

GaussianGrasper: 3D Language Gaussian Splatting for Open-vocabulary Robotic Grasping

Yuhang Zheng, Xiangyu Chen, Yupeng Zheng, Songen Gu, Runyi Yang, Bu Jin, Pengfei Li, Chengliang Zhong, Zengmao Wang, Lina Liu, Chao Yang, Dawei Wang, Zhen Chen, Xiaoxiao Long*, Meiqing Wang*

Abstract—Constructing a 3D scene capable of accommodating open-ended language queries, is a pivotal pursuit, particularly within the domain of robotics. Such technology facilitates robots in executing object manipulations based on human language directives. To tackle this challenge, some research efforts have been dedicated to the development of language-embedded implicit fields. However, implicit fields (e.g. NeRF) encounter limitations due to the necessity of processing a large number of input views for reconstruction, coupled with their inherent inefficiencies in inference. Thus, we present the *GaussianGrasper*, which utilizes 3D Gaussian Splatting to explicitly represent the scene as a collection of Gaussian primitives. Our approach takes a limited set of RGB-D views and employs a tile-based splatting technique to create a feature field. In particular, we propose an Efficient Feature Distillation (EFD) module that employs contrastive learning to efficiently and accurately distill language embeddings derived from foundational models. With the reconstructed geometry of the Gaussian field, our method enables the pre-trained grasping model to generate collision-free grasp pose candidates. Furthermore, we propose a normal-guided grasp module to select the best grasp pose. Through comprehensive real-world experiments, we demonstrate that *GaussianGrasper* enables robots to accurately query and grasp objects with language instructions, providing a new solution for language-guided manipulation tasks. Data and codes can be available at <https://github.com/MrSecant/GaussianGrasper>.

Index Terms—Language-guided Robotic Manipulation, 3D Gaussian Splatting, Language Feature Field

I. INTRODUCTION

Recently, there has been an increasing scholarly focus on language-guided robotic manipulation due to its vast potential in facilitating human-robot interaction, enabling robotic home services, and enhancing flexible manufacturing. Imagine that a robot is asked to pick up the water cup in a cluttered, unstructured environment, it needs to (1) locate the water cup via responding to language description; (2) be aware of the geometry to execute a stable grasp. In this process, understanding the diverse objects with different shapes and material properties in the open world is the pivotal challenge.

Although many works have proposed various solutions, existing capabilities of scene understanding are insufficient

Yuhang Zheng is with the School of Mechanical Engineering and Automation, Beihang University and EncoSmart. Xiangyu Chen and Zhen Chen are with the EncoSmart. Yupeng Zheng and Bu Jin are with the Institute of Automation, Chinese Academy of Sciences. Songen Gu, Pengfei Li and Chengliang Zhong are with the AIR, Tsinghua University. Runyi Yang is with the Imperial College London. Lina Liu is with the China Mobile Research Institute. Zengmao Wang is with the School of Computer Science, Wuhan University. Chao Yang is with the Shanghai AI Laboratory. Dawei Wang and Xiaoxiao Long are with the Department of Computer Science, the University of Hong Kong. Meiqing Wang is with the School of Mechanical Engineering and Automation, Beihang University.

* Xiaoxiao Long and Meiqing Wang are the corresponding authors. Email: xxlong@connect.hku.hk, wangmq@buaa.edu.cn

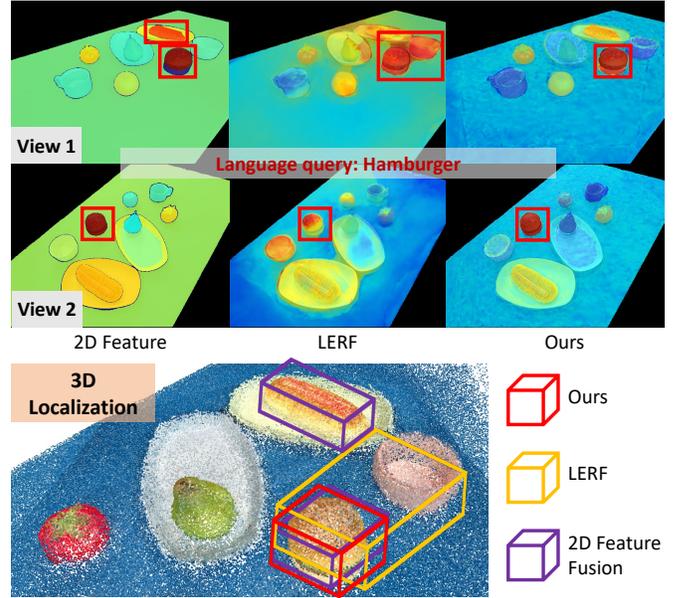


Fig. 1. We present a comparison between our method, 2D feature fusion, and LERF. When given the language query "hamburger", the features extracted by the 2D foundation models exhibit inconsistencies between two viewpoints, and LERF lacks clear segmentation boundaries. Consequently, they both suffer from imprecise 3D localization, as depicted by the yellow and purple 3D bounding boxes. In contrast, our method reconstructs a consistent feature field and achieves more precise 3D localization.

to afford language-guided manipulation. Most existing works are based on 2D images [1], [2], [3], [4] which are efficient but have limitations for robotic manipulation as robots can not easily infer visual occlusion and spatial relation from multi-view misaligned images.

To obtain precise 3D positions for robotic manipulation, recent works have focused on 3D representations. A straightforward approach is leveraging 2D visual models [5], [6], [7] to extract semantics and then fuse the 2D semantics into 3D points or volume. However, this fusion strategy suffers from semantic inconsistency in 3D, as the semantics provided by the visual model are not consistent across multi-views. Other methods [8], [9], [10], [11], [12], [13] that use 3D backbone to extract features and are supervised by 3D annotation or manipulation feedback can effectively make robots explicitly understand 3D scenes and learn skills but are difficult to apply to the real world due to the challenges in data acquisition and annotation. Recently, distilled feature fields (DFFs) [14], [15], [16] which reconstruct 3D feature fields from 2D images via implicit representation were introduced. Based on 2D-to-3D distillation, recent works [17], [18] have made impressive

progress in improving 3D scene understanding and enabling robots to interact with the physical world according to natural language. However, DFFs can not be inflexibly applied for robotics manipulation as most of these methods suffer from (1) imprecise localization as these methods extract patch-level features, resulting in ambiguous boundaries; (2) high costs of collecting dense training views (e.g. 50 views in F3RM [17]); (3) slow inference speed, hindering robots from responding to language instructions in time and (4) weak ability to cope with scene changes caused by manipulation.

