , the activation map M_c for enrichment, the contrasted activation map M_d between the two modalities, the attention map $Q \cdot K^T$ produced by the SA block of the RNGE module, and the enhanced features (\bar{X}_n^0) by the RNGE module.

4.2 ILLUMINATION DIFFERENCE ESTIMATION (IDE) MODULE

The IDE module, as shown in Fig. 6, aims to guide the fusion of the features of the RGB and NIR modalities, by modeling the difference in illumination components between the two modalities through estimating two gating matrices W_r and W_n .

Figure 5: Intermediate feature visualization of the RNGE module.

Specifically, to accurately measure the difference between L_n and L_r , both of them are first fed into the 1×1 Average Pooling (AP) block followed by two 5×5 convolution layers to produce refined features L'_n and L'_r . The difference L_d between them is then computed using element-wise subtraction to capture illumination variations across different areas of the two modalities. We further leverage L_d to enhance the different areas in the original RGB-NIR image pair by concatenating L_d with L'_n and L'_r

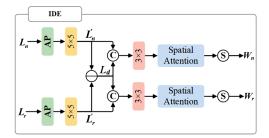


Figure 6: The proposed IDE module.

separately, and subsequently adjusted to 4 channels via two 3×3 convolution layers, matching the four scales of the Decoder. For each branch, we use a spatial-attention block (Woo et al., 2018) to further enhance the corresponding features and a $\operatorname{sigmoid}(\cdot)$ function to normalize these feature values to the range of [0,1], producing two weight matrices W_n and $W_r \in R^{W \times H \times 4}$, where [W,H] is the spatial resolution of the input images. The whole process can be described as:

$$L_{d} = \operatorname{Conv}(\operatorname{AP}(L_{n})) - \operatorname{Conv}(\operatorname{AP}(L_{r})),$$

$$W_{n} = \operatorname{sigmoid}(\operatorname{SPA}(\operatorname{Conv}(\operatorname{Concat}(L'_{n}, L_{d})))),$$

$$W_{r} = \operatorname{sigmoid}(\operatorname{SPA}(\operatorname{Conv}(\operatorname{Concat}(L'_{n}, L_{d})))),$$
(2)

where $SPA(\cdot)$ denotes the spatial attention block (Woo et al., 2018).

As the weight matrices W_r and W_n visualized in Fig. 7, the IDE module learns to prioritize deeper RGB features and shallower NIR features. This is because RGB contains richer semantic information, while NIR provides more distinct edge and structural details.

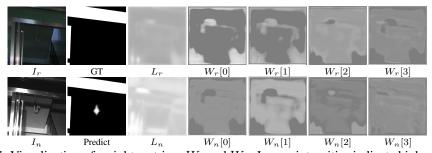


Figure 7: Visualization of weight matrices W_r and W_n . Larger intensities indicate higher weights.

4.3 RGB-NIR FUSION AND LOCALIZATION (RNFL) MODULE

The proposed RNFL module, as shown in Fig. 8, aims to effectively aggregate multi-modal features to localize glass surfaces, which contains two groups of cross-attention (CA) mechanisms, i.e., the left and right parts. The left two cross-attention mechanisms aim to extract the shared glass features from X_r^i and \bar{X}_n^i . Specifically, \bar{X}_n^i and X_r^i are used as Q in turn. When \bar{X}_n^i serves as Q, X_r^i acts as K and K to query features in K_r^i that are similar to K_r^i . A similar process applies when K_r^i is used as K0. Subsequently, at the top branch of the right cross-attention group, since the decoded features K_{de}^{i+1} produced by the deeper decoder layer contains semantic information about glass surfaces,

we use X_{de}^{i+1} as Q, while the detailed features \tilde{X}_r^i and \tilde{X}_n^i obtained from the first group serve as K and V, respectively, to perform crossattention for locating glass features. The situation at the bottom branch is similar. Finally, we apply the weights $W_r[i]$ and $W_n[i]$ obtained from the IDE module to the results of the right cross-attention group, termed as \hat{X}_r^i and \hat{X}_n^i , respectively. The weighted features from the top and bottom branches are then summed to produce the decoded feature X_{de}^i . This process can be formulated as:

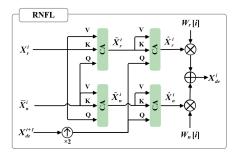


Figure 8: The proposed RNFL module.

$$\tilde{X}_{r}^{i} = \operatorname{CA}(\bar{X}_{n}^{i}, X_{r}^{i}), \quad \tilde{X}_{n}^{i} = \operatorname{CA}(X_{r}^{i}, \bar{X}_{n}^{i}),
\hat{X}_{r}^{i} = \operatorname{CA}(X_{de}^{i+1}, \tilde{X}_{r}^{i}), \quad \hat{X}_{n}^{i} = \operatorname{CA}(X_{de}^{i+1}, \tilde{X}_{n}^{i}),
X_{de}^{i} = W_{r}[i] * \hat{X}_{r}^{i} + W_{n}[i] * \hat{X}_{n}^{i}.$$
(3)

4.4 Training Strategy

We first train the Retinex Decomposition Stage of our network. We directly use the pre-trained CTDN module (Jiang et al., 2024) for RGB decomposition, and finetune the CTDN module for NIR decomposition using the NIR images of our dataset. We then train the Encoder and Decoder parts, using the BCE and IoU loss for supervising the glass surface detection and the Dice loss for glass boundary detection. Refer to the Supplemental for more training and implementation details.

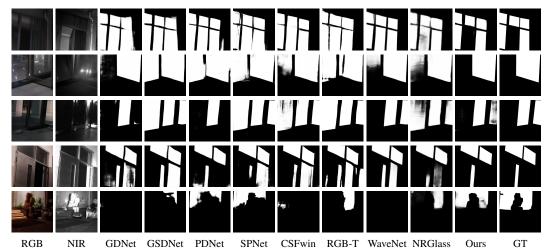


Figure 9: Visual comparison of our method with 8 competing methods.

