

RFTrans: Leveraging Refractive Flow of Transparent Objects for Surface Normal Estimation and Manipulation

Tutian Tang^{1*}, Jiyu Liu^{1*}, Jieyi Zhang¹, Haoyuan Fu¹, Wenqiang Xu¹ and Cewu Lu²

Abstract—Transparent objects are widely used in our daily lives, making it important to teach robots to interact with them. However, it's not easy because the reflective and refractive effects can make depth cameras fail to give accurate geometry measurements. To solve this problem, this paper introduces RFTrans, an RGB-D-based method for surface normal estimation and manipulation of transparent objects. By leveraging refractive flow as an intermediate representation, the proposed method circumvents the drawbacks of directly predicting the geometry (*e.g.* surface normal) from images and helps bridge the sim-to-real gap. It integrates the RFNet, which predicts refractive flow, object mask, and boundaries, followed by the F2Net, which estimates surface normal from the refractive flow. To make manipulation possible, a global optimization module will take in the predictions, refine the raw depth, and construct the point cloud with normal. An off-the-shelf analytical grasp planning algorithm is followed to generate the grasp poses. We build a synthetic dataset with physically plausible ray-tracing rendering techniques to train the networks. Results show that the proposed method trained on the synthetic dataset can consistently outperform the baseline method in both synthetic and real-world benchmarks by a large margin. Finally, a real-world robot grasping task witnesses an 83% success rate, proving that refractive flow can help enable direct sim-to-real transfer. The code, data, and supplementary materials are available at <https://rftrans.robotflow.ai>.

Index Terms—Perception for Grasping and Manipulation, RGB-D Perception

I. INTRODUCTION

TRANSSPARENT objects like glass goblets and plastic bottles are widely used in our daily lives. While humans can easily interact with transparent objects, teaching robots to manipulate them is not straightforward. One of the main barriers lies in the perception part, for transparency implies several special physical properties, including reflection, refraction, and the absence of color and texture. These properties

Manuscript received: October 20, 2023; Revised December 12, 2023; Accepted January 29, 2024.

This paper was recommended for publication by Editor Markus Vincze upon evaluation of the Associate Editor and Reviewers' comments. This work was supported by the National Key R&D Program of China (No. 2021ZD0110704), Shanghai Municipal Science and Technology Major Project (2021SHZDZX0102), Shanghai Qi Zhi Institute, and Shanghai Science and Technology Commission (21511101200)

*Equal contribution.

¹Tutian Tang, Jiyu Liu, Jieyi Zhang, Haoyuan Fu and Wenqiang Xu are with School of Electronic Information and Electrical Engineering, Shanghai Jiao Tong University, Shanghai, China {tutang, waterloo-sunset, yi_eagle, simon-fuhaoyuan, vinjohn}@sjtu.edu.cn

²Cewu Lu is the corresponding author, a member of Qing Yuan Research Institute and MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai, China lucewu@sjtu.edu.cn

Digital Object Identifier (DOI): see top of this page.

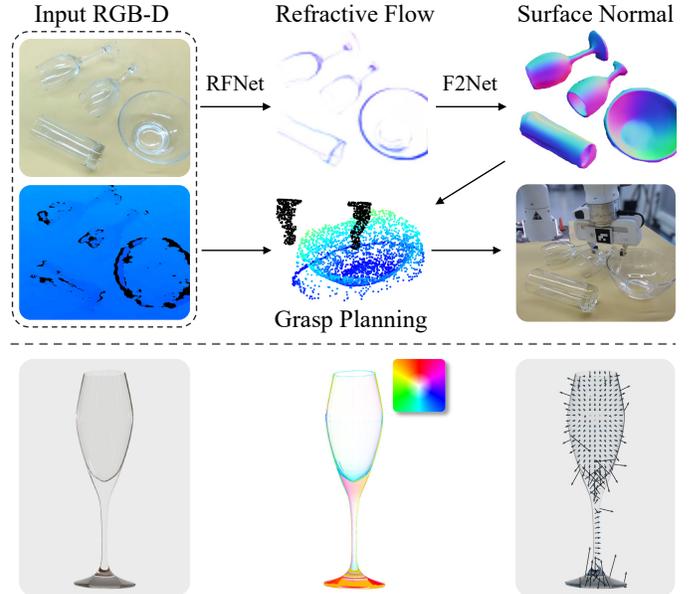


Fig. 1. **Top:** Transparency can cause inaccurate and missing depth captured by those widely-used RGB-D cameras. We utilize refractive flow to recover the surface normal and finally get the point cloud for robot manipulation. **Bottom:** (Left) A common wine glass. (Middle) We visualize the refractive flow by color map. The color represents the direction and magnitude. White indicates no refraction on the pixel. (Right) We sample some points on the image and show the corresponding refractive flow as arrows, which start from the foreground pixels on the glass to their corresponding pixels on the background.

often cause most 3D sensors, including LiDAR and RGB-D cameras, to fail to produce accurate geometry measurements for transparent objects [1]. Consequently, mainstream methods [2], [3] for transparent object manipulation typically adopt a two-stage pipeline. In the first stage, neural networks are used to recover the geometry of transparent objects based on noisy and inaccurate readings from commercial RGB-D cameras. Then, in the second stage, manipulation algorithms can be applied to the point cloud derived from the estimated geometry. Since these two stages are carried out sequentially, ensuring high-quality reconstructed geometry becomes crucial for reliable manipulation.

The recovery of transparent object geometry has been studied for decades. Previous approaches mostly relied on tracking the delicate light path using images from multiple viewpoints [4]–[6]. ClearGrasp [3] pioneered single image-based surface normal estimation, which directly predicts the surface normal from the RGB images. However, such direct

prediction can be problematic in two aspects. First, the correlation between RGB patterns and the underlying geometry is insignificant, especially when the background is complex. Second, obtaining the accurate surface normal in the real world is hard, so researchers must generate large amounts of synthetic data to train and evaluate the networks. In this case, the sim-to-real transferability becomes a challenging problem.