To tackle problems, we introduce **GaussianGrasper**, an open-world robotic manipulation system based on 3D Gaussian Splatting (3DGS) [19], which models the 3D scene as a set of 3D Gaussian primitives. Our insight is that we (1) reconstruct a 3D feature field via efficient feature distillation to support language-guided localization; (2) render depth and surface normal to provide detailed local geometry, enabling the generation of feasible grasping poses; (3) operate Gaussian primitives and fine-tune 3DGS to update the changed scene.

More specifically, our method enables language-guided manipulation via the following steps: (1) Initialization: we scan RGB-D images of a few viewpoints to initialize the 3DGS, reducing the cost of data collection. (2) Feature field reconstruction: we propose an efficient feature distillation (EFD) module that employs SAM [5] and CLIP [6] to extract dense and shape-aware 2D descriptors and leverage contrastive learning to efficiently optimize the distilled features. (3) Localization and grasp: we use open-vocabulary queries to locate the target object and use a pre-trained grasping module to provide grasp poses where rendered normal is used to filter out unfeasible proposals based on *Force-closure* theory. (4) Scene updating: After executing manipulation, we update the scene by operating corresponding Gaussian primitives and fine-tuning 3DGS with images from fewer views.

In summary, the contributions of this paper are as follows:

- We introduce **GaussianGrasper**, a robot manipulation system implemented by a 3D Gaussian field endowed with open-vocabulary semantics and accurate geometry that is capable of rapid updates to support open-world manipulation tasks guided by language.
- We propose EFD that leverages contrastive learning to efficiently distill CLIP features and augment feature fields with SAM segmentation prior, addressing computational expense and boundary ambiguity challenges.
- We propose a normal-guided grasp module that uses rendered normal to filter out unfeasible grasp poses.
- We demonstrate the system’s zero-shot generalization capability for manipulation tasks in multiple real-world household tabletop scenes and common objects.

II. RELATED WORK

A. Grasp Pose Detection

Grasp pose detection is the pivotal part of robot grasping, which plays a critical role in enabling the robot to interact with objects in the physical world. Previous 3-DoF Grasping

methods treated the grasping task as 2D pose detection [20], [21], [22], [23], [24]. They define the grasp pose as a fixed-height oriented rectangle and predict the orientation and the width of the rectangle. However, their predicted grasp poses are limited to 3-DoF due to the lack of 3D geometry. To allow robots to plan higher dexterous grasps, extensive works focuses on 6-DoF grasping which use depth to augment the grasp pose detection [25], [26] or leverage point cloud as input to provide local geometry [27], [28], [29], [30], [31]. These methods exhibit high success rate when the depth is accurate but suffering from performance drops when encountering photometrically challenging objects such as transparent objects. To further improve the performance, some work fuses RGB with depth [32], [33], [34] as RGB can alleviate the depth-missing problem.

In this paper, we use the RGB-D based method to generate grasp poses. To achieve stable grasp, we further explicitly utilize *Force-closure* theory [35] to enhance grasp pose detection where estimated normal is used to filter out unfeasible grasp poses (more details in *Normal-guided grasp* of Sec. III-C3).

B. Reconstructing 3D Feature Field For Manipulation

A number of recent work integrate 2D foundation models with 3D feature fields in contexts other than robotic manipulation [36], [14], [37], [38], [16], [39]. Based on the implicit representation, methods [14], [39], [16] leverage feature distillation via neural rendering to reconstruct 3D feature field. While as explicit representations, methods [36], [37], [38] re-project 2D features and optimize the feature field in 3D.

For robot manipulation, recent works like F3RM [17] and LERF-TOGO [18] leverage feature distillation to improve the robot understanding of 3D scenes and enable language-guided grasping. However, these methods need images from many viewpoints as input and can not quickly update the scene to cope with changes (movement or rotation) of objects. As an explicit representation, SparseDFF [40] re-projects 2D features to 3D and leverages fusion strategy to optimize the feature field, leading to a significant reduction in usage of the number of viewpoints. However, explicit representation is not efficient for carrying high-dimension language features due to memory usage and computation time.

Different from the methods above, we propose an efficient feature distillation module based on explicit 3D Gaussians representation which uses segmentation priors provided by SAM to speed up feature field reconstruction and reduce the usage of memory (more details in Sec. III-B).

III. METHODOLOGY

In general terms, our method aims to pick up objects or place objects in specified locations according to the language instructions. The pipeline is shown in Fig. 2 (a) where our method (1) collects multi-view RGB-D images as input to initialize 3D Gaussian field; (2) reconstructs 3D feature field via efficient feature distillation module and (3) achieves language-guided manipulation. Specifically, we first introduce how to initialize 3DGS and the differentiable rasterizer of 3DGS in

