

VISO-Grasp: Vision-Language Informed Spatial Object-centric 6-DoF Active View Planning and Grasping in Clutter and Invisibility

Yitian Shi*, Di Wen*, Guanqi Chen, Edgar Welte, Sheng Liu, Kunyu Peng, Rainer Stiefelhagen, Rania Rayyes

Abstract— We propose VISO-Grasp, a novel vision-language-informed system designed to systematically address visibility constraints for grasping in severely occluded environments. By leveraging Foundation Models (FMs) for spatial reasoning and active view planning, our framework constructs and updates an instance-centric representation of spatial relationships, enhancing grasp success under challenging occlusions. Furthermore, this representation facilitates active Next-Best-View (NBV) planning and optimizes sequential grasping strategies when direct grasping is infeasible. Additionally, we introduce a multi-view uncertainty-driven grasp fusion mechanism that refines grasp confidence and directional uncertainty in real-time, ensuring robust and stable grasp execution. Extensive real-world experiments demonstrate that VISO-Grasp achieves a success rate of 87.5% in target-oriented grasping with the fewest grasp attempts outperforming baselines. To the best of our knowledge, VISO-Grasp is the first unified framework integrating FMs into target-aware active view planning and 6-DoF grasping in environments with severe occlusions and entire invisibility constraints.

I. INTRODUCTION

Robotic grasping in unstructured, cluttered environments remains a significant challenge, particularly for executing target-oriented grasps under partial or complete occlusions [1]. Humans instinctively overcome occlusion during targeted searches by adjusting their viewpoints and intuitive reasoning about spatial relationships to infer potential target locations. In comparison, leveraging the recent success of deep neural networks and contributions to large-scale training data [2]–[4], existing learning-based robotic grasp detectors [1] rely primarily on static viewpoint observations, assuming sufficient visibility of the target object from single or pre-defined multi-view perspectives [5]–[7]. This constraint reduces flexibility in scenarios where the target object is highly occluded or entirely unobservable.

Meanwhile, recent active grasping frameworks [8]–[10] have made initial attempts to integrate Next-Best-View (NBV) planning into grasping, which prioritize 3D reconstruction and rendering over target-driven view search for grasping, relying heavily on pre-defined search spaces. However, without such prior knowledge, these methods may fail to handle heavily cluttered scenarios and target invisibility. We argue that effectively reaching the target, even under

Karlsruhe Institute of Technology, Karlsruhe, Germany. Email: yitian.shi@kit.edu. This work is sponsored by the DFG SFB-1574 Circular Factory project and Baden-Württemberg Ministry of Science, Research and the Arts within InnovationCampus Future Mobility.

*Equal contribution

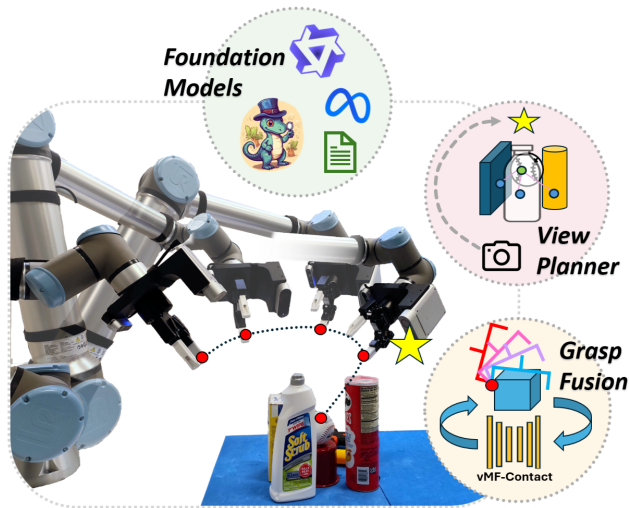


Fig. 1: VISO-Grasp, a unified system integrating Foundation Models (FMs) into target-aware active view planning and uncertainty-driven real-time 6-DoF grasp fusion.

invisibility, requires a structured active vision and sequential grasping strategy. This shall leverage object relationships to systematically optimize sensor viewpoints, remove obstructions when necessary, and incrementally reveal the target for efficient and successful grasp execution.

To this end, we propose VISO-Grasp, a Vision-language Informed Spatial Object-centric grasping framework that integrates off-the-shelf Foundation Models (FMs) [11]–[13] with active vision and occlusion-aware SE(3) grasp planning, which leverages the inherent prior of VLMs to achieve human-like decision-making. In contrast to recent works [14]–[16], which focus on planar 4-DoF grasp reasoning with VLMs while overlooking their fundamental limitations in accurate 3D spatial understanding and real-time execution, VISO-Grasp introduces evolving spatial reasoning that continuously integrates spatial and occlusion relationships. By incorporating an online grasp fusion mechanism, our approach dynamically refines target visibility and substantially improves grasp efficiency in heavy occlusions. Specifically, when direct grasping is infeasible, a sequential grasp strategy guided by VLMs is employed, leveraging the observable relational hierarchy to determine the optimal grasp sequence for decluttering and removing potential obstructions. In summary, our main contributions are threefold:

- **Vision-language-driven unified scene understanding for robotic perception**, which integrates scene reasoning from the VLM [11] to infer spatial relationships between objects. This facilitates occluder identification through Visual-In-Context Learning (VICL), and determines optimal sequential grasp orders in highly occluded environments.