In this work, we propose **RFTrans**, which achieves surface normal estimation and manipulation of transparent objects based on a single RGB-D image. The main character distinguishing RFTrans from other works is that it uses the physical property, refraction, by adopting refractive flow [7] as an intermediate representation to mitigate the challenges of predicting surface normal directly from RGB images. Illustrated in Figure 1, the refractive flow is a per-pixel offset map between the transparent foreground and non-transparent background pixels to model the refractive effect of transparent objects. On each pixel, the refractive flow forms a 2D offset vector. We find refractive flow a good intermediate representation, for it features several merits, such as *small sim-to-real gap*, *stable under different ambient lights*, and *insensitive to complex background*. The details will be discussed in Section III-B.

In RFTrans, the **RFNet** first predicts refractive flow, object mask, and boundary based on RGB images. The **F2Net** then estimates the surface normal based on refractive flow. Then, these geometry-related elements, along with the original depth map, will be fed into a global optimization module as in [3], which will output the point cloud of transparent objects with normal to be fed into an off-the-shelf grasp planning algorithm. Here, we use ISF [8], an analytical grasp planning algorithm built on top of a heuristic surface matching metric. It can generate the grasp pose and guide the robot’s execution.

To train and evaluate the proposed method, we build a synthetic dataset of 62 transparent objects with RFUniverse [9], which features the latest simulation and rendering technologies and can thus generate physically plausible images. In addition, we directly test our model trained with synthetic data in a real-world benchmark [3]. Results show that our method outperforms the baseline method consistently. In the real-world grasp task, the success rate increases to 83% from 35% after the proposed method is applied.

Our main contributions can be summarized as follows:

- We propose RFTrans, a pipeline for surface normal estimation and manipulation of transparent objects based on RGB-D images. Refractive flow is used as an intermediate representation to help reconstruct accurate surface geometry.
- We construct a synthetic dataset of 62 transparent objects. The data generation pipeline is fully open-sourced.
- We set up a real-world grasping task to prove that RFTrans can enable direct sim-to-real transfer.

II. RELATED WORK

Our proposed method is most closely related to those approaches that utilize the refractive properties for transparent object geometry estimation and manipulation.

A. Refractive Property Estimation for Transparent Objects

Refractive flow describes the refractive properties of transparent objects. Initial work in this domain [10]–[13] involved the reconstruction of water surfaces by positioning a pre-defined pattern beneath a water tank and leveraging an optical flow-based algorithm to find the corresponding points between the camera pixels and the points from the pattern. Later methods [14], [15] remove the dependence on the water tank but impose some assumptions on the geometry of the objects. Some other correspondence-based methods [16], [17] measure the refractive flow without any prior of the geometry of the transparent objects. A more recent work, TOM-Net [7], estimates the refractive flow with neural networks trained on synthetic data, but it’s for image composition and matting instead of geometry recovery. There are many different methods to get refractive flow [14]–[16]. In this work, we model the refractive flow in a gray code-based approach [18] for its ease of use.

B. Estimating Geometry of Transparent Objects

Estimating the surface geometry of transparent objects has long been a challenge [1], [19]. Some methods [20]–[22] necessitate direct, intrusive interaction with the objects, keeping them unsuitable for robotic manipulation tasks. Murali *et al.* [23] introduce the non-intrusive tactile modality for reconstruction without damaging the objects.

Conversely, vision-based methods exploit the phenomena of reflection and refraction to deduce the underlying geometric attributes. Since specular reflection usually only happens in a small area of the transparent object body, reflection-based approaches traditionally either involve much human labor [24], require the periodic movement of the lighting source [25], or use multiple camera views [12]. Refraction, however, can be observed through nearly the whole body of a transparent object, providing much more information about surface geometry than reflection. Therefore, Kutulakos and Steger [15] utilize a combination of both phenomena to help reconstruct the objects.

The past decade has seen tremendous progress in neural networks. Stets *et al.* [26] and Sajjan *et al.* [3] are the pioneers of applying deep learning techniques in this field. Both methods utilize neural networks to estimate the surface normal of transparent objects directly from a single RGB image. Later, researchers push these methods into end-to-end depth restoration pipelines by using CNN-based networks [27], transformer-based networks [28], or implicit representations [29]. In multi-view vision systems, a recent trend is to adopt neural radiance fields (NeRF) to represent the scenes and objects implicitly. Li *et al.* [30] are the pioneers in applying NeRF in transparent object reconstruction, followed by Dex-NeRF [31]. However, a minimum of 3×3 camera array system is required, which hinders its application in robotics. Kerr *et al.* [32] further improve its speed and remove the need for the camera array. Dai *et al.* [33] propose GraspNeRF, which leverages generalizable NeRFs to reduce the number of images required to reconstruct one single scene. However, the NeRF-based methods still require multiple sparse-view

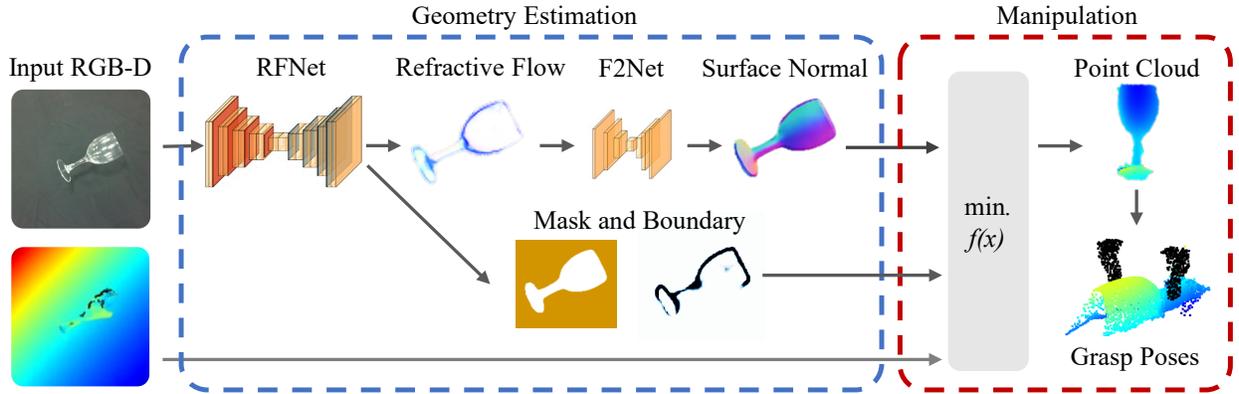


Fig. 2. Given an RGB-D image, RFNet first predicts the mask, the boundary, and the refractive flow of transparent objects. Next, F2Net will predict the surface normal based on the refractive flow. The global optimization will generate the singulated point cloud with normal. Finally, we apply the off-the-shelf manipulation algorithm, ISF, to generate grasp poses. The black points represent the fingers of the Franka Emika Panda robot.