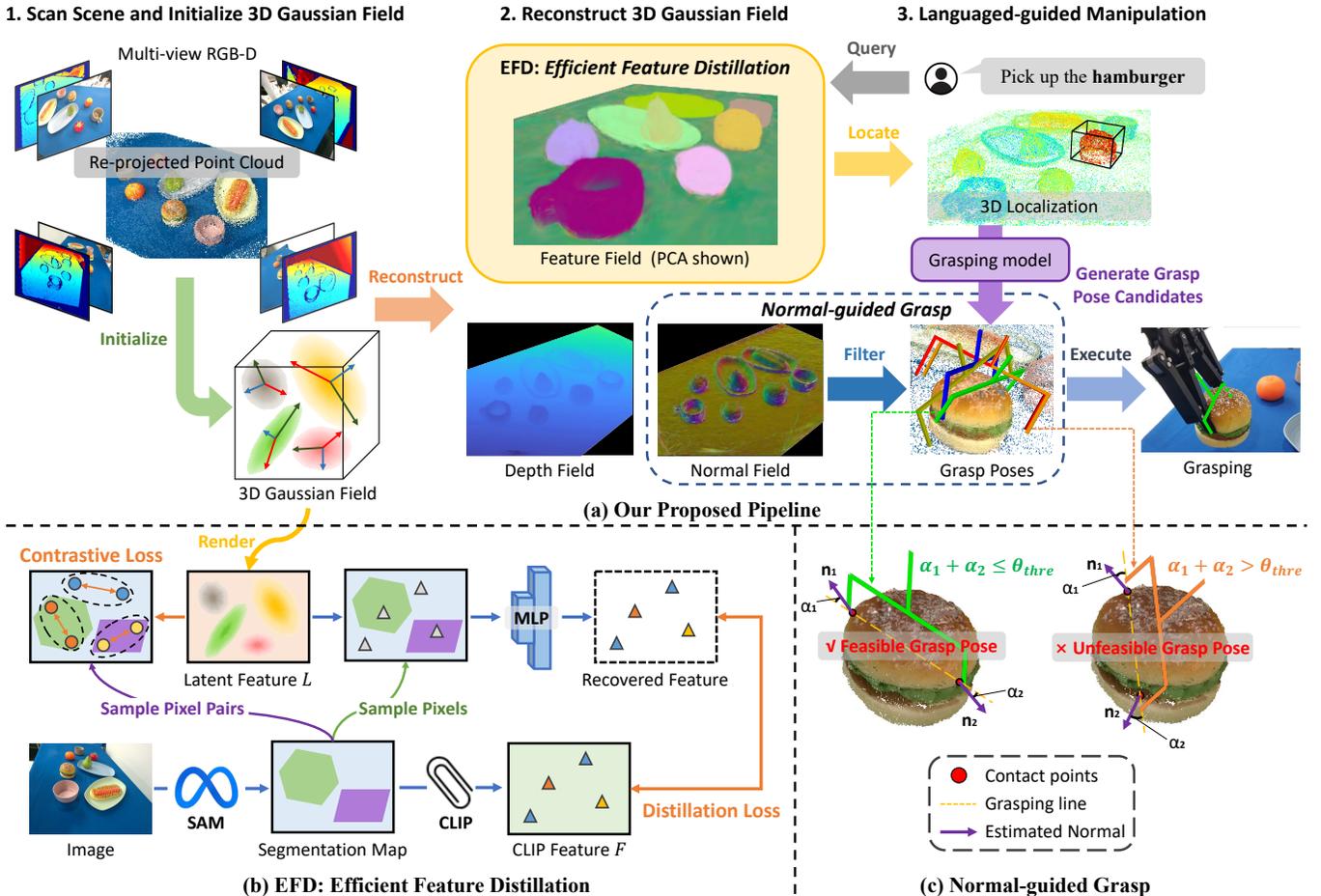


Fig. 2. The architecture of our proposed method. (a) is our proposed pipeline where we scan multi-view RGBD images for initialization and reconstruct 3D Gaussian field via feature distillation and geometry reconstruction. Subsequently, given a language instruction, we locate the target object via open-vocabulary querying. Grasp pose candidates for grasping the target object are then generated by a pre-trained grasping model. Finally, a normal-guided module that uses surface normal to filter out unfeasible candidates is proposed to select the best grasp pose. (b) elaborates on *EFD* where we leverage contrastive learning to constrain rendered latent feature L and only sample a few pixels to recover features to the CLIP space via an MLP. Then, the recovered features are used to calculate distillation loss with the CLIP features. (c) shows the normal-guided grasp that utilizes Force-closure theory to filter out unfeasible grasp poses.

section III-A. Next, we elaborate on the *EFD* module in section III-B. Finally, we introduce how to achieve language-guided manipulation in detail in section III-C, including locating objects through language queries, using a pre-trained grasping model to generate grasp poses, proposing a normal-guided grasp strategy to select feasible poses and updating the scene after manipulation.

A. Preliminaries: 3D Gaussian Splatting

1) *Gaussian Primitive Initialization*: 3DGS uses Structure from Motion (SfM) [41] as the initialization part, which takes a collection of RGB images as input and outputs a sparse point cloud. These points construct a set of Gaussian primitives, each defined by mean μ and a 3D covariance matrix $\Sigma = RSS^T R^T$, where R and S represent the rotation matrix and the scaling matrix. To reduce the viewpoints and speed up initialization, we employ multi-view depth to re-project pixels of RGB images into the world coordinate to initialize Gaussian field. Each point is the center of a Gaussian primitive and the rotation matrix and scaling matrix are initialized randomly.

2) *Differentiable Rasterizer for 3D Gaussians*: 3DGS renders the Gaussian primitives into images in a differentiable manner to optimize the parameters. Given a set of 3D Gaussian primitives $\mathcal{G} = \{g_i \mid i = 1, 2, 3, \dots, n\}$, each 3D Gaussian primitive g_i is first projected onto the corresponding 2D plane and is then sorted based on its depth d_i from the viewpoint plane. By tile-based rasterization, we sum up the pixel color $C(u)$ after sorting each primitive:

$$C(u) = \sum_i c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (1)$$

where c_i is a feature vector represented by spherical harmonics (SH) and α_i is obtained by multiplying Gaussian weight with opacity α associated to Gaussian primitives.

B. Efficient Feature Distillation

We reconstruct 3D open-vocabulary feature field via extracting dense CLIP features and efficiently distilling them into 3D,

based on 3DGS. The open-vocabulary reconstruction enables the scene to respond to natural language instructions.

1) *Instance-level Segmentation Prior and Open-vocabulary Features*: To extract dense and shape-aware open-vocabulary features, we first use SAM to produce a set of instance-level masks. Then, we leverage CLIP to obtain open-vocabulary features for each mask. Concretely, we process the input images through SAM to obtain a set of mask proposals and corresponding scores. Based on these scores, a non-maximum suppression strategy [42] is then implemented to filter superfluous masks. The resultant filtered set of masks constitutes a segmentation map of the image, representing instance-level priors. After filtration, we process each valid mask-aligned image region into the CLIP model to extract open-vocabulary features. Finally, we incorporate CLIP features of all masks into a feature map, referred to as F .