- **Target-aware view planning for object-centric grasping**, a novel NBV planner that constructs object-centric continuous velocity fields using 3D instance segmentations from existing foundation models [12], [13] to optimize sensor viewpoints and maximize target visibility in cluttered environments.

- **Uncertainty-driven online grasp fusion**, a real-time grasp fusion framework that integrates multi-view grasp predictions through pose-centric categorization and real-time Bayesian updates. This leverages a pre-trained uncertainty-aware grasp generator [17], trained exclusively on simulated data, to enhance grasp robustness and reliability.

II. RELATED WORKS

A. Active View Search for Grasp Detection

Recent advancements in scene-level SE(3) grasp detection frameworks prioritize generalization to unseen, unstructured environments with diverse objects. While end-to-end approaches [6], [18]–[20] employ deep neural networks to infer grasp configurations from point clouds [18], [19] or volumetric representations [6], [20], the robustness diminishes in the presence of sim-to-real gap and heavy occlusions, where incomplete geometric information significantly impairs prediction reliability. To this end, multi-view strategies [5]–[7], [21] integrate observations across multiple viewpoints, while, either bringing the cost of target-agnostic view exploration or relying on pre-defined perspectives by sacrificing flexibility.

Within this context, active vision [22] facilitated by NBV strategies, emerges as a principled approach in vision-based decision-making such as navigation [23], [24] and 3D reconstruction [25], [26]. By incorporating historical observations, active vision employs a closed-loop optimization process to iteratively select viewpoints and maximize scene visibility for downstream tasks. Recent active grasping frameworks [8]–[10] primarily formulate the NBV problem as a correlated reconstruction and rendering task. For instance, they leverage strategies such as maximization of volumetric Information Gain (IG) derived from ray casting [8] and minimization of neural graspness or rendering inconsistency [9], [10]. These approaches are constrained by discrete viewpoint sampling, pre-defined search space prior (e.g., the bounding box around the targets), target-agnostic grasp planning, and the prioritization of reconstruction or rendering over grasping. Consequently, these limitations hinder adaptability for optimal target-oriented grasp and view planning in unstructured environments or even invisibility. In comparison, we consider involving spatial reasoning between objects to realize object-centric active view planning.

B. FMs in Robotic Perception and Planning

Foundation Models (FMs), including VLMs, have significantly improved robotic perception and planning by zero-shot object detection [27], [28], multi-modal reasoning [16], [29], [30], and decision-making [31], [32], reducing reliance on task-specific data collection while enhancing adaptability in real-world environments [33], [34]. Recent works leverage FMs to enable open-vocabulary perception, integrating zero-shot segmentation for novel object recognition [14], [15], while uncertainty-aware frameworks refine scene understanding by dynamically calibrating perception and decision uncertainties [35]. In robotic planning, multimodal models facilitate goal-conditioned decision-making by aligning spatial and semantic representations [36], while vision-language systems enhance strategic scene exploration and adaptive reasoning in unstructured scenarios [37]. Yet, most FM-driven approaches still rely on 2D representations, lacking explicit 3D spatial priors and depth reasoning [38]. Furthermore, they assume minimal occlusions and primarily operate under SE(2) planar constraints, with limited real-time viewpoint adaptation and active perception to dynamically refine scene understanding. To overcome these limitations, we propose the first FMs-driven framework that unifies multi-view 3D spatial reasoning and active viewpoint adaptation for occlusion-aware 6-DoF grasping, enabling dynamic scene refinement and grasp execution in cluttered environments.

III. PROBLEM FORMULATION

We consider a robotic manipulation system equipped with a wrist-mounted (eye-in-hand) RGB-D camera and a parallel-jaw gripper, operating in an unknown, cluttered environment containing a set of *objects* $\mathcal{O} = \{O_1, O_2, \dots, O_m\}$. Among these, a distinct *target object* O_* is subject to significant occlusions, often in *pile* or *packed* [6] configurations. The target object may be heavily or completely occluded due to *coverage* — where the target is fully enclosed by surrounding objects, or due to *blocking* — where objects positioned in front obscure direct access of the sensor’s line of sight.

To achieve target-oriented grasping in these scenarios, we incorporate FMs [11]–[13] for semantic scene understanding and adaptive view planning, allowing zero-shot generalization without task-specific training. In extreme cases such as complete occlusions or no valid grasp poses available on the target, the VLM either determines an optimal grasp execution sequence from current observed objects \mathcal{O}' to prioritize the removal of occluders, or infers the potential occluders by the VLM with VICL from historical observations.