5 RESULTS AND DISCUSSION

5.1 NIGHTTIME GSD RESULTS

Quantitative Results. We compare our method to 17 SOTA methods, including Glass Surface Detection (GSD) methods (Mei et al., 2020; Lin et al., 2021; Fan et al., 2023; Lin et al., 2022; Yan et al., 2025; Lin et al., 2025; Huo et al., 2023; Yan et al., 2024), Mirror Detection (MD) methods (Lin et al., 2020; Xie et al., 2024; He et al., 2023; Mei et al., 2021a), and Salient Object Detection (SOD) methods (Zhou et al., 2023a; 2021; Tu et al., 2021; Cong et al., 2022; Zhou et al., 2023b). All methods are re-trained on the proposed dataset. Tab. 1 reports the results. While multi-modal methods generally perform better than RGB-based methods, our method achieves the best performance on all five evaluation metrics with reasonable computational overheads.

Tab. 2 reports the results of four best-performing methods (according to Tab. 1, refer to the Supplemental for full results) taking as input the enhanced RGB images from the SOTA low-light enhancement

Table 1: Quantitative comparison between our method and 17 state-of-the-art methods on our proposed dataset. GSD, MD, and SOD indicate glass surface detection, mirror detection, and salient object detection, respectively. Best detection results are marked in **bold**.

Modal(s)	Method	Backbone	IoU \uparrow F $_{\beta}\uparrow$	ACC↑	MAE↓	BER↓	#Params	FLOPs(G)	Γime(ms)
RGB	GDNet (Mei et al., 2020)	ResNeXt	78.78 0.881	0.874	0.084	0.101	201.72M	207.95	229.7
RGB	PMD (Lin et al., 2020)	ResNeXt	79.50 0.871	0.893	0.079	0.094	147.66M	119.26	203.7
RGB	GSDNet (Lin et al., 2021)	ResNeXt	78.87 0.876	0.891	0.085	0.100	83.72M	41.27	195.2
RGB	GlassSemNet (Lin et al., 2022)	ResNet	80.12 0.886	0.907	0.088	0.093	361.33M	1412.03	215.2
RGB	RFENet (Fan et al., 2023)	ResNeXt	78.85 0.881	0.882	0.083	0.099	152.65M	756.91	26.3
RGB	HetNet (He et al., 2023)	ResNeXt	79.33 0.870	0.899	0.082	0.093	49.59M	38.73	161.1
RGB	CSFwinformer (Xie et al., 2024)	Swin	82.08 0.898	0.905	0.080	0.085	230.86M	188.62	187.2
RGB	GhostingNet (Yan et al., 2025)	SwinV2	81.37 0.889	0.910	0.077	0.091	271.53M	321.70	90.76
RGB-D	PDNet (Mei et al., 2021a)	ResNet	81.47 0.891	0.910	0.074	0.084	80.54M	69.85	197.9
RGB-D	SPNet (Zhou et al., 2021)	Res2Net	83.24 0.909	0.918	0.066	0.069	175.29M	81.05	216.4
RGB-D	CIRNet (Cong et al., 2022)	ResNet	82.38 0.903	0.904	0.069	0.082	103.15M	50.70	40.2
RGB-D	RGB-Depth (Lin et al., 2025)	ResNext	81.41 0.895	0.909	0.075	0.084	53.28M	32.13	33.8
RGB-T	MIDD (Tu et al., 2021)	ResNet	77.27 0.865	0.867	0.102	0.117	79.75M	169.70	92.1
RGB-T	RGBT GSD (Huo et al., 2023)	ResNet	81.50 0.891	0.908	0.069	0.083	85.02M	85.55	48.5
RGB-T	PRLNet (Zhou et al., 2023a)	Swin	83.01 0.893	0.906	0.065	0.079	570.66M	277.01	144.3
RGB-T	WaveNet (Zhou et al., 2023b)	WaveMLP	85.66 0.922	0.925	0.057	0.066	84.88M	64.02	154.3
RGB-NIR	NRGlassNet (Yan et al., 2024)	SwinV2	84.54 0.917	0.914	0.066	0.073	245.30M	265.35	163.2
RGB-NIR	Ours	SwinV2	87.98 0.934	0.936	0.047	0.055	234.88M	469.98	109.2
RGB-NIR	Ours (ablation)	ResNet	82.63 0.901	0.905	0.070	0.081	158.42M	428.81	65.1
RGB	Ours*	SwinV2	83.51 0.906	0.835	0.065	0.075	159.60M	284.27	23.6

method (Jiang et al., 2024). Since the enhanced images differ from normal daytime images in terms of noise distribution and reflection and transmission information in the glass region, most methods cannot achieve significant improvements. Please refer to the Supplemental for more quantitative comparisons conducted on GDD Mei et al. (2020) and GSD Lin et al. (2021) datasets.

Table 2: Comparison with competing methods taking the RGB images enhanced by LightenDiffusion (Jiang et al., 2024) as input. Results better than the corresponding ones in Tab. 1 are underlined.

Modal(s)	Methods	IoU↑	$F_{eta} \uparrow$	ACC↑	MAE↓	BER↓
RGB	GhostingNet (Yan et al., 2025)	81.29	0.887	0.911	0.079	0.092
RGB	CSFwinformer (Xie et al., 2024)	82.75	0.899	0.908	0.066	0.079
RGB-T	WaveNet (Zhou et al., 2023b)	85.73	0.921	0.927	0.055	0.067
RGB-NIR	NRGlassNet (Yan et al., 2024)	84.47	0.913	0.916	0.067	0.075
RGB-NIR	Ours	87.98	0.934	0.936	0.047	0.055

Visual Results. We compare our method with 8 state-of-the-art methods (Mei et al., 2020; Lin et al., 2021; Huo et al., 2023; Yan et al., 2024; Xie et al., 2024; Mei et al., 2021a; Zhou et al., 2021; 2023b) in Fig. 9. These examples demonstrate that our method can accurately detect glass surfaces under various challenging night-time conditions (*e.g.*, with opening doors) by exploiting complementary patterns on the glass surfaces between NIR and RGB images, while the competing methods often fail.