RGB images to reconstruct the objects. All these methods consider the relationship between RGB images and the surface geometry as a black box without explicitly exploiting the phenomena of reflection and refraction. Our method first adopts the refractive flow as an intermediate representation, making it less sensitive to backgrounds (Sec. III-B3), less data-hungry (Sec. IV-B), and have better sim-to-real transferability (Sec. IV-A).

C. Transparent Object Manipulation

Transparent object manipulation is an important application area for geometry estimation. The upstream geometry estimation and the downstream manipulation algorithm can be either tightly or loosely coupled. Apparently, a loose coupling design allows the whole manipulation framework to benefit from the latest advances from both sides easily. ClearGrasp [3] is a typical loose coupling framework, followed by Jiang *et al.* in A4T [2], which sets a good example of bridging depth completion and manipulation via affordance. Later, Fang *et al.* [27] and Dai *et al.* [28] witness the direct improvement by applying the latest CNNs and transformers in depth completion. In Dex-NeRF [31], although the scene is represented implicitly via NeRF, the depth map is still extracted explicitly, leading the whole framework into a loose coupling manner. However, GraspNeRF [33] shows the tight coupling way, which introduces the Truncated Signed Distance Field (TSDF) as the bridge between implicit scene representation and grasping. There are also many data-driven approaches [34]–[36] aimed at directly generating grasping poses for transparent objects from noisy RGB-D images without explicit geometry estimation or depth completion. These deep-learning-based manipulation algorithms are usually data-hungry and not flexible enough to be a downstream grasp planner. First, they may fail to generalize between different models of RGB-D cameras. Second, if the user wants to add some new objects or adopt a new gripper of a different configuration, the networks must be retrained. Therefore, we decide to go for an analytical grasp planning algorithm, ISF [8], for the real-world grasping task.

III. METHOD

In this section, we will first give the definition and acquisition method of refractive flow (Sec. III-A). Next, we discuss some properties of the refractive flow to show the reason why we adopt it as the intermediate representation (Sec. III-B). Then, we describe how to estimate the geometry of transparent objects with neural networks (Sec. III-C). To train and evaluate the neural networks, we construct a synthetic dataset based on RFUniverse [9] (Sec. III-D). Finally, we will introduce our real-world grasping system in Section III-E. The whole pipeline is illustrated in Figure 2.

A. Refractive Flow: Definition and Acquisition

RGB-D cameras usually produce two kinds of errors on transparent objects [1]. Type I error happens when reflection happens and the camera fails to detect the depth value, resulting in incomplete depth maps. Type II error is highly related to the refractive effect, which occurs when the light refracts through the object’s surface and is reflected back by the non-transparent background. Figure 3 illustrates the relationship between refraction and Type II error.

Refractive flow is used to model the refractive effect of a transparent object. Intuitively, each *pixel* on the refractive flow is a 2D vector $(\Delta x, \Delta y)$, which indicates the offset between the foreground pixel and its refraction correspondence on the background image [7]. For those commonly-used thin-shell transparent objects, we can usually observe significant refraction near the edges. Therefore, the magnitude of the refractive flow is usually large near the boundaries, while the direction depends on a lot of factors including the structure of the object and the direction of the camera.

We use the gray code-based calibration method [18] to acquire the refractive flow. As illustrated in Figure 3, the transparent object is placed between a high-resolution LCD monitor and a fixed camera. The LCD monitor displays a sequence of 20 gray-coded images, including 10 vertical and 10 horizontal patterns. The camera is calibrated before the acquisition process [37], and it captures the corresponding images for the calibration process to get the refractive flow. Strong direct light should be avoided to prevent large reflective

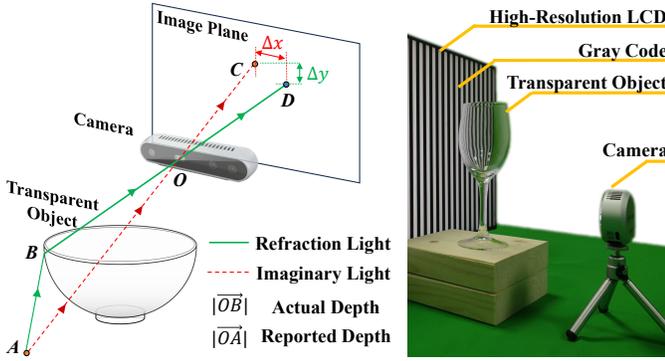


Fig. 3. **Left:** Point O is the optical center of the pin-hole camera. Point A is the point on the non-transparent background, *e.g.* a table. The refractive effect takes place at point B . \overline{AB} is the incident ray and \overline{BO} is the refracted ray. Point D is the image of A on the image plane. Without the transparent object, an imaginary ray will be directly from A to O , intersecting the image plane at point C . The orthogonal distances between C and D on the image plane, $(\Delta x, \Delta y)$, is the refractive flow at point D . Due to transparency, when the Type II error happens, RGB-D cameras usually report $proj|\overline{OA}|$ as the depth value of point D , while the actual depth should be $proj|\overline{OB}|$, where $proj \cdot |$ denotes the projected length on the principal axis. **Right:** The data acquisition system to capture refractive flow.

areas on the object. Additionally, for generating synthetic data, the system can be cloned into a digital twin in simulation, ensuring a small sim-to-real gap, as discussed in Section III-B1.

B. Properties of Refractive Flow

Apparently, the refractive flow will change according to the transparent objects' viewpoint and geometry. Therefore, it can encode rich information about the geometry. We find refractive flow features several merits to be a good intermediate representation for neural networks.

1) *Small Sim-to-Real Gap:* We adopt the gray-code calibration process to generate the refractive flow. Thanks to the latest ray tracing-based rendering technology, we can get photo-realistic and physically plausible images. Figure 4 shows one of the required images in the calibration process and the resulting refractive flow, both in simulation and in the real world. The sim-to-real gap is small. Also, we will show metrics on a real-world benchmark later in Section IV-A.