IV. METHODS

A. System Overview

In general, VISO-Grasp integrates vision-based physical grasping and semantic scene understanding as core functional capabilities. Fig. 2 illustrates the fundamental components of our system, which comprises:

i) Adaptive Multi-view Open-Vocabulary 3D Object Detector (AMOV3D, Sec. IV-B), implemented via the VLM [11], refines 3D oriented bounding boxes (BBs)

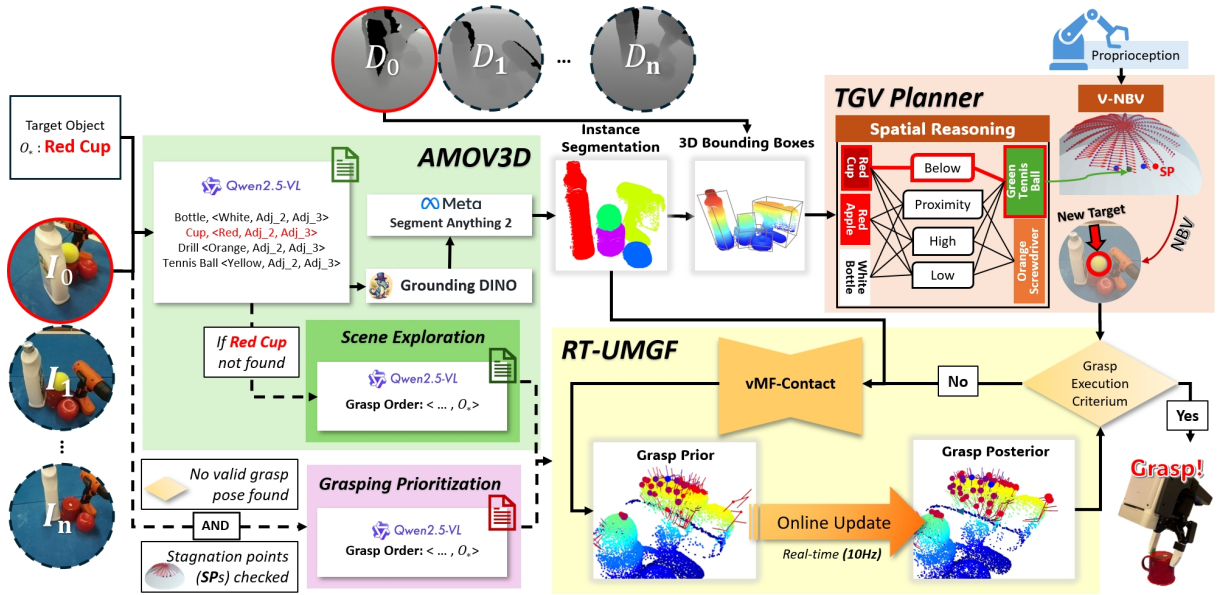


Fig. 2: Overview of VISO-Grasp. Our system consists of: 1) Adaptive Multi-View Open-Vocabulary 3D Object Detector (AMOV3D, Sec. IV-B): It leverages VLM (Qwen2.5-VL [11]), Grounding DINO [13] and Segment Anything 2 [12] to extract open-vocabulary 3D object representations. 2) Target-Guided View Planning (TGV-Planner, Sec. IV-C): The TGV-Planner dynamically adjusts viewpoints using spatial reasoning and *Velocity-field-based NBV* (V-NBV) to enhance object discovery. 3) Real-Time Uncertainty-guided Multi-view Grasp Fusion (RT-UMGF, Sec. IV-D) to enable occlusion-aware object-centric grasping, RT-UMGF fuses inferred grasps through real-time uncertainty-aware updates in 10Hz.

and associated attributes using historical *RGB-D images*, $\{\mathbf{I}^t, \mathbf{D}^t\}_{t=0}^T$, which is aided by instance segmentations from foundation models [12], [13]. The inferred attributes from AMOV3D further inform the following active view planner.

ii) Target-guided View Planner (TGV-Planner, Sec. IV-C) formulates a continuous velocity field $V(\mathbf{x}, \mathcal{O}')$, which is inferred together from the current position of the camera’s optical frame \mathbf{x} . This aims to guide the camera toward viewpoints that enhance the visibility of the target object.

iii) Real-Time Uncertainty-guided Multi-view Grasp Fusion (RT-UMGF, Sec. IV-D) provides 6-DoF grasp candidates G_t , which are inferred from the historical *depths* $\{\mathbf{D}^t\}_{t=0}^T$ by the grasp generator, which is pre-trained purely from simulated data. This aims to grasp the target object O_* and remove the occluders when necessary. Below, we detail each component.

B. Adaptive Multi-View Open-Vocabulary 3D Object Detection (AMOV3D)

In general, the AMOV3D module leverages FMs to achieve robust 3D object detection and segmentation in cluttered environments.