Table 3: We test the generalization ability of our method on existing multi-modal (*i.e.*, RGB-NIR (Yan et al., 2024), RGB-Thermal (Huo et al., 2023), and RGB-Depth (Lin et al., 2025)) daytime GSD datasets. All methods are re-trained on the corresponding datasets. Refer to the Supplemental for full comparisons. Best results are in **bold**.

Modal(s)	Methods	Venue	IoU↑	$F_{\beta} \uparrow$	MAE↓	BER↓
RGB-NIR RGB-NIR	NRGlassNet (Yan et al., 2024) Ours	KBS'24	90.03 88.97	0.955 0.947	0.033 0.050	0.036 0.044
RGB-T RGB-T	RGB-Thermal (Huo et al., 2023) Ours	TIP'23	93.80 94.19	0.965 0.969	0.027 0.029	0.040 0.041
RGB-D RGB-D	RGB-D GSD (Lin et al., 2025) Ours	AAAI'25	74.20 77.96	0.853 0.857	0.043 0.034	0.093 0.080

Figure 10: Comparison of our method with competing methods on daytime scenes.

5.2 DAYTIME GSD RESULTS

We now evaluate the generalization ability of our method on existing daytime GSD datasets (*i.e.*, RGB-NIR (Yan et al., 2024), RGB-Thermal (Huo et al., 2023), and RGB-Depth (Lin et al., 2025) GSD datasets). The results are shown in Tab. 3, where all methods are re-trained on the corresponding datasets. The comparison shows that our method generalizes reasonably well to daytime scenes, despite the existing domain discrepancies.² Fig. 10 shows some visual comparisons, where we can see that our method can accurately detect the glass surfaces in daytime scenes. Refer to the Supplemental for full comparisons.

5.3 ABLATION RESULTS

We report ablation results on the proposed dataset in Tab. 4. The first two rows show that using either single modality (*i.e.*, "w/o NIR" and "w/o RGB") significantly decreases the performance. The 3rd to 6th rows show the results when we individually ablate the Retinex decomposition, RNGE, IDE, and RNFL modules. The 7th to 9th rows show the ablation results obtained by removing specific input features and the subtraction operation from our RNGE model. Moreover, we replace the backbone with ResNet-50, termed as "Ours(ablation)", and ablate our network to a single-branch network (no NIR input, no NIR feature encoder), termed as "Ours*", as shown in Tab. 1. The results demonstrate that the "SwinV2" is more effective and our complete model obtains the best performance.

Refer to the Appendix for more implementation details, results, and analysis of ablations.

Table 4: Ablation study on multi-modal input and our proposed modules.

Methods	IoU↑	$F_{\beta} \uparrow$	ACC↑	MAE↓	BER↓
w/o NIR	84.19	0.910	0.924	0.062	0.072
w/o RGB	82.68	0.903	0.908	0.072	0.083
w/o Retinex Dec.	87.03	0.927	0.929	0.054	0.061
w/o RNGE	87.21	0.931	0.935	0.051	0.059
w/o IDE	87.35	0.935	0.933	0.048	0.057
w/o RNFL	87.09	0.915	0.931	0.049	0.060
$w/o X_r^3$	87.34	0.932	0.936	0.049	0.058
w/o X_r^i	87.42	0.933	0.931	0.049	0.059
w/o subtraction	87.67	0.933	0.936	0.048	0.056
Ours	87.98	0.934	0.936	0.047	0.055

6 Conclusion

We have proposed a novel method for night-time glass surfaces detection, by modeling the complementary patterns of glass surface regions between RGB and NIR image pairs. Our network has a novel RNGE module for RGB-to-NIR guiding feature enhancement and a novel RNFL module for glass surface detection based on the guided multimodal feature aggregation. We have also constructed the first large-scale night-time glass surface detection dataset. Extensive evaluations show that our proposed method outperforms the SOTA methods, and can generalize to daytime scenarios well.

²Note that the NIR images of the RGB-NIR dataset (Yan et al., 2024) are captured by covering an NIR filter on the DSLR camera len, which are different from ours.