2) *Open-vocabulary Feature Distillation*: We propose a novel and efficient open-vocabulary feature distillation method that starts with enhancing each 3D Gaussian primitive with embedded open-vocabulary feature. As an explicit representation, a 3D Gaussian field can be composed of millions of primitives. Directly incorporating high-dimensional CLIP features (over 500 dimensions) into all primitives will result in unacceptable memory costs and computation time. Therefore, we compress the embedded open-vocabulary feature of Gaussian primitives from high-dimension CLIP space to low-dimension latent space. After initializing 3D Gaussian primitives with low-dimension latent feature l , we employ a feature rasterizer to render the feature map L :

$$L(u) = \sum_i l_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (2)$$

where l_i is the open-vocabulary feature embedding of i^{th} 3D gaussian primitives and $L(u)$ represents the rendered open-vocabulary feature embedding at pixel u .

To distill the 2D open-vocabulary feature to 3D field, we need to (1) recover the dimension of L to that of F and (2) minimize the feature distance between the recovered L and the F . However, high-dimensional vector computation for dense feature maps leads to a catastrophic increase in computation time and memory usage. To tackle this problem, we propose a contrastive-learning-based distillation strategy, as shown in Fig. 2 (b). Specifically, having instance-level segmentation masks extracted by SAM, we impose constraints for the consistency of each pixel’s rendered feature within the same mask. To ensure runtime efficiency, we randomly sample a few pixel pairs within each mask and only minimize the distance of latent features between pixels of per sampled pair. The number of pairs for each mask is proportion to the mask’s area and the number of total pairs n is fixed. We calculate the average feature distances between all pairs as contrastive loss, written as:

$$\mathcal{L}_{contr.} = 1 - \frac{1}{n} \sum_{i=1}^n L(u_i) \cdot L(v_i) \quad (3)$$

where u_i and v_i are pixels in the i^{th} sampled pair. As the contrastive loss homogenizes the features within each mask, we only need to recover latent features of per mask to the CLIP space and subsequently minimize the distance between the recovered feature and the CLIP feature. In practice, we randomly sample the same number of pixels within each mask, whose latent features are then recovered via a trainable decoder Ψ composed of two fully connected layers. We calculate the distillation loss between the recovered feature and the CLIP feature of all sampled k pixels, defined as:

$$\mathcal{L}_{distill} = 1 - \frac{1}{k} \sum_{i=1}^k \Psi(L(i)) \cdot F(i) \quad (4)$$

By enhancing the low-dimension open-vocabulary feature embedding of 3D Gaussian and using contrastive learning which leverages the segmentation prior derived from SAM, our method provides a powerful and efficient solution for reconstructing 3D open-vocabulary representation.

C. Language-guided Robotic Manipulation

We use the reconstructed feature field to conduct robotic manipulation. Given a language instruction, our method begins with employing open-vocabulary queries to locate the target object. Subsequently, we render the depth and normals of 3D Gaussian primitives to obtain the object’s detailed geometry. Then, a point cloud-based grasping module is used to generate grasp poses and the rendered normal is used to filter out unfeasible ones. After manipulating objects, we quickly update the scene using observations from fewer viewpoints.

1) *Open-vocabulary Querying*: As our reconstructed feature field is aligned with natural language, we can locate the object described by language instructions via open-vocabulary querying. We first follow the approach of LERF [16] to compute the relevance score s for each textual query:

$$s = \min_i \frac{\exp(\Psi(L) \cdot T^q)}{\exp(\Psi(L) \cdot T^q) + \exp(\Psi(L) \cdot T_i^{canon})} \quad (5)$$

where T^q is the CLIP embedding of the text query and T_i^{canon} is a set of canonical phrase embeddings T_i^{canon} selected from "object", "things", "stuff", and "texture".

As a result, for each textual query, we obtain a relevance heatmap where the points with relevance scores below a pre-determined threshold will be filtered out. Thus, the remaining region forms a mask for predicting the queried object. After obtaining the mask of the queried object, we locate the object with a bounding box and convex hull.

2) *Geometry Restriction*: To obtain dense point cloud representations of objects and surface normal which is closely related to robotic grasping, we render depth and surface normal to multiple viewpoints.

Depth rendering: Similar to the rendering of RGB, we compute the depth value for each pixel using the rasterizer:

$$D(u) = \sum_i d_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (6)$$

where $D(u)$ indicates the rendered depth map at pixel u .

We use the depth map obtained from the depth camera for the corresponding viewpoint to supervise the rendered depth map where we calculate the L1 loss:

$$\mathcal{L}_{depth} = \frac{1}{m} \sum_{i=1}^m |\hat{D}(i) - D(i)| \quad (7)$$

where \hat{D} is the depth map obtained from the depth camera and m is the number of pixels with valid depth value.

Normal rendering: As surface normals are directional vectors that should exhibit rotational equivariance, normals cannot be rendered as semantic features. Therefore, we follow [43], [44] that use the shortest axis direction of the 3D Gaussian primitives to serve as surface normal. As a result, the normals are geometric properties of Gaussian primitives, related to their orientations. We render the normal map by using the rasterizer:

$$\mathbf{N}(u) = \sum_i \mathbf{n}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (8)$$

where $\mathbf{N}(u)$ indicates the rendered surface normal map at pixel u and \mathbf{n}_i is the normal of the i^{th} 3D Gaussian primitive.

The rendered normal map represents the per-pixel surface normal in the robot base coordinate. We use a Sobel-like operator to compute the normals for each pixel in the acquired depth map and transform the calculated normals from the camera coordinate to the robot base coordinate. After normalizing all normal maps to unit vectors, we supervise the rendered normals within the valid depth region:

$$\mathcal{L}_{normal} = \frac{1}{m} \sum_{i=1}^m \left(\left(\hat{\mathbf{N}}(i) - \mathbf{N}(i) \right)^2 + 1 - \hat{\mathbf{N}}(i) \cdot \mathbf{N}(i) \right) \quad (9)$$

3) *Feasible Grasp Pose Generation:* After obtaining the localization and geometry of the object, we (1) employ a grasp detection method to propose initial grasp poses and (2) utilize our rendered normal to augment the grasp pose proposals according to the force closure theory.

