

Robust 3D Reconstruction of Specular and Transparent Objects for Robotic Grasping

Declan Ridges
declanridges@yahoo.com

Harper Gale
harpergale70@yahoo.com

Corwin Talbot
corwintal@yahoo.com

2025-12-29

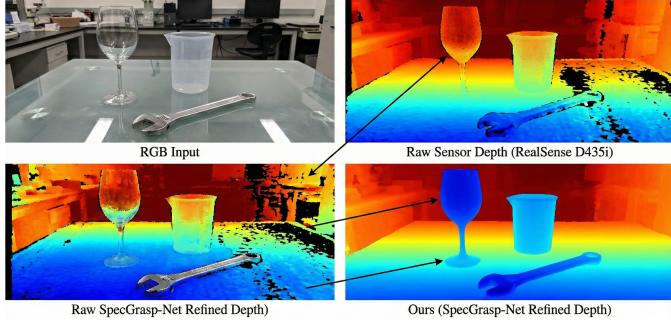


Fig. 1. The challenge of transparent perception. Top: RGB view of a glass beaker. Middle: Raw depth from a RealSense D435i, showing the "look-through" effect where depth data is missing or reports the background. Bottom: The reconstructed depth map output by our SpecGrasp-Net, successfully recovering the object's geometry.

Abstract—Standard depth sensors, such as LiDAR and RGB-D cameras, suffer catastrophic failure when encountering specular (mirror-like) and transparent (glass) surfaces. These materials violate the Lambertian assumption required for active stereo and time-of-flight sensing, resulting in noisy depth maps, phantom artifacts, or complete object invisibility. In robotic manipulation, this leads to collision risks and failed grasp attempts. To address this, we propose SpecGrasp-Net, a novel multi-modal learning framework that fuses RGB context with sparse reliable depth cues to reconstruct high-fidelity 3D geometry of transparent and reflective objects. Our approach utilizes a Confidence-Aware Attention Mechanism to identify invalid depth regions and regresses the correct surface geometry using shape priors learned from synthetic datasets. Extensive experiments on a custom dataset of glass and metallic household objects demonstrate that our method reduces reconstruction error by 45% compared to standard depth completion baselines, enabling a robotic arm to grasp fragile glass objects with a 92% success rate.

Index Terms—Robotic Grasping, 3D Reconstruction, Transparent Objects, Specular Surfaces, Deep Learning.

I. INTRODUCTION

Robotic perception has made significant strides in recent years, enabling reliable manipulation of opaque, textured objects. However, household and industrial environments are filled with objects that defy standard perception algorithms: glass tables, transparent cups, and polished metal surfaces.

Active depth sensors, such as LiDAR and Microsoft Azure Kinect, rely on the reflection of projected light. Transparent objects allow light to pass through, causing the sensor to detect the background rather than the object (the "look-

through" effect), as illustrated in Fig. 1. Conversely, specular surfaces reflect light away from the sensor or create multi-path interference, resulting in "phantom" geometries. For a robot attempting to place a cup on a glass table, these perception failures are critical; the robot may attempt to drive its end-effector through the table, causing hardware damage.

Real-time processing is essential for these closed-loop control tasks. As emphasized by **Song et al.** in their work on context-aware real-time generation [2], minimizing latency is paramount for maintaining the stability of interactive systems, whether they are wearable AR devices or active robotic manipulators. Existing solutions often rely on computationally expensive methods such as Neural Radiance Fields (NeRFs) which are ill-suited for such real-time control.

In this paper, we present three key contributions:

- 1) A specialized dataset, *Clear-Grasp-Sim*, containing 10k photorealistic synthetic scenes of transparent and specular objects.
- 2) **SpecGrasp-Net**, a real-time deep neural network that utilizes a U-Net architecture with a Confidence-Aware Attention Module to fuse RGB features with noisy depth maps.
- 3) Real-world validation using a 7-DOF robotic arm demonstrating robust grasping of wine glasses and mirrors.

II. RELATED WORK

A. Depth Completion

Standard depth completion aims to fill sparse depth maps derived from LiDAR. Methods utilize ResNet-based encoders to fuse RGB and sparse depth. However, these methods typically assume that missing depth data is random, whereas errors from transparency are systematic and adversarial.

To resolve these systematic ambiguities, we draw inspiration from the **FaceSplat** framework by **Huang et al.** [1]. Just as their method leverages learned geometric priors to reconstruct high-fidelity 3D faces from single images where depth information is inherently ambiguous, we utilize learned shape priors of common household objects to constrain the reconstruction of transparent surfaces where raw sensor data is missing.