REFERENCES

- Mohd Farid Mohd Ariff, Mohammad Ehsan Kosnan, Zulkepli Majid, Albert Chong, and Khairulnizam Idris. A study of near-infrared (nir) filter for surveillance application. *Jurnal Teknologi*, 2015.
- Runmin Cong, Qinwei Lin, Chen Zhang, Chongyi Li, Xiaochun Cao, Qingming Huang, and Yao Zhao. Cir-net: Cross-modality interaction and refinement for RGB-D salient object detection. *IEEE Trans. Image Processing*, 2022.
- Ke Fan, Changan Wang, Yabiao Wang, Chengjie Wang, Ran Yi, and Lizhuang Ma. Rfenet: Towards reciprocal feature evolution for glass segmentation. In *IJCAI*, 2023.
 - Hao He, Xiangtai Li, Guangliang Cheng, Jianping Shi, Yunhai Tong, Gaofeng Meng, Véronique Prinet, and Lubin Weng. Enhanced boundary learning for glass-like object segmentation. In *ICCV*, 2021.
 - Ruozhen He, Jiaying Lin, and Rynson W.H. Lau. Efficient mirror detection via multi-level heterogeneous learning. In *AAAI*, 2023.
 - Tianyu Huang, Bowen Dong, Jiaying Lin, Xiaohui Liu, Rynson WH Lau, and Wangmeng Zuo. Symmetry-aware transformer-based mirror detection. In *AAAI*, 2023.
 - Dong Huo, Jian Wang, Yiming Qian, and Yee-Hong Yang. Glass segmentation with rgb-thermal image pairs. *IEEE Trans. on Image Processing*, 2023.
 - Hai Jiang, Ao Luo, Xiaohong Liu, Songchen Han, and Shuaicheng Liu. Lightendiffusion: Unsupervised low-light image enhancement with latent-retinex diffusion models. In *ECCV*, 2024.
 - Agastya Kalra, Vage Taamazyan, Supreeth Krishna Rao, Kartik Venkataraman, Ramesh Raskar, and Achuta Kadambi. Deep polarization cues for transparent object segmentation. In *CVPR*, 2020.
 - Jiaying Lin, Guodong Wang, and Rynson W.H. Lau. Progressive mirror detection. In CVPR, 2020.
 - Jiaying Lin, Zebang He, and Rynson W.H. Lau. Rich context aggregation with reflection prior for glass surface detection. In *CVPR*, 2021.
 - Jiaying Lin, Yuen-Hei Yeung, and Rynson W.H. Lau. Exploiting semantic relations for glass surface detection. In *NeurIPS*, 2022.
 - Jiaying Lin, Xin Tan, and Rynson W.H. Lau. Learning to detect mirrors from videos via dual correspondences. In CVPR, 2023.
 - Jiaying Lin, Yuen-Hei Yeung, Shuquan Ye, and Rynson W.H. Lau. Leveraging rgb-d data with cross-modal context mining for glass surface detection. *AAAI*, 2025.
 - Fang Liu, Yuhao Liu, Jiaying Lin, Ke Xu, and Rynson WH Lau. Multi-view dynamic reflection prior for video glass surface detection. In *AAAI*, 2024.
 - Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, Furu Wei, and Baining Guo. Swin transformer V2: scaling up capacity and resolution. In *CVPR*, 2022.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
 - Haiyang Mei, Xin Yang, Yang Wang, Yuanyuan Liu, Shengfeng He, Qiang Zhang, Xiaopeng Wei, and Rynson W.H. Lau. Don't hit me! glass detection in real-world scenes. In *CVPR*, 2020.
- Haiyang Mei, Bo Dong, Wen Dong, Pieter Peers, Xin Yang, Qiang Zhang, and Xiaopeng Wei. Depth-aware mirror segmentation. In *CVPR*, 2021a.
- Haiyang Mei, Bo Dong, Wen Dong, Pieter Peers, Xin Yang, Qiang Zhang, and Xiaopeng Wei. Depth-aware mirror segmentation. In *CVPR*, 2021b.
 - Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *3DV*, 2016.

- Gorazd Planinsic. Infrared thermal imaging: Fundamentals, research and applications. *European Journal of Physics*, 2011.
- Fulin Qi, Xin Tan, Zhizhong Zhang, Mingang Chen, Yuan Xie, and Lizhuang Ma. Glass makes blurs:
 Learning the visual blurriness for glass surface detection. *IEEE Trans. on Industrial Informatics*,
 2024.
- Xi Shen, François Darmon, Alexei A. Efros, and Mathieu Aubry. Ransac-flow: Generic two-stage image alignment. In *ECCV*, 2020.
- Xin Tan, Jiaying Lin, Ke Xu, Pan Chen, Lizhuang Ma, and Rynson W.H. Lau. Mirror detection with the visual chirality cue. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2023.
 - Zhengzheng Tu, Zhun Li, Chenglong Li, Yang Lang, and Jin Tang. Multi-interactive dual-decoder for rgb-thermal salient object detection. *IEEE Trans. on Image Processing*, 2021.
 - A Vaswani. Attention is all you need. In NeurIPS, 2017.
- Alex Warren, Ke Xu, Jiaying Lin, Gary K.L. Tam, and Rynson W.H. Lau. Effective video mirror detection with inconsistent motion cues. In *CVPR*, 2024.
- Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, 2018.
 - Enze Xie, Wenjia Wang, Wenhai Wang, Mingyu Ding, Chunhua Shen, and Ping Luo. Segmenting transparent objects in the wild. In *ECCV*, 2020.
 - Zhifeng Xie, Sen Wang, Qiucheng Yu, Xin Tan, and Yuan Xie. Csfwinformer: Cross-space-frequency window transformer for mirror detection. *IEEE Trans. Image Processing*, 2024.
 - Ke Xu, Tsun Wai Siu, and Rynson WH Lau. Zoom: learning video mirror detection with extremely-weak supervision. In *AAAI*, 2024.
 - Tao Yan, Shufan Xu, Hao Huang, Helong Li, Lu Tan, Xiaojun Chang, and Rynson W.H. Lau. NRGlassNet: Glass surface detection from visible and near-infrared image pairs. *Knowledge-Based Systems*, 2024.
 - Tao Yan, Jiahui Gao, Ke Xu, Xiangjie Zhu, Hao Huang, Helong Li, Benjamin Wah, and Rynson W.H. Lau. Ghostingnet: A novel approach for glass surface detection with ghosting cues. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2025.
 - Xin Yang, Haiyang Mei, Ke Xu, Xiaopeng Wei, Baocai Yin, and Rynson W.H. Lau. Where is my mirror? In *ICCV*, 2019.
 - Kai Zhang, Guoyang Zhao, Jianxing Shi, Bonan Liu, Weiqing Qi, and Jun Ma. Monoglass3d: Monocular 3d glass detection with plane regression and adaptive feature fusion. *arXiv* preprint *arXiv*:2509.05599, 2025.
 - Heng Zhou, Chunna Tian, Zhenxi Zhang, Chengyang Li, Yuxuan Ding, Yongqiang Xie, and Zhongbo Li. Position-aware relation learning for rgb-thermal salient object detection. *IEEE Trans. on Image Processing*, 2023a.
 - Tao Zhou, Huazhu Fu, Geng Chen, Yi Zhou, Deng-Ping Fan, and Ling Shao. Specificity-preserving RGB-D saliency detection. In *ICCV*, 2021.
- Wujie Zhou, Fan Sun, Qiuping Jiang, Runmin Cong, and Jenq-Neng Hwang. Wavenet: Wavelet network with knowledge distillation for rgb-t salient object detection. *IEEE Trans. on Image Processing*, 2023b.