1) *3D Scene Representation Generation*: Given an input image \mathbf{I} and a prompt-specified target object O_* , a VLM generates structured descriptions for each detected object, prioritizing the identification of O_* to mitigate the risk of occlusion-induced omission. Each object instance is described using its label and three attributes: color (A_c), pattern (A_p), and spatial relation (A_s):

$$\mathcal{L} = \{(l_i, A_c^i, A_p^i, A_s^i) \mid i \in \mathcal{I}\}, \quad (1)$$

where \mathcal{I} represents the indices of detected objects, l_i denotes the instance label, and (A_c^i, A_p^i, A_s^i) are its associated attributes. This structured representation serves as a text prompt for Grounding DINO [13], which predicts the corresponding 2D bounding boxes. Each bounding box is refined by SAM2 [12] to obtain pixel-wise instance segmentation masks. The transformation from 2D pixel coordinates (u, v) to 3D world coordinates (X, Y, Z) is performed using the depth map \mathbf{D} and the camera intrinsics K^1 :

$$(X, Y, Z)^T = \mathbf{D}(u, v) * K^{-1} \cdot (u, v, 1)^T. \quad (2)$$

The 3D oriented Bounding Box (BB) is then computed using Principal Component Analysis (PCA) [39], aligning the box with the object’s principal axes. This structured scene representation serves as the foundation for downstream planning and grasp execution.

2) *Multi-View Object Verification and Fusion*: To provide TGV-Planner with a coherent scene representation, object detection across multiple viewpoints must maintain consistency. A historical object list \mathcal{O}' , which stores all objects currently perceived within the scene, is maintained and updated at each new viewpoint to refine detected object attributes and resolve occlusions. This list persists throughout the planning process, dynamically evolving as new observations improve object descriptions and scene understanding. Each viewpoint captured in the history $\{\mathbf{I}^t\}_{t=0}^T$ provides additional observations of the scene. During Next-Best-View (NBV) planning, AMOV3D updates the detected object set \mathcal{O}' by integrating new observations \mathcal{L}^t and verifying consistency

¹We denote $*$ as scalar multiplication and \cdot as the dot product.

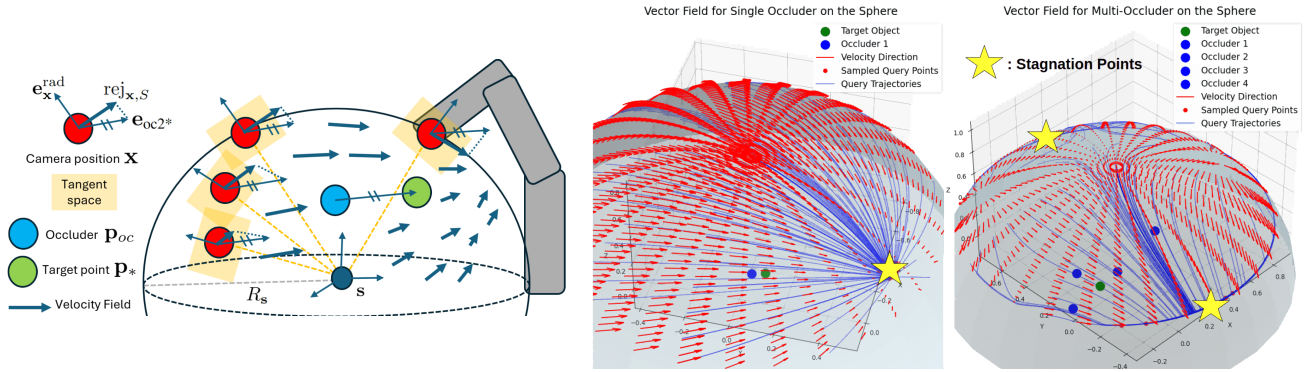


Fig. 3: Principle for TGV-Planner (left). Given single or multiple occluders (blue) and target center (green), the sampled camera viewpoints (red) follow the trajectories (blue curves) to expose the target within the spherical velocity field (red arrows) to achieve stagnation point(s) for single(middle)-/multi(right)-occluders. Note: To prevent unplannable movements, we truncate the tangential component of the field that aligns with the negative z -axis for elevation angles below 45° , ensuring no further downward movement.

with prior data:

$$\mathcal{O}^{t+1} = \mathcal{O}^t \cup \{O_i^t \mid O_i^t \in \mathcal{L}^t, \phi(O_i^t, \mathcal{O}^t, \mathbf{I}^t)\}. \quad (3)$$

$$\phi = \begin{cases} \text{Merge}(O_i^t, O_j^t), & \text{if } \exists O_j^t \in \mathcal{O}^t, \\ \text{Add } O_i^t, & \text{if } O_i^t \notin \mathcal{O}^t. \end{cases} \quad (4)$$

By continuously refining \mathcal{O}' and selecting viewpoints that maximize the visibility of *target object* O_* , AMOV3D improves the accuracy of object detection and occlusion reasoning, thereby enhancing the spatial reasoning accuracy of the TGV-Planner.