2) *Stable Under Different Ambient Light:* According to Snell's law [38]: $v_2 \sin \theta_1 = v_1 \sin \theta_2$, where v_1, v_2 are the phase velocities of light in two different media, and θ_1, θ_2 are the incidence and refraction angles respectively, we can conclude that different ambient lights can cause the gray code pattern shift. To quantitatively inspect the shift, we put a camera, an LED lamp, and a wine glass in front of a checkerboard. The color temperature of the LED lamp can be tuned from 2600K to 5000K. We compare the rooted mean squared error (RMSE) in pixels of the refractive patterns under the same viewpoint with different color temperatures. Results show that the RMSE value always stays around 0.1px on a change of color temperature. Therefore, the refractive flow calculated based on these patterns should also be stable. For details, please refer to the supplementary video.

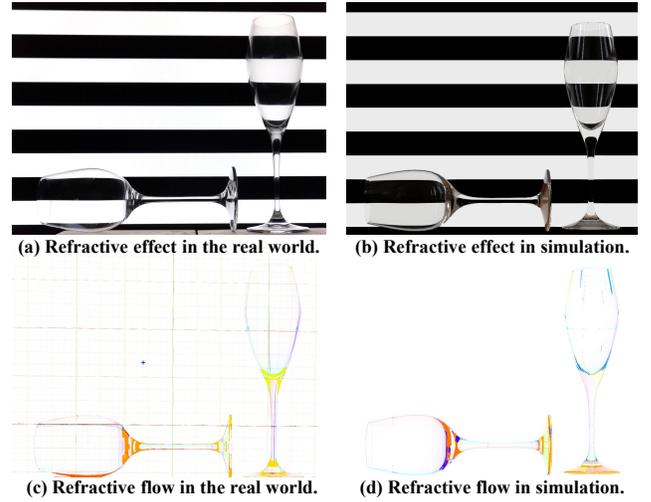


Fig. 4. The refractive effect and the corresponding refractive flow in the real world and our simulation environment.

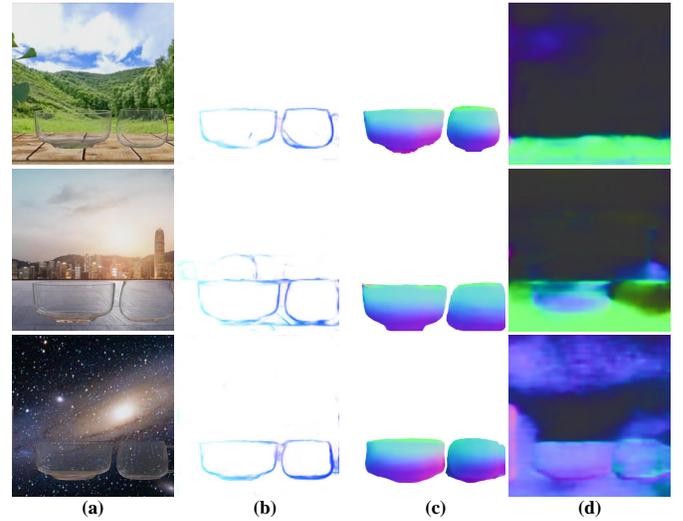


Fig. 5. (a): The same transparent objects are placed in front of different, complex backgrounds. (b): The refractive flow predicted by RFNet. (c): The surface normal derived from the refractive flow. (d): The direct surface normal prediction from RGB images *i.e.*, ClearGrasp [3]. Please note that RFTrans only predicts the surface normal of transparent objects, while ClearGrasp also predicts the background's normal.

3) *Insensitive to Complex Background:* Refractive flow models the pixel-to-pixel correspondence between the foreground and background, making it possible to handle the cases even when the background texture is complex. As shown in Figure 5, the surface normal from a refractive flow can still be rather stable with a complex background, while the surface normal predicted directly from RGB images as in [3] is not.

C. Geometry Estimation for Transparent Objects

As shown in Figure 2, taking in a single RGB-D image, **RFNet** first predicts the mask, boundary, and refractive flow for transparent objects. Next, the refractive flow is fed into another smaller, **F2Net** (Flow-to-Normal-Network) which estimates the surface normal. Following previous works [2], [3],

RFNet is composed by the DeepLabv3+ networks [39] with the DRN-D-54 backbone [40]. The output head is connected to the final feature layer, so the output size is identical to the input size. To note, although the RFNet can consist of any kind of CNN with the encoder-decoder structure, we stick with DeepLabv3+ because we want to keep the network architecture the same as the baseline method for a fair comparison. We use cross-entropy loss \mathcal{L}_{CE} for mask and boundary and mean squared-error loss \mathcal{L}_{flow} for refractive flow in the training process. Then, in the F2Net, since refractive flow already encodes rich information about the geometry, we can use a much smaller network like the original U-Net [41] to estimate the surface normal from the predicted refractive flow. We use the cosine-similarity loss \mathcal{L}_{norm} to supervise the surface normal. The loss is calculated over the pixels representing the transparent objects. Apparently, the two networks can be jointly optimized, which we will discuss in Section IV-C2.

D. Synthetic Dataset Construction

We use RFUniverse [9] to construct the synthetic training data. RFUniverse features accurate, physics-based rendering, which is confirmed in Figure 4. Figure 6 shows some randomly chosen samples of our dataset. Specifically, we collected 438 HDR sky-boxes and 40 textures for the tables. As for object models, the quality of the current large-scale object dataset [42] varies significantly from object to object, where many objects lack internal structures. Other datasets [27], [43] built by scanning real objects usually face the problem of wrecks and uneven surfaces. Such defects can lead to unrealistic reflective and refractive effects, making them unsuitable for synthetic dataset construction. By carefully examining the surface quality, geometry, and the resulting reflective and refractive effects, we manually select a total of 62 different CAD models of 5 categories, including glass bottles, wine glasses, glass cups, bowls, and plates. To generate one data sample, we first randomly select the sky-box and the texture of the table. Then, we randomly pick 1 to 5 objects and drop them onto the table. The object poses, surface normal, segmentation masks, and boundaries are recorded. As for depth, we collect the active infrared stereo-based synthetic depth [44], which simulates the Realsense camera and serves as the input. The ideal depth from the depth buffer serves as the ground truth. Finally, we use the gray-code calibration process to get the refractive flow. With one NVIDIA RTX 4090 GPU, the data generation pipeline runs at $\sim 5,000$ samples per hour, which is quite efficient. The assets are publicly available on our website.