A APPENDIX AND SUPPLEMENTARY MATERIALS

In this appendix, we provide implementation details of our method, more analysis of the impact of pre-processing operations on input data. We also provide more qualitative visual comparisons between existing state-of-the-art methods from relevant fields and our model. Finally, we analyzed the impact of different modal inputs on nighttime glass detection problems and also examined the practical role of Retinex decomposition.

A.1 IMPLEMENTATION DETAILS AND LOSS FUNCTIONS

Our network is implemented using PyTorch on a Nvidia RTX 4090 GPU. Swin Transformer V2 (Liu et al., 2022) pre-trained on ImageNet-1K is adopted as the backbone of our network. The resolution of the image input is set to 384×384 . For data augmentation, we follow the previous work (Lin et al., 2021) to use random cropping, random rotation, and random horizontal flip. We use the AdamW (Loshchilov & Hutter, 2019) optimizer, while the initial learning rate is set to 1×10^{-5} , and the batch size is set to 2. We train our model with 100 epochs, which takes about 10 hours.

As for loss function, we first train the Retinex Decomposition Stage of our network. We directly use the pre-trained CTDN module (Jiang et al., 2024) for RGB decomposition, and finetune the CTDN module for NIR decomposition using the NIR images of our dataset via the reconstruction loss L_{rec} as follow:

$$L_{rec} = ||I_n - R_n \times L_n||_2. \tag{4}$$

We then train the Encoder and Decoder parts, using the BCE and IoU loss for supervising the glass surface detection and the Dice loss for glass boundary detection. The BCE loss and IoU loss are adopted to supervise the predicted glass surface mask P_i at i-th scale, and the Dice loss (Milletari et al., 2016) for boundary predictions B_i . A 3×3 convolution is used as the boundary detection head. The prediction loss function L_{pred} and boundary loss function L_{bound} can be defined as follow:

$$L_{pred} = \sum_{i=0}^{4} (L_{bce}(P_i, P_{gt}) + L_{iou}(P_i, P_{gt})), \tag{5}$$

$$L_{bound} = \sum_{i=1}^{4} L_{dice}(B_i, B_{gt}), \tag{6}$$

where P_{gt} and B_{gt} are the ground truth glass masks and boundary maps, $L_{bce}(\cdot)$, $L_{iou}(\cdot)$, and $L_{dice}(\cdot)$ are the BCE, IoU, and Dice Losses, respectively. i=4 is used to index the glass surface and boundary supervisions at the top of the Encoder. The total loss for the second stage can be defined as follow:

$$L_{stage2} = L_{pred} + \lambda L_{bound}, \tag{7}$$

where λ is a hyper-parameter, empirically set to 0.1.

A.2 SUPPLEMENTAL QUANTITATIVE COMPARISON

Moreover, Tab. 5 shows the evaluation results of the competing methods taking Reflectance components of RGB images (and that of NIR images) decomposed by CTDN (Jiang et al., 2024) as input. In another word, the reflectance components of RGB and NIR images takes place of the original RGB and NIR images. The results demonstrate that most competing methods can not directly benefit from the learning-based Retinex Decomposition.

Tab. 6 reports the results of the competing methods taking as input the enhanced RGB images from the SOTA low-light enhancement method (Jiang et al., 2024). Existing methods may not be easily benefited, as the enhancement may introduce noisy patterns that affect the glass surface detection.

We conducted additional experiments on the glass surface detection benchmarks GDD (Mei et al., 2020) and GSD (Lin et al., 2021) to validate the robustness of our model. As shown in the Tab. 7, we achieve state-of-the-art performance across most metrics, demonstrating that our model is well-suited for addressing the problem of glass surface detection.

Table 5: Study on competing methods directly taking Reflectance component R_r (and R_n) as input. Each value better than its corresponding value shown in Table 1 of our paper is marked in <u>underlined</u>.

Methods	IoU↑	$F_{\beta} \uparrow$	Acc↑	MAE↓	BER↓
GDNet Mei et al. (2020)	78.22	0.872	0.878	0.085	0.103
GSDNet Lin et al. (2021)	77.98	0.876	0.873	0.085	0.105
GlassSemNet Lin et al. (2022)	79.05	0.879	0.887	0.088	0.097
RFENet Fan et al. (2023)	78.53	0.872	0.887	0.086	0.102
GhostingNet Yan et al. (2025)	81.23	0.884	0.908	0.080	0.093
PMD Lin et al. (2020)	78.89	0.866	0.889	0.083	0.097
HetNet He et al. (2023)	78.59	0.863	0.895	0.086	0.096
CSFwinformer Xie et al. (2024)	82.14	0.896	0.907	0.070	0.083
PDNet Mei et al. (2021a)	80.77	0.889	0.903	0.077	0.086
SPNet Zhou et al. (2021)	82.53	0.904	0.907	0.072	0.078
CIRNet Cong et al. (2022)	81.34	0.897	0.898	0.076	0.086
RGB-Depth Lin et al. (2025)	80.94	0.887	0.903	0.078	0.086
MIDD Tu et al. (2021)	76.96	0.863	0.866	0.105	0.118
RGB-Thermal Huo et al. (2023)	79.09	0.877	0.891	0.084	0.091
PRLNet Zhou et al. (2023a)	82.37	0.887	0.901	0.067	0.083
WaveNet Zhou et al. (2023b)	84.59	0.912	0.902	0.066	0.077
NRGlassNet Yan et al. (2024)	83.84	0.911	0.910	0.070	0.078
Ours	87.98	0.934	0.936	0.047	0.055

Table 6: Comparisons with competing methods taking the enhanced RGB images produced by LightenDiffusion (Jiang et al., 2024) as input. Results better than the corresponding ones in Tab. 1 of our paper are <u>underlined</u>.