3) *Occlusion Reasoning and Scene Exploration*: If the target object O_* is not detected in the current view, the system initiates an occlusion reasoning process to infer potential obstructing objects. We employ VICL within the VLM to analyze the scene’s spatial configuration and identify potential occluders based on structured object descriptions \mathcal{L}^t and the current view image \mathbf{I}^t . To ensure accurate inference without additional computational cost, we enforce a Mixture-of-Reasoning-Experts (MoRE) [40] paradigm, guiding the VLM to concurrently assess spatial relations, material properties, and geometric constraints via structured prompting. The outputs are aggregated using voting to enhance occlusion identification robustness. The system then updates the grasp order, designating one potential occluder as the current grasp target. Subsequent modules adjust their reasoning and execution to remove this object, incrementally revealing the prompt-specified target object.

C. Target-Guided View Planner (TGV-Planner)

The TGV-Planner performs spatial reasoning to determine whether the target object O_* can be directly grasped or if an occluder must be removed first. If occlusions prevent direct grasping, the system prioritizes occluder removal. However, when direct removal is infeasible or does not sufficiently expose O_* , the *Velocity-field-based NBV (V-NBV)* module optimizes the view planning. *V-NBV* is applied to either enhance the visibility of O_* for grasp execution or to improve occluder removal efficiency by selecting a more informative camera perspective. To ensure robust grasp execution, the

TGV-Planner operates in parallel with the grasp inference from vMF-Contact. At the same time, the grasp execution criterion evaluates the grasp confidence of O_* continuously. If this exceeds a predefined threshold, the grasp is executed immediately. Otherwise, the grasps are updated in real-time as the viewpoint is adjusted to improve visibility.

1) *Rule-Based Spatial Reasoning*: Given the BBs obtained from AMOV3D, objects in the scene are assigned spatial relationships based on predefined geometric criteria:

- **Proximity**: Two objects are in proximity if their BBs intersect after isotropic expansion.

- **Below**: O_i is below O_j if the constraint $\min(\text{BB}_z(O_j)) - \max(\text{BB}_z(O_i)) > \gamma_{\text{Below}}$ is satisfied and their projections onto the xy -plane overlap. γ_{Below} is the height threshold.

- **High / Low**: An object O_i is considered **High** relative to O_j if $h(O_i, O_j) = \max(\text{BB}_z(O_i)) - \max(\text{BB}_z(O_j))$ satisfies $h(O_i, O_j) > \gamma_{HL}$ (High) or $h(O_i, O_j) < -\gamma_{HL}$ (Low), where γ_{HL} is a predefined height threshold.

The grasp execution strategy is determined by the assigned spatial relations: i) If O_* is **Below** any object, it is prioritized for grasping to expose O_* ; ii) If an object in **Proximity** to O_* is also classified as **High**, the TGV-Planner evaluates whether grasping O_* is feasible or if occluder removal is required. iii) **Low** occluders are ignored as they do not obstruct grasp execution. iv) Crucially, **High** occluders trigger the *V-NBV* module to optimize viewpoint adjustment.

2) *V-NBV with Single Occluder*: Given the current camera position \mathbf{x} , the system queries a continuous velocity field constructed from the spatial relationships between all **High** occluders and O_* to guide viewpoint adjustment by: $V : \mathbb{R}^3 \times \mathcal{O} \rightarrow \mathbb{R}^3, V(\mathbf{x}, \mathcal{O}_h) = \dot{\mathbf{x}}$.

We begin with single occluder scenario for a better understanding of our approach, where the principle is illustrated in Fig. 3 (left). For instance, a camera is constrained on a hemisphere centered at \mathbf{s} with radius R_s by position \mathbf{x} . Given two center points of the target object \mathbf{p}_* and single occluder \mathbf{p}_{oc} inside the hemisphere respectively, the objective is to generate camera motions that naturally follow the geodesic that improves the visibility of the target. To achieve this, the

velocity at \mathbf{x} is computed as follows:

First, the normalized vector \mathbf{e}_{oc2^*} from \mathbf{p}_{oc} to \mathbf{p}_* and the radial vector from the sphere center \mathbf{s} to \mathbf{x} is given by:

$$\mathbf{e}_{oc2^*} = \frac{\mathbf{p}_* - \mathbf{p}_{oc}}{\|\mathbf{p}_* - \mathbf{p}_{oc}\|}, \quad \mathbf{e}_{\mathbf{x}}^{\text{rad}} = \frac{\mathbf{x} - \mathbf{s}}{R_s}. \quad (5)$$