E. Manipulation

After the geometry estimation step, we must construct the point cloud of transparent objects for downstream manipulation algorithms. The original depth map, along with the mask, boundary, and surface normal, are used as the input to a global optimization as in [3]. The global optimization will finally output the refined, singulated point cloud with normal. Since grasping is the foundation of prehensile manipulation, we follow the previous works [2], [3] to build a real-world grasping environment to demonstrate manipulation. We use

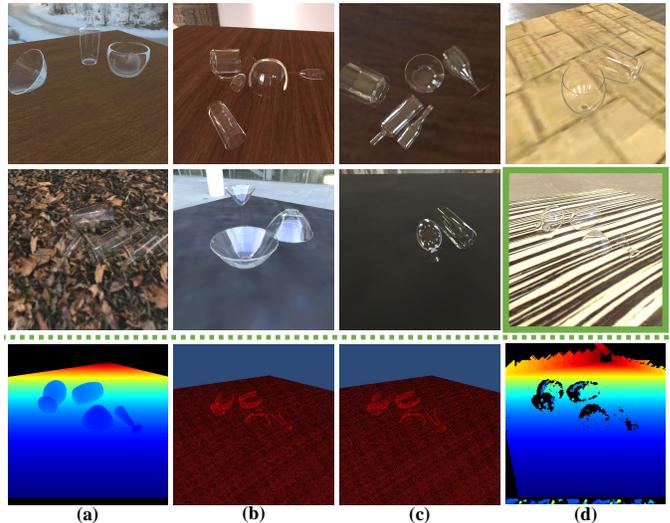


Fig. 6. **Top:** Samples of our synthetic dataset. **Bottom:** For the RGB image in the green frame (last image in the second row), (a) shows its ideal depth from the depth buffer, (b) and (c) show the simulated left and right IR images, and (d) shows the active depth generated by RFUniverse.

ISF [8] as a standalone and efficient grasp planner that can work with grippers of different DoFs. Taking in the singulated point cloud, the ISF will generate valid grasp poses and the corresponding energy-based metric. A lower energy suggests a better grasp pose, so we pick the one with the lowest energy. In our experiment, we adopt the top-down manner with a Franka Emika panda robot and its original gripper.

IV. EXPERIMENT

To evaluate RFTrans, we test it on both our proposed synthetic dataset, and the real-world benchmark proposed in ClearGrasp [3]. ClearGrasp is considered the baseline model and ImplicitDepth [29] is the previous *state-of-the-art* depth completion method. To prove that the RFTrans trained on synthetic datasets can be directly transferred into real-world applications, we design real-world transparent objects grasping environment and report success rates.

A. Surface Normal Estimation

Following previous works [2], [3], we evaluate the surface normal estimation performance by three kinds of metrics: the mean and median errors in degrees, and the percentages of pixels with errors below certain thresholds θ of 11.25, 22.5, and 30 degrees. Unless specified, metrics are calculated only over the pixels representing transparent objects.

The networks are trained with batch size 16 on a single RTX Titan GPU. The input and output sizes are fixed to 256 by 256. We use SGD optimizer with momentum 0.9 and weight decay $5e-4$. The learning rate is fixed to $1e-4$ during the whole 100 epochs. It takes less than 48 GPU hours to train the networks.

In Table I, test set *Syn* indicates the models are trained on our synthetic dataset of 5,000 images and tested on the synthetic test set of 1,000 images. The proposed method achieves the best result on all metrics. ImplicitDepth is an end-to-end depth completion method, which does not predict

TABLE I

COMPARISON WITH OTHER METHODS ON SURFACE NORMAL ESTIMATION. TEST SET SYN MEANS THE SYNTHETIC TEST SET AND REAL MEANS THE REAL-KNOWN TEST SET FROM THE CLEARGRASP BENCHMARK. \downarrow SUGGESTS THE LOWER THE BETTER AND VICE VERSA.

Test Set	Method	mean \downarrow	med. \downarrow	$\theta_{11.25} \uparrow$	$\theta_{22.5} \uparrow$	$\theta_{30} \uparrow$
Syn	Ours	11.10	6.94	71.18	85.51	89.18
	ClearGrasp	22.89	17.12	37.56	63.17	72.15
	ImplicitDepth	30.56	22.10	23.89	52.23	63.23
Real	Ours	29.16	24.87	22.94	48.68	61.32
	ClearGrasp	34.76	30.64	17.74	41.54	53.34
	ImplicitDepth	33.33	28.75	12.17	36.61	53.15

surface normal directly, but calculates the normal from the predicted depth (i.e., point cloud). This may explain why it does not perform well. We also measure the average speed. On a single RTX Titan GPU, our method runs at 31 fps. In comparison, ClearGrasp runs at 33 fps and ImplicitDepth runs at 10 fps.

To evaluate the sim-to-real performance, we benchmark the proposed method on the *real-known* test set proposed in ClearGrasp. Instead of using the synthetic train set from ClearGrasp, we have to build our own synthetic train set due to the necessity of refractive flow. Although ClearGrasp does not make its data-generation pipeline open-sourced, it makes the 9 CAD models publicly available, which enables us to rebuild a train set containing 5,000 images of the 9 models. The data generation pipeline is as described in Section III-D. In Table I, test set *Real* indicates the models are retrained on the synthetic dataset of the 9 objects and are tested on ClearGrasp *real-known* test set. Again, the proposed method achieves the best result on all metrics, which further proves that the refractive flow is a good intermediate representation that can enable a direct sim-to-real transfer.

B. Depth Completion

We benchmark the depth completion result to see if the improved surface normal estimation can result in better depth completion in the real world. For fair comparison, all models are trained on the proposed synthetic dataset and tested on ClearGrasp *real-known* test set. Following previous works [3], [29], we evaluate the depth completion performance by four kinds of metrics: the root mean squared error (RMSE), the absolute relative difference (REL), the mean absolute error (MAE) and the percentages of pixels with errors below certain thresholds δ of 1.05, 1.10, and 1.25. In Table II, results show that the proposed method can produce the lowest mean error and it performs significantly better than others on the most strict $\delta_{1.05}$ metric.