Methods	IoU↑	$F_{\beta} \uparrow$	ACC↑	MAE↓	BER↓
GDNet Mei et al. (2020)	78.82	0.878	0.880	0.082	0.100
GSDNet Lin et al. (2021)	78.92	0.878	0.880	0.081	0.100
GlassSemNet Lin et al. (2022)	80.20	0.885	0.900	0.080	0.094
RFENet Fan et al. (2023)	79.22	0.882	0.878	0.080	0.099
GhostingNet Yan et al. (2025)	81.29	0.887	0.911	0.079	0.092
PMD Lin et al. (2020)	79.54	0.870	0.893	0.080	0.094
HetNet He et al. (2023)	78.82	0.869	0.887	0.083	0.097
CSFwinformer Xie et al. (2024)	82.75	<u>0.899</u>	<u>0.908</u>	<u>0.066</u>	<u>0.079</u>
PDNet Zhou et al. (2021)	81.16	0.892	0.910	$0.076 \\ 0.075 \\ \underline{0.067} \\ \underline{0.073}$	0.085
SPNet Zhou et al. (2021)	82.88	0.896	0.914		0.071
CIRNet Cong et al. (2022)	82.43	0.903	0.900		0.082
RGB-Depth Lin et al. (2025)	81.46	0.891	<u>0.911</u>		0.086
MIDD Tu et al. (2021)	77.53	0.867	0.865	0.098	0.102
RGB-Thermal Huo et al. (2023)	81.70	0.896	<u>0.904</u>	0.080	0.087
PRLNet Zhou et al. (2023a)	83.07	0.891	<u>0.909</u>	0.068	0.080
WaveNet Zhou et al. (2023b)	85.73	0.921	<u>0.927</u>	0.055	0.067
NRGlassNet Yan et al. (2024) Ours	84.47	0.913	0.916	0.067	0.075
	87.98	0.934	0.936	0.047	0.055

A.3 SUPPLEMENTAL QUALITATIVE COMPARISON

Fig. 11, 12 and 13 show more qualitative comparison results produced by our method and 8 SOTA competing methods Mei et al. (2020); Lin et al. (2021); Huo et al. (2023); Yan et al. (2024); Xie et al. (2024); Mei et al. (2021a); Zhou et al. (2021; 2023b).

The Fig. 11 shows that smooth and flat non-glass surfaces can easily lead to misjudgment in glass detection methods due to the reflection of visible light, mistakenly identifying it as a glass area. Our method utilizes the characteristic of similar reflection phenomena on non-glass surfaces and

Table 7: Quantitative comparison between our method and 9 state-of-the-art methods on the benchmark datasets GDD and GSD. Best detection results are marked in **bold**.

		GDD				G	SD		
Methods	Venue	IoU↑	$F_{\beta} \uparrow$	$MAE \downarrow$	$BER{\downarrow}$	IoU↑	$F_{eta} \uparrow$	$MAE{\downarrow}$	BER↓
MINet	CVPR'20	84.35	0.919	0.077	0.074	77.29	0.879	0.077	0.095
SINet	CVPR'20	83.27	0.912	0.101	0.084	77.04	0.875	0.077	0.092
TransLab	ECCV'20	82.93	0.891	0.091	0.089	74.05	0.837	0.088	0.114
GDNet	CVPR'20	87.63	0.937	0.063	0.056	79.01	0.869	0.069	0.077
GSDNet	CVPR'21	88.07	0.932	0.059	0.057	83.64	0.903	0.055	0.061
EBLNet	ICCV'21	88.16	0.939	0.059	0.056	85.04	0.916	0.053	0.064
RFENet	IJVAI'23	88.72	0.940	0.055	0.054	86.50	0.931	0.048	0.062
VBNet	TII'24	90.58	0.944	0.048	0.047	85.90	0.915	0.043	0.054
GhostingNet	TPAMI'25	89.30	0.943	0.054	0.051	83.77	0.904	0.055	0.061
Ours	-	91.27	0.946	0.044	0.045	87.91	0.921	0.039	0.047

significant differences in reflection phenomena on glass surfaces in two modals to distinguish glass regions.

The Fig. 12 shows that open doors and windows have the same edge features as closed glass doors and windows, leading to the failure of edge dependent detection methods. Our method focuses on the differences between the two modes in the door and window area, so it can accurately distinguish between open doors and windows and closed glass doors and windows.

The Fig. 13 shows that the curved glass boundary features pose challenges to glass surface detection algorithms. Ordinary glass typically exhibits regular geometric shapes, and curved glass boundaries can mislead this prior knowledge.

The Fig. 14 shows more results on the day-time NIR dataset (Yan et al., 2024).

These qualitative results demonstrate that our method has better performance than the competing methods.

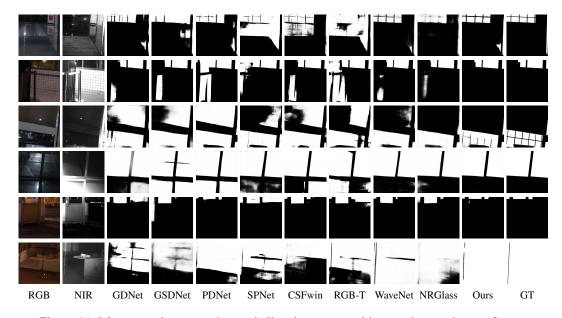


Figure 11: More experiment results on challenging scenes with smooth non-glass surfaces.

A.4 STUDY ON DIFFERENT MULTI-MODAL INPUTS

We have also conducted an experiment to study the different combinations of multimodal image data (such as RGB, Depth, NIR and thermal images) for glass surface detection. The hybrid imaging

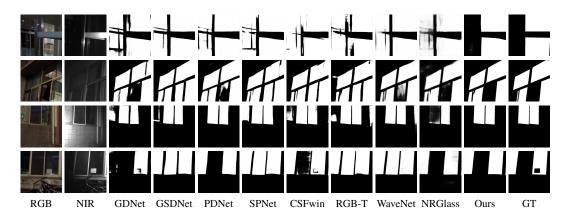


Figure 12: More experiment results on challenging scenes with opened windows and doors.