Grasp pose generation: In our work, we employ the AnyGrasp [34], the state-of-the-art grasp detection network, which takes colorful point cloud as input and generates a set of collision-free grasp proposals for parallel two-finger grippers. Each grasp proposal is represented by grasp position, width, height, depth and a grasping score. We generate the grasp-pose proposals for the object queried by language with the following steps: (1) Re-project the rendered depth of the queried object from each viewpoint to produce a dense point cloud representation. (2) Combine the object's point cloud with the point cloud derived from 3D Gaussian primitives because the derived point cloud depicts the localization and approximate shape of other objects, which is beneficial for generating collision-free grasp poses. (3) Employ AnyGrasp to generate grasp poses, which are restricted in the aforementioned bounding box, obtained from steps in III-C1.

Normal-guided grasp: Although AnyGrasp provides a "grasping score" for each grasp proposal, the grasp pose with the

highest score may not be the suitable one for stable grasp. Thus, we explicitly utilize the force closure theory as a screening mechanism. For each grasp pose proposal with two contact points, we calculate the angle between the grasping line and the surface normal of each contact point and get the sum of two angles. If the sum of the two angles is less than or equal to a pre-defined threshold θ_{thre} , we regard the grasp pose as a feasible proposal, otherwise it is unfeasible, as shown in Fig. 2 (c). We choose the pose with the highest "grasping score" in feasible proposals as the final grasp pose.

4) *Gaussian Field Updating:* As an explicit representation composed of 3D Gaussian primitives, it is easy to operate Gaussian primitives. Therefore, after manipulating an object, we will operate all Gaussian primitives within this object's convex hull with the same rotation and translation, which can be calculated from the transform of the end-effector of the robot arm. After operating 3D Gaussian primitives to a new position, we use fewer views and time to fine-tune 3D Gaussian field to update the scene.

IV. EXPERIMENT

In this section, we first introduce the setup of the experimental environment. Next, we conduct experiments to validate our proposed EFD module where we report both quantitative results and qualitative results. Subsequently, we show the results of geometry reconstruction and conduct ablation study to demonstrate the effectiveness of our proposed normal-guided grasp. Finally, we conduct grasp-update-grasp experiments on scenes to verify the effectiveness and efficiency of our scene update module. These experiments fully demonstrate the performance of our method in open-scene understanding and language-guided grasping.

A. Experimental Setup

1) *Scenes, Objects and Devices:* We built a $140 \times 70 \times 30 \text{cm}^3$ desktop scene with common objects in the kitchen including various food and tableware as well as office supplies including staplers, mouse and decorative ornaments. We use a UR5 robot arm equipped with a ROBOTIQ gripper to execute robotic manipulation. We set up our system in 10 open desktop scenes with a total of 44 objects (40 are graspable) where we execute language-guided manipulation 120 times. In terms of computing resources, we use an NVIDIA RTX-3090 GPU to reconstruct the feature field and reconstruct geometry. The reconstruction process only requires approximately 6GB of memory in total.

2) *Data Collection and Processing:* We first use the robot arm equipped with a Realsense D455 to scan the desktop scene from 16 viewpoints. At the same time, we will also record the camera extrinsic parameter, calculated through the transformation of the end effector to the base. After obtaining images and depth maps, we re-project the depth map and convert all re-projected points to the base coordinate. We downsample the number of these points to about 300k to initialize the 3D Gaussian primitives. The collected images are processed by SAM and CLIP to generate segmentation maps and open-vocabulary feature maps.

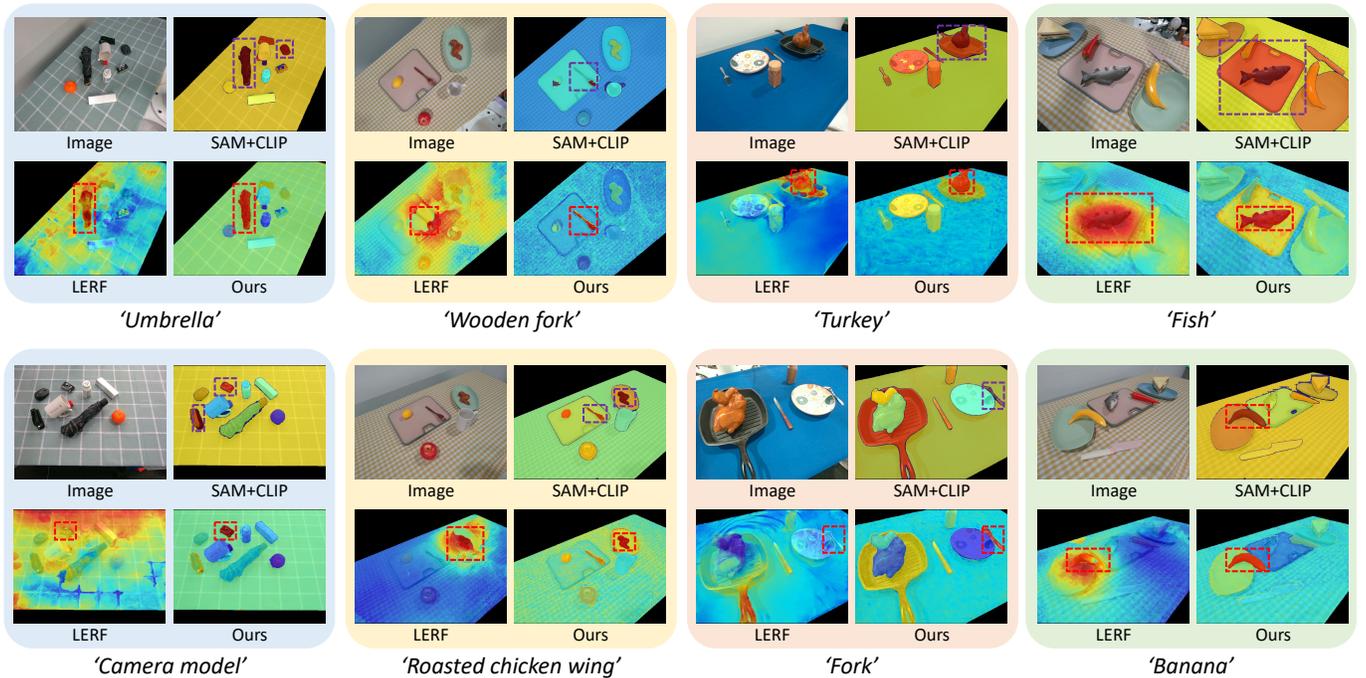


Fig. 3. Relevance map of the given language instructions. Our method exhibits clearer segmentation boundaries compared to LERF, which can be used to obtain more accurate localization. Compared with SAM + CLIP, our approach exhibits more consistent open-vocabulary features across multi-views. For instance, in *'Roasted chicken wing'*, the response of SAM + CLIP is the chicken wing and the fork while our method makes the correct response.