To note, the authors of ClearGrasp report its performance with a heavy pretraining process on 80k images from the Scannet [45] and Matterport3D [46] datasets in the paper [3]. The proposed method can produce comparable results with only 5k images for training, which proves that introducing refractive flow as the intermediate representation can help make the model significantly less data-hungry. Also, the original paper of ImplicitDepth reports that it can outperform

TABLE II

COMPARISON WITH OTHER METHODS ON DEPTH COMPLETION. \downarrow SUGGESTS THE LOWER THE BETTER AND VICE VERSA. ALL MODELS ARE RETRAINED ON THE SYNTHETIC DATASET AND TESTED ON THE CLEARGRASP REAL-KNOWN BENCHMARK FOR FAIR COMPARISON.

Method	RMSE \downarrow	REL \downarrow	MAE \downarrow	$\delta_{1.05} \uparrow$	$\delta_{1.10} \uparrow$	$\delta_{1.25} \uparrow$
Ours	0.038	0.059	0.031	71.46	89.69	96.93
ClearGrasp	0.045	0.071	0.036	53.23	79.01	95.94
ImplicitDepth	0.048	0.080	0.043	28.79	62.55	98.45

TABLE III

ABLATION STUDY ON DIFFERENT INPUT SOURCES TO ESTIMATE SURFACE NORMAL. \downarrow SUGGESTS THE LOWER THE BETTER AND VICE VERSA.

Source	mean \downarrow	med. \downarrow	$\theta_{11.25} \uparrow$	$\theta_{22.5} \uparrow$	$\theta_{30} \uparrow$
RGB	22.89	17.12	37.56	63.17	72.15
Boundary	21.98	16.12	41.24	65.53	73.94
Flow	11.10	6.94	71.18	85.51	89.18

ClearGrasp with its Omniverse Object training dataset of 60k images, while in our experiment with only 5k images, it cannot. This may indicate the size of training set is important for ImplicitDepth.

C. Ablation Study

In this section, we evaluate the effect of adopting refractive flow, analyze the effects of jointly optimizing the RFNet and F2Net, and check the sensitivity of different viewpoints.

1) *Refractive Flow*: To evaluate the effect of refractive flow as the intermediate representation, we design three different networks that predict surface normal from different input sources. In Table III, source *RGB* means we directly use the RFNet to predict the surface normal from RGB input, just like what ClearGrasp does. It achieves the worst result. Source *Boundary* means we use the boundary of the objects as the intermediate representation. To be specific, we use RFNet to first predict the boundary, followed by the modified F2Net to predict surface normal. It can be treated as the variant of RFTrans. Source *Flow* means the proposed RFTrans. Results show the boundary can also act as the intermediate representation for surface normal prediction, for the related model outperforms ClearGrasp slightly. Refractive flow achieves the best, which proves that it can encode more information than the plain boundary.

2) *Joint Optimization*: In surface normal estimation via refractive flow, the RFNet and the F2Net can be optimized jointly or separately. Recall that separately, we use \mathcal{L}_{flow} for refractive flow and \mathcal{L}_{norm} for surface normal. In the joint manner, we use $\mathcal{L} = \alpha\mathcal{L}_{flow} + \mathcal{L}_{norm}$ as the final loss. α is the coefficient to balance the two terms. Here, we quantitatively evaluate the effectiveness of joint optimization of RFNet and F2Net. The networks are trained and tested on the synthetic dataset. Table IV shows the joint optimization (denoted by End2End) can significantly boost the performance, compared with the separate optimization. Further, we also analyze the effects of different values of α . We vary the value of α from 0.001 to 1.5 in the log scale. Results show that we can slightly benefit from tuning this parameter and $\alpha = 0.01$ is the best.

TABLE IV
ABLATION STUDY ON JOINT OPTIMIZATION AND THE VALUE OF α . \downarrow
SUGGESTS THE LOWER THE BETTER AND VICE VERSA.

Method	α	mean \downarrow	med. \downarrow	$\theta_{11.25}$ \uparrow	$\theta_{22.5}$ \uparrow	θ_{30} \uparrow
Separate	-	18.03	12.29	56.51	75.47	81.19
End2End	0.001	12.23	7.79	68.21	83.59	87.65
End2End	0.01	11.10	6.94	71.18	85.51	89.18
End2End	0.1	11.65	7.36	70.08	84.45	88.24
End2End	1.0	12.33	7.88	67.86	83.35	87.53
End2End	1.5	12.5	8.08	67.75	83.10	87.29

TABLE V
SENSITIVITY OF VIEWPOINTS. \downarrow SUGGESTS THE LOWER THE BETTER AND
VICE VERSA.

Train Set	Test Set	mean \downarrow	med. \downarrow	$\theta_{11.25}$ \uparrow	$\theta_{22.5}$ \uparrow	θ_{30} \uparrow
High-View	Low-View	22.78	18.23	36.92	62.68	72.57
Low-View	Low-View	10.74	6.66	72.09	86.07	89.62

3) *Sensitivity of Viewpoints*: In order to test the model’s sensitivity on different viewpoints, we select 2,000 images from the synthetic train set where the pitch angle of the camera stays between 75 and 90 degrees. The derived subset is named *High-Viewpoint* train set. We also select images with the pitch angle between 30 and 75 degrees to get the *Low-Viewpoint* train set. Similarly, we derive the *Low-Viewpoint* test set. In Table V, results show that on the *Low-Viewpoint* test set, the network trained on the *High-Viewpoint* set performs significantly worse than the *Low-Viewpoint* one. This experiment confirms that refractive flow is subjective to viewpoints. Therefore, in order to get best performance, the training set should cover the viewpoint for inference.

D. Manipulation

To test RFTrans in real-world manipulation scenarios, we set up a robotic system consisting of a 7-DoF robot arm with a parallel gripper by Franka Emika and an Intel Realsense D415 RGB-D camera. We use the *easy-handeye* software package to calibrate hand-eye. We use the success rate as the evaluation metric. We collect 10 glass objects for grasping, one of which is color tinted. We set up 10 scenes for evaluation. In each scene, we randomly place 3 to 5 transparent objects. In each attempt, the robot should pick up an object according to the grasp pose generated by the grasp planning algorithm, ISF [8]. An attempt is regarded as successful if the object is lifted at least 20 cm above the table. We shift to a new scene once all objects in the current scene are picked up, or a maximum of 5 trails is reached. Table VI reports the success rate of our method compared with the baseline methods. Method *Raw Realsense* means we use the raw depth value to directly construct the point cloud that ISF estimates grasp poses based on. Method *ClearGrasp* means the point cloud is generated by ClearGrasp [3]. With the help of RFTrans, the success rate is increased from 35% to 83%, outperforming the baseline method. There are numerous factors leading to the failure of grasping. First, most transparent objects are composed of glass - a material renowned for its smoothness that often results in grippers losing hold of the object. Second, the global

TABLE VI
SUCCESS RATE OF MANIPULATION.