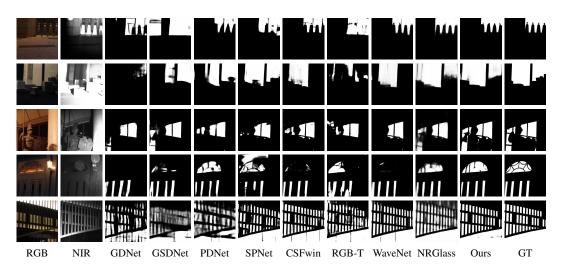


Figure 13: More experiment results on challenging scenes with complex boundaries.

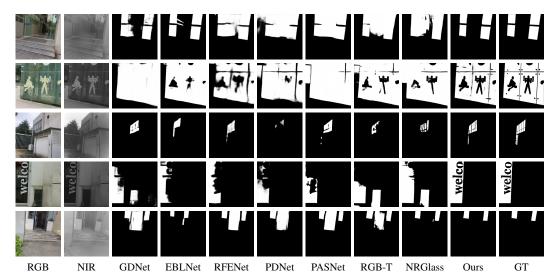


Figure 14: Comparison of our method with competing methods on daytime scenes.

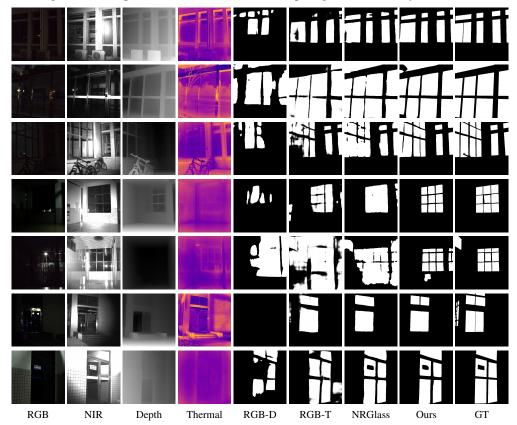


Figure 15: Study on different multi-modal inputs.

systems we used to capture multimodal image data are shown in Fig. 19. As shown in Fig. 15, 16 and 18 We captured RGB images, NIR images, thermal images, and depth maps of the target scenes. Then, we compare our method with the competing multi-modal methods including RGB-D based method (Lin et al., 2025), RGB-T based method(Huo et al., 2023) on the capture multimodal image data.

The 1st scene of Fig. 18 shows that, in low-light scenes, due to the lack of the sun as a thermal radiation source, the intensity difference between different regions in the thermal images is very

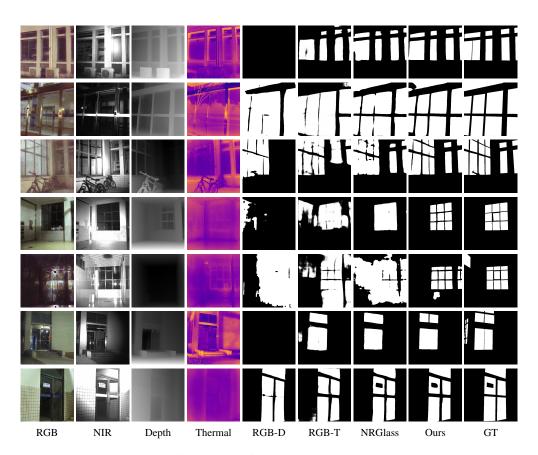


Figure 16: Use the Reflectance component as input.

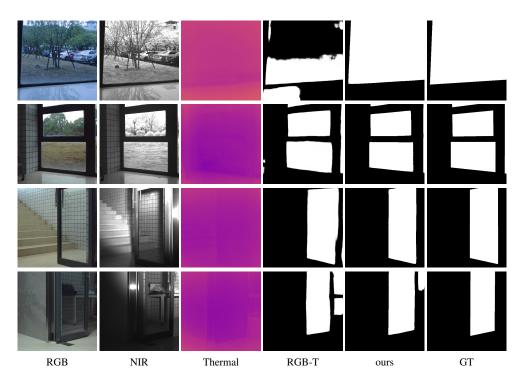


Figure 17: Some indoor low heat radiation difference scenes.

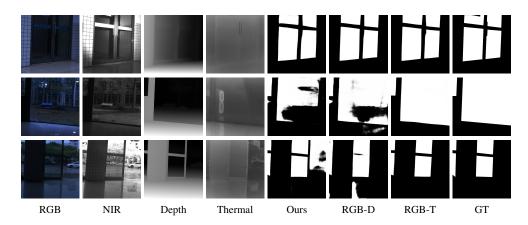


Figure 18: Study on different multi-modal inputs in low-light environment.





(a) Our designed hybrid imaging system.

(b) RGB-T and RGB-D imaging system.

Figure 19: (a) shows our designed hybrid imaging system for capturing RGB-NIR image pairs. It consists of a DSLR camera and a NIR camera. (b) shows the hybrid imaging system consisting of a thermal infrared camera and a stereo camera we used for capturing thermal images and depth maps.

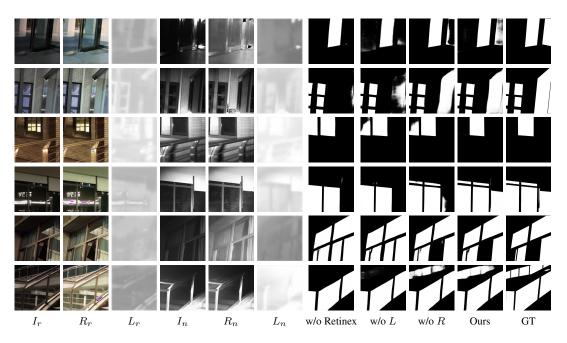


Figure 20: Ablation Study on Retinex Decomposition.

small. The 2nd scene of Fig. 18 shows that thermal image is susceptible to interference from external strong thermal radiation sources, such as the air conditioner outdoor unit reflected on glass surface in the thermal image. The 3rd scene of Fig. 18 shows that depth values of the glass door is incorrect, which are close to the depth values of the outdoor scene. Thus, RGB-D based method (Lin et al., 2025) cannot distinguish glass door and open glass door.