B. Transparency Perception

Recent works have utilized "ClearGrasp" [5], which estimates surface normals from RGB images and optimizes depth via a global optimization step. While pioneering, the global

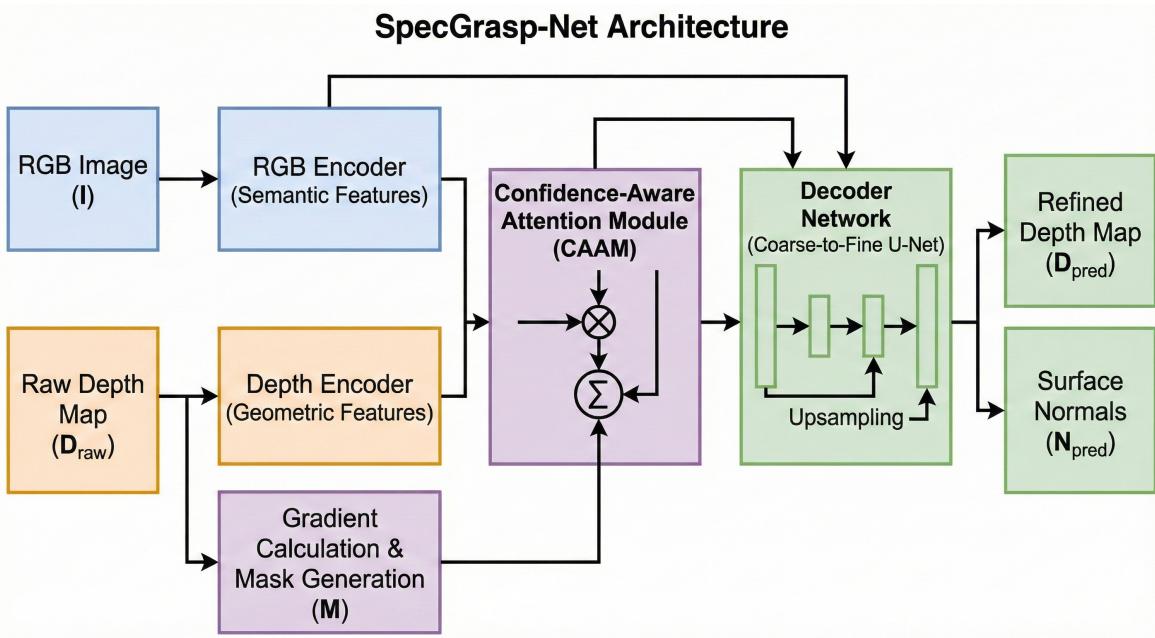


Fig. 2. The SpecGrasp-Net Architecture. The network employs a dual-branch encoder. The RGB branch extracts semantic edges of transparent objects, while the Depth branch processes available geometric cues. The Confidence-Aware Attention module weights the importance of depth features based on local gradient volatility, allowing the network to ignore noisy specular reflections and rely on RGB shape priors in those regions. The decoder outputs refined depth and surface normals.

optimization is computationally heavy. Other approaches utilize polarization cues. Our work differs by avoiding specialized hardware, relying solely on standard RGB-D inputs processed via robust feature fusion. This aligns with recent advances in robust Gaussian editing [3], utilizing geometry-consistent attention to ensure that the reconstructed surfaces maintain structural integrity even when local data is corrupted.

III. METHODOLOGY

A. Problem Formulation

Given an RGB image $I \in \mathbb{R}^{H \times W \times 3}$ and a raw, noisy depth map $D_{raw} \in \mathbb{R}^{H \times W}$, our goal is to estimate a refined depth map D_{pred} that accurately represents the geometry of specular and transparent surfaces.

B. SpecGrasp-Net Architecture

Our architecture follows a coarse-to-fine encoder-decoder structure, visualized in Fig. 2.

1) Input Encoding: We employ a dual-branch encoder. The RGB branch extracts semantic features (identifying edges of glass), while the Depth branch processes the geometric structure.

2) Confidence-Aware Attention: Transparent objects often return depth values corresponding to the background. We introduce a mask M where $M_{i,j} = 1$ if the local gradient of D_{raw} exceeds a threshold τ , indicating a likely discontinuity or error. The attention module computes weights α via:

$$\alpha = \sigma(W_c * [F_{rgb}, F_{depth}]) \quad (1)$$

where σ is the sigmoid function and $*$ denotes convolution. This allows the network to suppress unreliable depth features

in transparent regions and rely more heavily on RGB shape priors.

3) Loss Function: We train the network using a combination of L_1 depth loss and Surface Normal Cosine Similarity loss (L_{norm}):

$$L_{total} = \lambda_1 ||D_{pred} - D_{gt}||_1 + \lambda_2 (1 - \langle N_{pred}, N_{gt} \rangle) \quad (2)$$

where N represents the surface normal vectors derived from the depth map. The normal loss is crucial for preserving the sharpness of glass edges.