B. Results of Efficient Feature Distillation

We show both qualitative results and quantitative results to demonstrate the effectiveness and efficiency of our proposed EFD module. Our baselines are Lseg [45] and LERF [16] (All mention of LERF in our experiments includes an **extra depth supervision** to ensure a fair comparison with our method.)

In qualitative results, we compare our method with SAM + CLIP (our 2D sudo labels) and LERF and show the relevance map of each given language instruction. As shown in Fig. 3, the purple boxes demonstrate that the features extracted from SAM and CLIP are not accurate such as the incomplete wooden fork in *'Wooden fork'* and the pot and turkey with similar semantic features in *'Turkey'*. In comparison, our relevance map is more accurate, proving that our reconstructed feature field solves the problem of feature inconsistency across multiple views. Besides, compared with LERF, our method exhibits better segmentation boundaries. Therefore, our method can help robots reduce the ambiguity of object perception.

We report the quantitative results of two tasks including segmentation and localization. In the segmentation task, as described in III-C1 we filter out the region whose relevance score is below 0.85 to form a predicted segmentation map. We calculate the mIoU metric between predicted segmentation maps and our manually annotated ground truth. In the localization task, following LERF, given a language instruction, if the point with the highest relevance score is in the target object, it is a successful localization. We calculate the average accuracy of all responses as the metrics. The results of segmentation and localization are shown in Table I where our method significantly outperforms other approaches.

TABLE I
QUANTITATIVE COMPARISONS OF SEMANTIC SEGMENTATION AND LOCALIZATION ACCURACY ON OUR SCENARIOS

Method	mIoU(%) \uparrow	Accuracy(%) \uparrow	Time per query (s) \downarrow
Lseg[45]	26.4	40.6	-
LERF*[16]	41.3	65.1	40.27
Ours	58.2	87.5	0.22

"*" represents LERF with an extra depth supervision

Besides, we test the query speed and report results in Table I. The metric is the time usage (s) per text query at a resolution of 640×480 . It can be seen that our method achieves an approximate 180 \times speedup over LERF. We also directly distill CLIP features into 3D Gaussian field, which takes over 70 GB of memory, making it hard to be applied to robots.

C. Results of Geometry Reconstruction

We show the visualization of our rendered depth and normal compared with ground truth (scanned by D455 camera), as shown in Fig. 4. The black region represents the invalid value of ground truth. It can be seen that the surface normal we rendered is smoother than the ground truth, as shown in red boxes. Furthermore, even in areas where the ground truth is invalid, we can still render accurate depth and surface normal. For example, in the third row, the camera can not capture the depth of the silver eyeglass case in this view because of the specular reflection of the case but our method renders accurate depth and normal of it.

D. Effectiveness of Normal-guided Grasp

In this subsection, we aim to validate the effectiveness of our proposed normal-guided grasp. We first give the qualitative

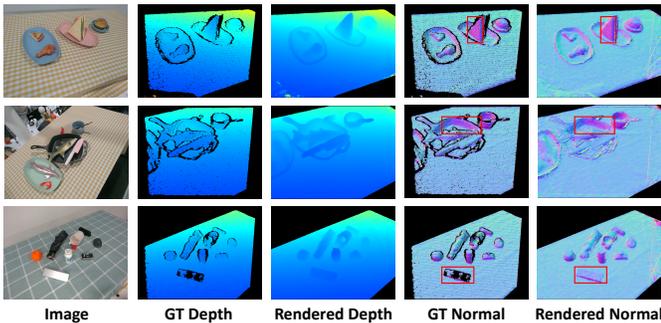


Fig. 4. Compared with scanned depth and surface normal, our rendered depth and surface normal is smoother. Our method renders accurate depth and surface normal even in areas where the ground truth is invalid.

TABLE II

GRASPING RESULTS: RESULTS ARE REPORTED ACROSS 40 DIFFERENT OBJECTS OF 10 SCENARIOS. EACH OBJECT IS GRASPED THREE TIMES.

Method	Grasping Success Rate (%)
Lseg + Depth[45]	26.7
LERF + AnyGrasp[16]	55.8
Ours w/o. Normal Filter	78.3
Ours w/ Normal Filter	85.0

result to validate that the surface normal can filter out unfeasible grasp poses. As shown in Fig. 5, the original top-ranked proposal (red) is filtered out as the angles between its grasping line and surface normal of contact points are too large. In contrast, the original second-ranked proposal is feasible. Thus, we execute a grasp based on this pose.

Besides, we report the quantitative results of the grasping success rate with and without the normal filter, as shown in Table II. Leveraging the normal filter significantly increases the success rate by 7.7%, further demonstrating the effectiveness of our proposed normal-guided grasp.

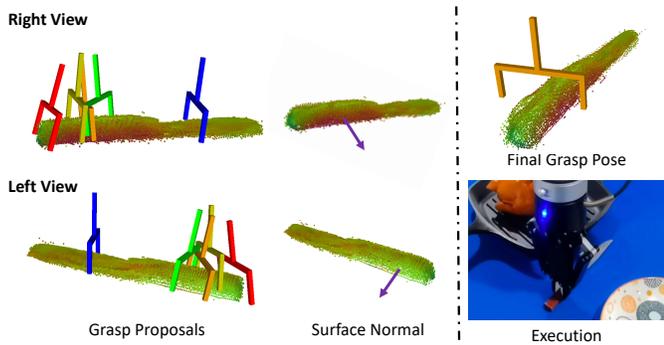


Fig. 5. Effectiveness of our proposed normal-guided grasp. The left column shows the top 5 grasp proposals provided by AnyGrasp. The redder the color, the higher the grasping score. The middle column displays the surface normal of the object, with purple arrows indicating the normal of the contact points. The right column demonstrates the successful execution of grasping the knife utilizing the final grasp pose after filtering out unreasonable proposals.