Method	#trials	#success	Rate
Raw Realsense	48	17	35%
ClearGrasp	48	35	72%
RFTrans	48	40	83%

optimization algorithm [3] may fail to accurately estimate the depth of the parts not in contact with the table, such as the stem of a wine glass, leading to incorrect point clouds and subsequent grasping failures. Lastly, the gripper may contact with the object prematurely as it approaches, causing it to shift and ultimately leading to failure. This issue is related to the trajectory planning algorithm of the robot. Please refer to the supplementary video for details.

V. CONCLUSIONS, LIMITATIONS, AND DISCUSSION

In this paper, we propose a framework that utilizes refractive flow to estimate the geometry and manipulate the transparent objects. Experiments show that refractive flow is a good intermediate representation that can lead to better surface normal estimation, which benefits manipulation. Besides, the networks trained on synthetic data can be used directly for manipulation tasks in the real world, which confirms refractive flow can help the direct sim-to-real transfer.

However, these merits only happen with those commonly used transparent objects featuring thin-shell structures in our daily lives. In an extreme case, with a triangular prism, it’s difficult to even get the precise refractive flow due to the dispersion effect. Moreover, our method may fail when objects heavily overlap with each other.

We hope our method can arouse people’s interest in utilizing physical properties in this track. Future works can extend this method to more complex objects and tasks or find other interesting physical properties that benefit manipulation. Also, since the ISF algorithm supports grippers with high DoFs, our proposed method can be easily extend to a dexterous grasping pipeline.

REFERENCES

- [1] J. Jiang, G. Cao, J. Deng, T.-T. Do, and S. Luo, “Robotic perception of transparent objects: A review.” [Online]. Available: <http://arxiv.org/abs/2304.00157>
- [2] J. Jiang, G. Cao, T. Do, and S. Luo, “A4t: Hierarchical affordance detection for transparent objects depth reconstruction and manipulation,” *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 9826–9833, Oct. 2022.
- [3] S. S. Sajjan, M. J. Moore, M. Pan, G. Nagaraja, J. Lee, A. Zeng, and S. Song, “Clear grasp: 3d shape estimation of transparent objects for manipulation,” *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3634–3642, 2019.
- [4] Y. Qian, M. Gong, and Y.-H. Yang, “3d reconstruction of transparent objects with position-normal consistency,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4369–4377.
- [5] K. Han, K.-Y. K. Wong, and M. Liu, “A fixed viewpoint approach for dense reconstruction of transparent objects,” *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [6] B. Wu, Y. Zhou, Y. Qian, M. Cong, and H. Huang, “Full 3d reconstruction of transparent objects,” *ACM Trans. Graph.*, vol. 37, no. 4, July 2018. [Online]. Available: <https://doi.org/10.1145/3197517.3201286>