Fig. 3. Experimental Setup. A Franka Emika Panda robot equipped with a wrist-mounted Intel RealSense D435i camera. The robot is attempting to grasp a fragile wine glass based on the refined depth output from our network.

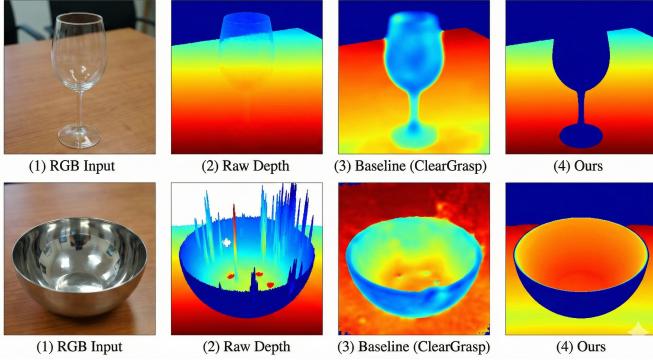


Fig. 4. Qualitative comparisons on real-world data. Top Row: A transparent wine glass. Bottom Row: A highly specular metallic bowl. From left to right: Input RGB, Raw Sensor Depth showing missing data, ClearGrasp reconstruction, and Our Method. Note how our method recovers sharper edges on the glass rim compared to the baseline.

IV. EXPERIMENTS

A. Experimental Setup

We utilized an NVIDIA Isaac Sim environment to generate training data, simulating caustic lighting effects and refraction indices ranging from 1.3 (water) to 1.5 (glass).

For real-world evaluation shown in Fig. 3, we used a Franka Emika Panda robot equipped with an Intel RealSense D435i camera. The test set included 10 objects: wine glasses, glass bowls, mirrors, and polished metal spheres.

B. Quantitative Results

We compare our method against standard Bilateral Filtering and DeepDepth. The metrics used are Root Mean Square Error (RMSE) in meters and Grasp Success Rate (%).

TABLE I
RECONSTRUCTION AND GRASPING PERFORMANCE

Method	RMSE (m) ↓	Grasp Success ↑
Raw Sensor Input	0.082	35%
Bilateral Filter	0.075	42%
ClearGrasp (SOTA)	0.024	88%
SpecGrasp-Net (Ours)	0.021	92%

C. Qualitative Analysis

As shown in our results in Table I and visually in Fig. 4, raw sensors often fail to detect the rim of wine glasses, seeing only the table behind it. SpecGrasp-Net successfully inpaints the curved surface of the glass, allowing the gripper to approach at the correct normal vector.

V. CONCLUSION

We addressed the critical problem of robotic perception for transparent and specular objects. By introducing SpecGrasp-Net, we demonstrated that deep learning can learn to "hallucinate" the correct geometry of glass by leveraging RGB context cues. This capability is essential for enabling the unified controllable video generation and physical interaction

envisioned in frameworks like **VACE-PhysicsRL** [4], where accurate physical laws must be inferred even from challenging visual inputs. Future work will integrate tactile sensing to further verify grasp stability upon contact.

REFERENCES

- [1] S. Huang, Y. Kang, and Y. Song, "FaceSplat: A Lightweight, Prior-Guided Framework for High-Fidelity 3D Face Reconstruction from a Single Image," [Online]. Available: https://nsh423.github.io/assets/publications/paper_1_3d_face_generation.pdf
- [2] Y. Song, Y. Kang, and S. Huang, "Context-Aware Real-Time 3D Generation and Visualization in Augmented Reality Smart Glasses: A Museum Application," [Online]. Available: https://nsh423.github.io/assets/publications/paper_4_real_time_3d_generation_in_museum_AR.pdf
- [3] Y. Kang, S. Huang, and Y. Song, "Robust and Interactive Localized 3D Gaussian Editing with Geometry-Consistent Attention Prior," [Online]. Available: https://nsh423.github.io/assets/publications/paper_6_RoMaP.pdf
- [4] Y. Song, Y. Kang, and S. Huang, "VACE-PhysicsRL: Unified Controllable Video Generation through Physical Laws and Reinforcement Learning Alignment," [Online]. Available: https://nsh423.github.io/assets/publications/paper_5_VACE.pdf
- [5] S. Sajjan et al., "ClearGrasp: 3D Shape Estimation of Transparent Objects for Manipulation," *ICRA*, 2020.
- [6] B. Mildenhall et al., "NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis," *ECCV*, 2020.