E. Results of Language-guided Manipulation

In this subsection, we first give the result of language-guided grasping. Then, we validate our proposed scene updating via continuous language-guided picking and placing.

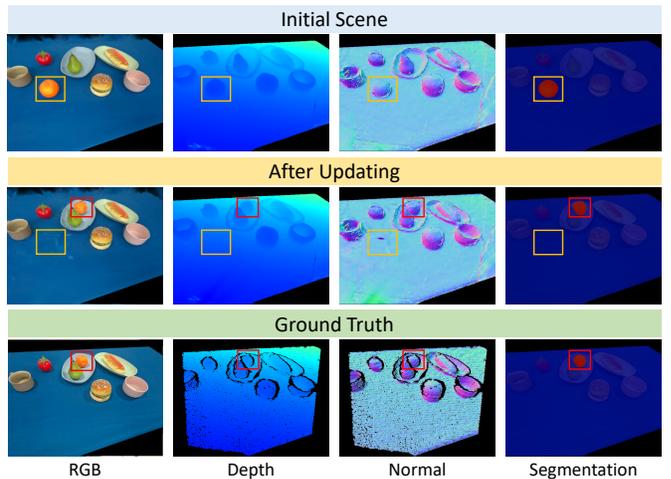


Fig. 6. Results of scene update. We show the RGB, depth, normal, and segmentation before and after the scene update based on the language query "orange". As indicated by the yellow boxes, our scene updating successfully moves the orange to the plate and restores the region that was previously obscured by the orange. As indicated by the red boxes, the updated orange still maintains accurate geometry and semantic features.

TABLE III

EFFICIENCY COMPARISON BETWEEN LERF AND OUR METHOD.

Method	Viewpoints	Memory	Time
LERF [16]	16	15GB	30min
Ours	5	4GB	1min

1) *Successful rate of manipulation*: In this subsection, we show the result of language-guided grasping where we tested 120 times on 40 objects. We compare our method with (1) Lseg + AnyGrasp and (2) LERF + AnyGrasp. To obtain the 3D point cloud, we use the rendered depth to re-project the segmentation masks of LERF and use scanned depth to re-project the segmentation masks of LSeg. We define a successful grasp as stably picking up the corresponding object according to the language instruction and raising it to a height of more than 10cm over 3 seconds. The result is shown in Table II, where our method far exceeds other methods in success rate.

2) *Scene updating*: To validate the effectiveness of our proposed efficient scene updating, we execute an experiment whose process is (1) picking up the object and placing it to the target position according to the language instruction, (2) capturing RGB-D images from 5 viewpoints to update the scene, (3) executing another manipulation on this object. The result is shown in Fig. 6, where our updated scene remains high quality RGB, geometry and semantic features, proving the effectiveness of the proposed scene updating. Besides, we also show the efficiency comparison between LERF and ours including viewpoint numbers, memory usage and reconstruction time for updating, as shown in Table IV-E2. Our scene update capability makes the reconstructed scene more capable of handling continuous grasping.

V. LIMITATION

One limitation is that our reconstructed scene remains static. Although we have proposed a scene-updating module to handle

continuous robotic grasping, we are unable to account for unrecordable scene changes, such as objects shifting positions caused by collisions or vibrations. Another limitation is that our method fails to estimate the depth and normal of transparent objects due to the lack of ground truth.

VI. CONCLUSION

In this paper, we introduce GaussianGrasper, a novel approach for open-world robotic grasping guided by natural language instructions from RGB-D inputs. Taking multi-view RGB-D images as input, our method efficiently reconstructs consistent feature fields through our proposed EFD module. The feature field enables the robot to understand the open world and make precise localization based on language instructions. Besides, we estimate the geometry and propose the normal-guided grasp to augment the robotic grasping. Furthermore, our scene can also be quickly updated to support continuous grasping. Sufficient real-world experiments have demonstrated the effectiveness of our method.