To obtain the velocity direction $\mathbf{e}_{\dot{\mathbf{x}}}$ within the tangent space of sphere (\mathbf{s}, R_s) allocated on \mathbf{x} , which needs to align with \mathbf{e}_{oc2^*} , we further calculate the normalized rejection of \mathbf{e}_{oc2^*} onto the radial vector $\mathbf{e}_{\mathbf{x}}^{\text{rad}}$, namely $\mathbf{e}_{\dot{\mathbf{x}}}$ by:

$$\text{rej}_{\mathbf{x},S} = \mathbf{e}_{oc2^*} - (\mathbf{e}_{oc2^*} \cdot \mathbf{e}_{\mathbf{x}}^{\text{rad}}) * \mathbf{e}_{\mathbf{x}}^{\text{rad}}, \quad \mathbf{e}_{\dot{\mathbf{x}}} = \frac{\text{rej}_{\mathbf{x},S}}{\|\text{rej}_{\mathbf{x},S}\|}. \quad (6)$$

Notably, we define the field strength coefficient β , which is computed based on the angular relationship between \mathbf{x} , \mathbf{p}_{oc} , and \mathbf{p}_* using:

$$\beta = \frac{1}{\pi} * \arccos\left(\frac{\mathbf{e}_{oc2^*} \cdot (\mathbf{x} - \mathbf{p}_{oc})}{\|\mathbf{x} - \mathbf{p}_{oc}\|}\right). \quad (7)$$

This ensures that the field's strength β smoothly decreases from 1 to 0 as the camera moves from the occluder's side to the opposite of the target, converging at the stagnation point ($\beta = 0$), which represents the final Next-Best View (NBV). Finally, the velocity field is calculated by: $\dot{\mathbf{x}} = V(\mathbf{x}, \mathcal{O}_h) = \beta \cdot \mathbf{e}_{\dot{\mathbf{x}}}$.

3) *Field Superposition for Multi-occluders*: For multiple occluders, the camera motion is determined by the weighted combination of individual occluder's influence using the superposition principle. Given a set of M occluders, the aggregated velocity is computed as the weighted sum by β following: $\dot{\mathbf{x}}_M = \sum_{i=1}^M \beta_m \cdot \mathbf{e}_{\dot{\mathbf{x}},m}$.

Here we illustrate two examples of single/multi-occluders in Fig. 3 (middle & right). In addition, to ensure stable and controlled motion, *both the center and the bounding box corners* of each potential occluder are incorporated as occluder points \mathbf{p}_{oc} s into the field generation.

D. Real-Time Uncertainty-guided Multi-view Grasp Fusion (RT-UMGF)

In general, RT-UMGF continuously refines grasp predictions by fusing multi-view observations. It applies Bayesian updates using a von Mises-Fisher (vMF) distribution, ensuring stable grasp selection in cluttered environments (Fig. 4). To maximize the utility of multi-view sensor data and enhance grasp generation robustness, we employ the vMF-Contact [17], which is adapted to real-time (10Hz) inference through pose-centric uncertainty modeling and buffered online Bayesian fusion. Our approach diverges from recent work on real-time grasp fusion [41], which rely solely on grasp quality scores and lack geometric grasp representation [18], by integrating directional von Mises-Fisher (vMF) distribution into our Bayesian update process [42].

In runtime, assume that the raw point cloud is captured in the time frame t , a set of contact grasps [18] is inferred as $G_t = \{\mathbf{g}_n^t = (\mathbf{c}, \mu_{\mathbf{c}}, \kappa_{\mathbf{c}}, \Delta_{\mathbf{c}}, w_{\mathbf{c}}, q_{\mathbf{c}})^t \mid n \in N\}$ that parameterizes: 1) The queried contact points $\mathbf{c} \in \mathbf{R}^3$; 2) Baseline vector distributions: $p(\mathbf{b}|\mathbf{c}) = \text{vMF}(\mathbf{b}|\mu_{\mathbf{c}}, \kappa_{\mathbf{c}}) =$

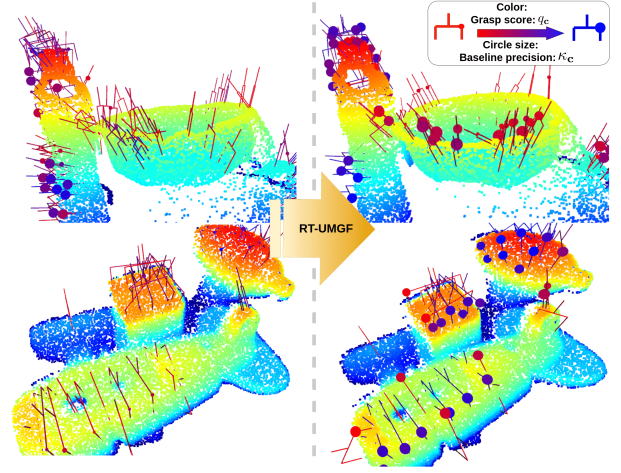


Fig. 4: RT-UMGF update, where grasps are merged according to corresponding categories. Confident grasps on high grasp quality $q_{\mathbf{c}}$ and directional precision $\kappa_{\mathbf{c}}$ can be observed consequently and fused to enhance their confidence further.