These examples demonstrate that our cpatured RGB-NIR image pair can provide more valuable information than RGB-D and RGB-T image pairs for glass surface detection, and our method outperforms those RGB-D and RGB-T based methods in various challenging night-time scenes.

A.5 STUDY ON MISALIGNMENT INPUTS

To evaluate the robustness of our model under practical conditions, we conducted experiments with synthetically misaligned RGB-NIR inputs. Specifically, we introduced random spatial shifts to the NIR images to simulate imperfect camera calibration and synchronization. As the offset increased, performance gradually degraded: a small misalignment (about 5 pixels) led to an IoU drop of only 1.1, while a larger shift (about10 pixels) caused a 2.4 decrease in IoU, as shown in Tab. 8. These results indicate that our model is reasonably tolerant to moderate spatial misalignment, which is a common issue in real-world dual-camera systems, and can still maintain acceptable performance under imperfect alignment.

Table 8: The impact of misaligned images.

pixel range	IoU↑	$F_{\beta} \uparrow$	MAE↓	BER↓	ACC↑
0-0 pixel	87.98	0.934	0.047	0.055	0.936
0-5 pixel	86.88	0.928	0.051	0.060	0.929
0-10 pixel	85.58	0.923	0.054	0.067	0.918

A.6 STUDY ON GLASS OBJECT DETECTION

The glass object detection (GOD) task substantially differs from the glass surface detection (GSD) task, due to the distinct properties of their respective targets. Specifically, GOD focuses on detecting glass objects characterized by their shapes and boundaries, whereas GSD emphasizes glass surfaces, which are defined by reflection and transmission phenomena.

We evaluate our method on the Trans-10K dataset (Xie et al., 2020), as shown in Tab. 9. This dataset contains two subsets of transparent objects: (1) Transparent things, referring to small-scale curved objects such as cups, bottles, and glasses; and (2) Transparent stuff, referring to large-scale surfaces such as windows, glass walls, and doors.

Table 9: Quantitative comparison between our method and TransLab (Xie et al., 2020) on the dataset Trans-10k. Best detection results are marked in **bold**.

Dataset	Methods	IoU↑	$F_{\beta} \uparrow$	MAE↓	BER↓	ACC↑
Trans10k-full	TransLab Ours	88.18 92.34	0.959	0.063 0.022	0.050 0.030	0.927 0.960
Trans10k-things	TransLab Ours	90.87 94.37	0.967	0.011	0.036 0.018	0.972
Trans10k-stuff	TransLab Ours	84.39 89.47	0.948	0.037	0.072 0.047	0.943

A.7 ABLATION STUDY OF RETINEX DECOMPOSITION

As shown in Fig. 20, apart from the case when Retinex Decomposition is not used (*i.e.*, decomposed components are replaced by the original input), we also test the cases when either one of the decomposed components is replaced by the original input image. We show more cases that are not shown in our paper. The reflection on glass surfaces is enhanced in the decomposed Reflectance component (R_T) of the RGB image, while there is no easily noticeable reflection on R_n . These cases

 demonstrate that Retinex Decomposition is helpful in distinguishing glass surfaces from non-glass regions.

A.8 ABLATION STUDY OF THE PROPOSED MODULES

The GSD task requires better global perception to capture glass surface properties. Previous methods always design modules by integrating features from multiple encoder layers or by applying convolutions with different kernel sizes. Our RNFL module adopts a multi-stage cross-attention design to capture long-range dependencies and modality-specific cues between RGB/NIR features. In contrast, simple fusion methods, e.g., convolutions, can only provide a limited receptive field, i.e., not very effective in exploiting such relationships.

We ablated RNFL by using a multi-stage convolutional fusion. Specifically, we first concatenate the RGB/NIR features (X_r^i/\bar{X}_n^i) and fuse them using a 3×3 convolution. The fused features are then concatenated with the decoder features X_{de}^{i+1} , followed by another 3×3 convolution for channel reduction, and finally combined through a residual connection.

This variant leads to an increased model size (50M parameters) and noticeable performance degradation: IoU decreased by 1.28 (87.98 \rightarrow 86.70), F_{β} by 0.021 (0.934 \rightarrow 0.913), ACC by 0.008 (0.936 \rightarrow 0.928), while MAE and BER increased by 0.005 (0.047 \rightarrow 0.052) and 0.006 (0.055 \rightarrow 0.061), respectively. These results highlight that our attention-based RNFL design is both more accurate and more parameter-efficient.

A.9 ABLATION STUDY OF THE RESNET BACKBONE

When the RNGE module is ablated from the variant of our network taking ResNet-50 as the backbone, the performance drop is more pronounced compared to our complete network taking Swin Transformer V2 as the backbone. As shown in Tab. 10, IoU decreased from 82.63 to 78.65, F_{β} dropped from 0.901 to 0.879, ACC declined from 0.905 to 0.883, while MAE increased from 0.070 to 0.086 and BER increased from 0.081 to 0.099.

Table 10: Results with ResNet-50 backbone and RNGE module ablation

Methods	IoU↑	$F_{\beta} \uparrow$	MAE↓	BER↓	ACC↑
w/o RNGE	78.65	0.879	0.086	0.099	0.883
full (taking ResNet-50 as the backbone)	82.63	0.901	0.070	0.081	0.905

A.10 Two Failure Cases

Our method does have limitations. The first row of Fig. 21 shows that our method over-detects the bottom-right non-glass region (an open window) as the glass region, as this region contains similar patterns to the upper-right glass region in both the RGB and NIR modalities. The second row of Fig. 21 shows that if a glass region is too dark for the retinal decomposition method to obtain the desired reflection component, our method may not detect it.

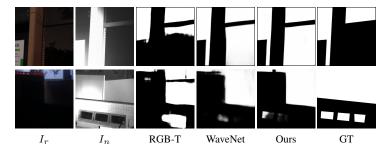


Figure 21: Two failure cases of our proposed method.