REFERENCES

- [1] C. Lynch and P. Sermanet, "Language conditioned imitation learning over unstructured data," *RSS*, 2020.
- [2] M. Shridhar, L. Manuelli, and D. Fox, "Cliport: What and where pathways for robotic manipulation," in *CoRL*, 2022.
- [3] X. Lin, J. So, S. Mahalingam, F. Liu, and P. Abbeel, "Spawnet: Learning generalizable visuomotor skills from pre-trained networks," *arXiv preprint arXiv:2307.03567*, 2023.
- [4] P.-L. Guhur, S. Chen, R. G. Pinel, M. Tapaswi, I. Laptev, and C. Schmid, "Instruction-driven history-aware policies for robotic manipulations," in *CoRL*, 2022.
- [5] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, *et al.*, "Segment anything," in *ICCV*, 2023.
- [6] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.
- [7] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *ICCV*, 2021.
- [8] D. Z. Chen, A. X. Chang, and M. Nießner, "Scanrefer: 3d object localization in rgb-d scans using natural language," in *ECCV*, 2020.
- [9] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu, "Synergies between affordance and geometry: 6-dof grasp detection via implicit representations," *RSS*, 2021.
- [10] S. Chen, R. G. Pinel, C. Schmid, and I. Laptev, "Polarnet: 3d point clouds for language-guided robotic manipulation," *ArXiv*, 2023.
- [11] M. Shridhar, L. Manuelli, and D. Fox, "Perceiver-actor: A multi-task transformer for robotic manipulation," in *CoRL*, 2022.
- [12] Y. Ze, G. Yan, Y.-H. Wu, A. Macaluso, Y. Ge, J. Ye, N. Hansen, L. E. Li, and X. Wang, "Gnfactor: Multi-task real robot learning with generalizable neural feature fields," in *CoRL*, 2023.
- [13] C. Zhong, Y. Zheng, Y. Zheng, H. Zhao, L. Yi, X. Mu, L. Wang, P. Li, G. Zhou, C. Yang, X. Zhang, and J. Zhao, "3d implicit transporter for temporally consistent keypoint discovery," in *ICCV*, 2023.
- [14] S. Kobayashi, E. Matsumoto, and V. Sitzmann, "Decomposing nerf for editing via feature field distillation," *NeurIPS*, 2022.
- [15] V. Tschernezki, I. Laina, D. Larlus, and A. Vedaldi, "Neural feature fusion fields: 3d distillation of self-supervised 2d image representations," in *3DV*, 2022.
- [16] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik, "Lerf: Language embedded radiance fields," in *ICCV*, 2023.
- [17] W. Shen, G. Yang, A. Yu, J. Wong, L. P. Kaelbling, and P. Isola, "Distilled feature fields enable few-shot language-guided manipulation," in *CoRL*, 2023.
- [18] A. Rashid, S. Sharma, C. M. Kim, J. Kerr, L. Y. Chen, A. Kanazawa, and K. Goldberg, "Language embedded radiance fields for zero-shot task-oriented grasping," in *CoRL*, 2023.
- [19] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, 2023.
- [20] J. Redmon and A. Angelova, "Real-time grasp detection using convolutional neural networks," in *ICRA*, 2015.
- [21] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. A. Ojea, and K. Goldberg, "Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics," *RSS*, 2017.
- [22] D. Guo, F. Sun, H. Liu, T. Kong, B. Fang, and N. Xi, "A hybrid deep architecture for robotic grasp detection," in *ICRA*, 2017.
- [23] X. Zhou, X. Lan, H. Zhang, Z. Tian, Y. Zhang, and N. Zheng, "Fully convolutional grasp detection network with oriented anchor box," in *IROS*, 2018.
- [24] Y. Zheng, Q. Wang, C. Zhong, H. Liang, Z. Han, and Y. Zheng, "Enhancing daily life through an interactive desktop robotics system," in *CICAI*, 2023.
- [25] X. Zhu, L. Sun, Y. Fan, and M. Tomizuka, "6-dof contrastive grasp proposal network," in *ICRA*, 2021.
- [26] G. Zhai, D. Huang, S.-C. Wu, H. Jung, Y. Di, F. Manhardt, F. Tombari, N. Navab, and B. Busam, "Monograspnet: 6-dof grasping with a single rgb image," in *ICRA*, 2023.
- [27] H. Liang, X. Ma, S. Li, M. Görner, S. Tang, B. Fang, F. Sun, and J. Zhang, "Pointnetgpd: Detecting grasp configurations from point sets," in *ICRA*, 2019.
- [28] C. Wu, J. Chen, Q. Cao, J. Zhang, Y. Tai, L. Sun, and K. Jia, "Grasp proposal networks: An end-to-end solution for visual learning of robotic grasps," *NeurIPS*, 2020.
- [29] M. Sundermeyer, A. Mousavian, R. Triebel, and D. Fox, "Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes," in *ICRA*, 2021.
- [30] B. Zhao, H. Zhang, X. Lan, H. Wang, Z. Tian, and N. Zheng, "Regnet: Region-based grasp network for end-to-end grasp detection in point clouds," in *ICRA*, 2021.
- [31] A. Alliegro, M. Rudorfer, F. Frattin, A. Leonardis, and T. Tommasi, "End-to-end learning to grasp via sampling from object point clouds," *RA-L*, 2022.
- [32] H.-S. Fang, C. Wang, M. Gou, and C. Lu, "Graspnet-1billion: A large-scale benchmark for general object grasping," in *CVPR*, 2020.
- [33] 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," *RA-L*, 2022.
- [34] H.-S. Fang, C. Wang, H. Fang, M. Gou, J. Liu, H. Yan, W. Liu, Y. Xie, and C. Lu, "Anygrasp: Robust and efficient grasp perception in spatial and temporal domains," *T-RO*, 2023.
- [35] A. Ten Pas and R. Platt, "Using geometry to detect grasp poses in 3d point clouds," *Robotics Research: Volume 1*, 2018.
- [36] H.-Y. F. Tung, R. Cheng, and K. Fragkiadaki, "Learning spatial common sense with geometry-aware recurrent networks," in *CVPR*, 2019.
- [37] H.-Y. F. Tung, Z. Xian, M. Prabhudesai, S. Lal, and K. Fragkiadaki, "3d-oes: Viewpoint-invariant object-factorized environment simulators," *arXiv preprint arXiv:2011.06464*, 2020.
- [38] S. Peng, K. Genova, C. Jiang, A. Tagliasacchi, M. Pollefeys, T. Funkhouser, *et al.*, "Openscene: 3d scene understanding with open vocabularies," in *CVPR*, 2023.
- [39] C. Wang, M. Chai, M. He, D. Chen, and J. Liao, "Clip-nerf: Text-and-image driven manipulation of neural radiance fields," in *CVPR*, 2022.
- [40] Q. Wang, H. Zhang, C. Deng, Y. You, H. Dong, Y. Zhu, and L. Guibas, "Sparsedff: Sparse-view feature distillation for one-shot dexterous manipulation," *arXiv preprint arXiv:2310.16838*, 2023.
- [41] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *CVPR*, 2016.
- [42] A. Neubeck and L. Van Gool, "Efficient non-maximum suppression," in *ICPR*, 2006.
- [43] Y. Jiang, J. Tu, Y. Liu, X. Gao, X. Long, W. Wang, and Y. Ma, "Gaussianshader: 3d gaussian splatting with shading functions for reflective surfaces," *arXiv preprint arXiv:2311.17977*, 2023.
- [44] X. Long, Y. Zheng, Y. Zheng, B. Tian, C. Lin, L. Liu, H. Zhao, G. Zhou, and W. Wang, "Adaptive surface normal constraint for geometric estimation from monocular images," *arXiv preprint arXiv:2402.05869*, 2024.
- [45] B. Li, K. Q. Weinberger, S. Belongie, V. Koltun, and R. Ranftl, "Language-driven semantic segmentation," in *ICLR*, 2022.