$Z(\kappa_{\mathbf{c}}) \exp(\kappa_{\mathbf{c}} \mu_{\mathbf{c}}^{\top} \mathbf{b})$, with $\mu_{\mathbf{c}}$ the mean direction of the baseline and $\kappa_{\mathbf{c}}$ as the directional precision. $Z(\kappa_{\mathbf{c}})$ is the normalization factor; 3) The quantized approach vector, represented by $\Delta_{\mathbf{c}} = \{\delta_{\mathbf{c}}^k\}_{k=0,\dots,K}$ represents discrete categorical bin scores that define a direction constrained to lie on a plane perpendicular to a given baseline. 4) The grasp width $w_{\mathbf{c}}$ and 5) The contact quality score $q_{\mathbf{c}}$.

1) *Grasp Categorization*: Technically, newly generated grasps G_t^t are categorized into two distinct groups based on their similarity to previously fused grasps G_{t-1} : The first contact point set C_t^{new} comprises grasps located significantly away from any existing fused grasps, which then undergo a process of *self-fusion* before being directly involved into the grasp buffer. This typically includes grasps associated with previously unseen objects. Conversely, if the new grasps are proximal to any previously fused grasps, they are fused with the closest existing grasp by *cross-fusion*, denoted as C_t^{pro} . As inspired by [41], the rule of distinction is:

$$C_t^{\text{pro}} = \{ \mathbf{c}_j^t \in C_t^t \mid \|\mathbf{c}_j^t - \mathbf{c}_i^{t-1}\|_2 < \gamma_d, \quad (8)$$

$$1 - \cos(\mu_{\mathbf{c}_i^{t-1}}, \mu_{\mathbf{c}_j^t}) > \gamma_{\theta} \}, \quad (9)$$

$$C_t^{\text{new}} = C_t^t \setminus C_t^{\text{pro}}. \quad (10)$$

To ensure the multi-modality of scene-level grasps, the similarity criterion between proximal and new grasps is determined based on 1) contact point distance limit γ_d and 2) cosine similarity of baseline mean directions constraint γ_{θ} . Notably, different from *cross-fusion*, which applies these criteria based on existing fused grasp \mathbf{g}_i^{t-1} , *self-fusion* does not rely on a predefined cluster center, where we applied DBSCAN [43] to identify the clusters.

2) *Contact-grasp Fusion*: Both *self/cross-fusion* processes are designed as follows: For contact point positions \mathbf{c}_j^t , we adopt the same regime as [41] utilizing weighted sum by grasp quality q : $\mathbf{c}_i^t = \frac{q_i^{t-1} \mathbf{c}_i^{t-1} + \sum_j (q_j \mathbf{c}_j^t)}{q_i^{t-1} + \sum_j (q_j)}$ with grouped element indices $j \in J$. For baselines and approach vectors,

the conjugate prior of vMF baseline distribution is initialized as: $\mu_{\mathbf{c}} \sim \text{vMF}(\mu_{\mathbf{c}}^0, \kappa_{\mathbf{c}}^0)$ for $t = 0$.

For the Bayesian inference in time frame t , we may update the posterior distribution following the rule of the exponential family [44] by:

$$\mu_{\mathbf{c}}^t = \frac{\kappa_{\mathbf{c}}^{t-1} \mu_{\mathbf{c}}^{t-1} + \sum_j \kappa_{\mathbf{c}_j}^t \mu_{\mathbf{c}_j}^t}{\kappa_{\mathbf{c}}^{t-1} + \sum_j \kappa_{\mathbf{c}_j}^t}, \kappa_{\mathbf{c}}^t = \kappa_{\mathbf{c}}^{t-1} + \sum_j \kappa_{\mathbf{c}_j}^t. \quad (11)$$

The approach categories are updated by:

$$\delta_{\mathbf{c}}^{t,i} = \delta_{\mathbf{c}}^{t-1,i} + \sum_j \delta_{\mathbf{c}_j}^{t,i}. \quad (12)$$

Here $\kappa_{\mathbf{c}}^t$ represents the precision on the observed mean likelihood $\mu_{\mathbf{c}}$. We refer interesting readers for the theoretical background to [17].

Throughout the runtime, we define three termination criteria for executing a grasp: i) The highest estimated grasp quality $q_{\mathbf{c}}^*$ must exceed a predefined threshold q_{\max} . ii) The directional uncertainty of the corresponding grasp from condition i), $\kappa_{\mathbf{c}}^*$, must surpass a given threshold κ_{\max} . iii) When the magnitude of the queried velocity field approaches a stagnation point, i.e., $\|\dot{\mathbf{x}}\| \approx 0$, the current best grasp is executed immediately if condition i) is fulfilled. Otherwise the grasping prioritization, where the optimal grasp order of potential obstructers will be queried from the VLM and the target will be switched.