- [7] G. Chen, K. Han, and K.-Y. K. Wong, "TOM-net: Learning transparent object matting from a single image," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, pp. 9233–9241. [Online]. Available: <https://ieeexplore.ieee.org/document/8579060/>
- [8] Y. Fan, H.-C. Lin, T. Tang, and M. Tomizuka, "Grasp planning for customized grippers by iterative surface fitting," in *2018 IEEE 14th International Conference on Automation Science and Engineering (CASE)*. IEEE, 2018, pp. 28–34.
- [9] H. Fu, W. Xu, R. Ye, H. Xue, Z. Yu, T. Tang, Y. Li, W. Du, J. Zhang, and C. Lu, "Demonstrating RFUniverse: A Multiphysics Simulation Platform for Embodied AI," in *Proceedings of Robotics: Science and Systems*, Daegu, Republic of Korea, July 2023.
- [10] H. Murase, "Surface shape reconstruction of an undulating transparent object," in *[1990] Proceedings Third International Conference on Computer Vision*, 1990, pp. 313–317.
- [11] S. Agarwal, S. P. Mallick, D. Kriegman, and S. Belongie, "On Refractive Optical Flow," in *Computer Vision - ECCV 2004*, T. Pajdla and J. Matas, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 483–494.
- [12] N. Morris and K. Kutulakos, "Dynamic refraction stereo," in *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, vol. 2, 2005, pp. 1573–1580 Vol. 2.
- [13] Y. Ding, F. Li, Y. Ji, and J. Yu, "Dynamic fluid surface acquisition using a camera array," in *2011 International Conference on Computer Vision*, 2011, pp. 2478–2485.
- [14] G. Wetzstein, D. Roodnick, W. Heidrich, and R. Raskar, "Refractive shape from light field distortion," in *2011 International Conference on Computer Vision*, 2011, pp. 1180–1186.
- [15] K. Kutulakos and E. Steger, "A theory of refractive and specular 3d shape by light-path triangulation," in *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, vol. 2, 2005, pp. 1448–1455 Vol. 2.
- [16] J. Xu, Z. Zhu, H. Bao, and W. Xu, "Hybrid Mesh-neural Representation for 3D Transparent Object Reconstruction," Mar. 2023, arXiv:2203.12613 [cs]. [Online]. Available: <http://arxiv.org/abs/2203.12613>
- [17] Y. Ji, J. Ye, and J. Yu, "Reconstructing gas flows using light-path approximation," 06 2013.
- [18] J. Lyu, B. Wu, D. Lischinski, D. Cohen-Or, and H. Huang, "Differentiable refraction-tracing for mesh reconstruction of transparent objects," *ACM Trans. Graph.*, vol. 39, no. 6, nov 2020.
- [19] I. Ihrke, K. Kutulakos, H. Lensch, M. Magnor, and W. Heidrich, "Transparent and specular object reconstruction," *Comput. Graph. Forum*, vol. 29, pp. 2400–2426, 12 2010.
- [20] B. Trifonov, D. Bradley, and W. Heidrich, "Tomographic reconstruction of transparent objects," in *ACM SIGGRAPH 2006 Sketches*, ser. SIGGRAPH '06. New York, NY, USA: Association for Computing Machinery, 2006, p. 55–es. [Online]. Available: <https://doi.org/10.1145/1179849.1179918>
- [21] M. B. Hullin, M. Fuchs, I. Ihrke, H.-P. Seidel, and H. P. A. Lensch, "Fluorescent immersion range scanning," *ACM Trans. Graph.*, vol. 27, no. 3, p. 1–10, aug 2008. [Online]. Available: <https://doi.org/10.1145/1360612.1360686>
- [22] K. Aberman, O. Katzir, Q. Zhou, Z. Luo, A. Sharf, C. Greif, B. Chen, and D. Cohen-Or, "Dip transform for 3d shape reconstruction," *ACM Trans. Graph.*, vol. 36, no. 4, jul 2017. [Online]. Available: <https://doi.org/10.1145/3072959.3073693>
- [23] P. K. Murali, B. Porr, and M. Kaboli, "Touch if it's transparent! ACTOR: Active Tactile-based Category-Level Transparent Object Reconstruction," July 2023, arXiv:2307.16254 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.16254>
- [24] N. J. W. Morris and K. N. Kutulakos, "Reconstructing the surface of inhomogeneous transparent scenes by scatter-trace photography," in *2007 IEEE 11th International Conference on Computer Vision*, 2007, pp. 1–8.
- [25] S.-K. Yeung, T.-P. Wu, C.-K. Tang, T. F. Chan, and S. Osher, "Adequate reconstruction of transparent objects on a shoestring budget," in *CVPR 2011*, 2011, pp. 2513–2520.
- [26] J. Stets, Z. Li, J. R. Frisvad, and M. Chandraker, "Single-shot analysis of refractive shape using convolutional neural networks," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2019, pp. 995–1003.
- [27] H. Fang, H.-S. Fang, S. Xu, and C. Lu, "Transcg: A large-scale real-world dataset for transparent object depth completion and a grasping baseline," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 7383–7390, 2022.
- [28] Q. Dai, J. Zhang, Q. Li, T. Wu, H. Dong, Z. Liu, P. Tan, and H. Wang, "Domain randomization-enhanced depth simulation and restoration for perceiving and grasping specular and transparent objects," in *European Conference on Computer Vision (ECCV)*, 2022.
- [29] 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," Apr. 2021, arXiv:2104.00622 [cs]. [Online]. Available: <http://arxiv.org/abs/2104.00622>
- [30] Z. Li, Y.-Y. Yeh, and M. Chandraker, "Through the Looking Glass: Neural 3D Reconstruction of Transparent Shapes," July 2020, arXiv:2004.10904 [cs]. [Online]. Available: <http://arxiv.org/abs/2004.10904>
- [31] J. Ichnowski*, Y. Avigal*, J. Kerr, and K. Goldberg, "Dex-NeRF: Using a neural radiance field to grasp transparent objects," in *Conference on Robot Learning (CoRL)*, 2020.
- [32] J. Kerr, L. Fu, H. Huang, Y. Avigal, M. Tancik, J. Ichnowski, A. Kanazawa, and K. Goldberg, "Evo-NeRF: Evolving NeRF for Sequential Robot Grasping of Transparent Objects," Aug. 2022. [Online]. Available: <https://openreview.net/forum?id=Bxr45keYrf>
- [33] Q. Dai, Y. Zhu, Y. Geng, C. Ruan, J. Zhang, and H. Wang, "Grasnerf: Multiview-based 6-dof grasp detection for transparent and specular objects using generalizable nerf," 2023.
- [34] T. Weng, A. Pallankize, Y. Tang, O. Kroemer, and D. Held, "Multi-modal transfer learning for grasping transparent and specular objects," *IEEE Robotics and Automation Letters*, vol. 5, no. 3, pp. 3791 – 3798, July 2020.
- [35] H. Cao, J. Huang, Y. Li, J. Zhou, and Y. Liu, "Fuzzy-depth objects grasping based on fsg algorithm and a soft robotic hand," 09 2021, pp. 3948–3954.
- [36] J. Chang, M. Kim, S. Kang, H. Han, S. Hong, K. Jang, and S. Kang, "Ghostpose: Multi-view pose estimation of transparent objects for robot hand grasping," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 5749–5755.
- [37] Z. Zhang, "A flexible new technique for camera calibration," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 11, pp. 1330–1334, 2000.
- [38] R. Feynman, R. Leighton, M. Sands, and E. Hafner, "The Feynman Lectures on Physics: Vol. I." AAPT, 1965, pp. "26–5"–"26–7".
- [39] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation," in *Computer Vision – ECCV 2018*, ser. Lecture Notes in Computer Science, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds. Cham: Springer International Publishing, 2018, pp. 833–851.
- [40] F. Yu, V. Koltun, and T. Funkhouser, "Dilated residual networks," in *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*. Institute of Electrical and Electronics Engineers Inc., Nov. 2017, pp. 636–644. [Online]. Available: <https://collaborate.princeton.edu/en/publications/dilated-residual-networks>
- [41] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, ser. LNCS, vol. 9351. Springer, 2015, pp. 234–241, (available on arXiv:1505.04597 [cs.CV]). [Online]. Available: <http://lmb.informatik.uni-freiburg.de/Publications/2015/RFB15a>
- [42] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, "ShapeNet: An Information-Rich 3D Model Repository," Stanford University — Princeton University — Toyota Technological Institute at Chicago, Tech. Rep. arXiv:1512.03012 [cs.GR], 2015.
- [43] L. Liu, W. Xu, H. Fu, S. Qian, Q. Yu, Y. Han, and C. Lu, "Akb-48: A real-world articulated object knowledge base," in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 14 789–14 798.
- [44] X. Zhang, R. Chen, A. Li, F. Xiang, Y. Qin, J. Gu, Z. Ling, M. Liu, P. Zeng, S. Han, Z. Huang, T. Mu, J. Xu, and H. Su, "Close the Optical Sensing Domain Gap by Physics-Grounded Active Stereo Simulation," in *T-RO 2023*, 2023.
- [45] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "ScanNet: Richly-annotated 3d reconstructions of indoor scenes," in *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017.
- [46] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3d: Learning from rgb-d data in indoor environments," *International Conference on 3D Vision (3DV)*, 2017.