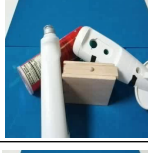




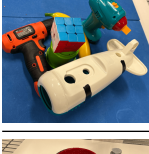



V. EXPERIMENTS

During the experiments, we aim to analyze the contributions of each component in VISO-Grasp. We investigate the advantages of TGV-Planner over fixed initial (*Init view*) and *Top-down* view setups in enhancing the semantic understanding of the VLM, improving the grasp update process, and refining grasp proposals. Furthermore, we include a baseline of identical settings as VISO-Grasp without grasp fusion (*w/o GF*), which aims to assess the effectiveness of RT-UMGF as another focus of our studies. In this setting, grasp predictions are inferred directly on the NBV in the experiments. Additionally, we incorporate reconstruction-based NBV planning approaches [6] (*Breyer's*) for comparison.

A. Experiment setups

Our experimental setup consists of a UR10e robotic arm integrated with an Orbbec Femto Mega RGB-D camera for scene perception and a Robotiq 2F85 parallel-jaw gripper for object manipulation (See Fig. 1). In AMOV3D, we employ Qwen2.5-VL-72B-Instruct-AWQ [11] as the Visual-Language Model (VLM) to enhance high-level reasoning and scene understanding. For RT-UMGF, we leverage the vMF-Contact model with PointNeXt-B [45] as the backbone, which is pre-trained on a purely simulated dataset according to [17], to infer the scene-level grasp and uncertainty prediction. The following setups are considered in our real-world experiments for ablation studies: *Top-down* and *Init view* are used to evaluate the performance *w/o V-NBV* in the TGV-Planner. In these experiments, all other components, including AMOV-3D, spatial reasoning in the TGV-Planner

TABLE I: Real-world experiment setups and results.

Setup	Initial View	Scene index: O_*		
		Method	FS(5)	#GA GSR(%)
(1): Red Cup				
		<i>Breyer's</i>	2	4.0 27.8
		<i>Top-down</i>	2	3.6 61.1
		<i>Init view</i>	5	3.4 88.2
		<i>w/o GF</i>	5	3.6 94.44
		<i>Our's</i>	5	3.4 88.2
(2): Jello Box				
		<i>Breyer's</i>	0	4.0 6.3
		<i>Top-down</i>	4	2.2 70.0
		<i>Init view</i>	3	3.0 33.3
		<i>w/o GF</i>	3	3.6 66.7
		<i>Our's</i>	4	3.0 80.0
(3): Pringles				
		<i>Breyer's</i>	0	4.0 11.8
		<i>Top-down</i>	2	2.8 30.8
		<i>Init view</i>	4	3.2 68.8
		<i>w/o GF</i>	5	2.0 90.0
		<i>Our's</i>	5	1.6 87.5
(4): Baseball				
		<i>Breyer's</i>	1	4.0 18.8
		<i>Top-down</i>	4	3.2 62.5
		<i>Init view</i>	2	3.2 62.5
		<i>w/o GF</i>	3	3.4 82.4
		<i>Our's</i>	5	2.8 85.7
(5): Green Pear				
		<i>Breyer's</i>	0.0	4.0 0.0
		<i>Top-down</i>	4	2.6 46.1
		<i>Init view</i>	3	3.6 88.9
		<i>w/o GF</i>	5	3.2 68.8
		<i>Our's</i>	5	2.4 100.0
(6): Purple Cup				
		<i>Breyer's</i>	2	4.0 85.0
		<i>Top-down</i>	1	4.0 70.0
		<i>Init view</i>	1	3.6 27.8
		<i>w/o GF</i>	1	4.0 45.0
		<i>Our's</i>	3	4.0 85.0
(7): Yellow Peach in the Green Cup				
		<i>Breyer's</i>	0	4.0 22.2
		<i>Top-down</i>	2	4.0 75.0
		<i>Init view</i>	3	3.8 63.2
		<i>w/o GF</i>	4	3.6 83.3
		<i>Our's</i>	4	3.2 75.0
(8): Tennis Ball in the Red Cup				
		<i>Breyer's</i>	0	4.0 21.1
		<i>Top-down</i>	2	4.0 45.0
		<i>Init view</i>	3	4.0 80.0
		<i>w/o GF</i>	2	3.8 55.0
		<i>Our's</i>	4	3.8 84.2

and RT-UMGF, remain identical as with our approach (denoted as *Our's*). Specifically, the *Init view* setup maintains a fixed viewpoint at a 45° elevation angle in the horizontal coordinate, identical to the configuration in *Ours*. Besides, to evaluate RT-UMGF, we include a comparison setting that excludes the grasp fusion process (*w/o GF*), where the best grasp is determined by the current inference after reaching the NBV. Finally, we incorporate *Breyer's* closed-loop NBV planner [8] as a baseline. For fair comparisons, we manually annotated the BBs involving all objects in proximity to the targets. These BBs are axis-aligned with the robot frame,

allowing the planner to perform targeted view planning.

To systematically evaluate the performance of VISO-Grasp, we conducted comprehensive experiments across eight distinct scenes, each varying in difficulty, occlusion severity, and occlusion type (Table. I). The following metrics are used in our experiments: Final Success (**FS**) quantifies the number of successful grasps of the target object; Grasp Attempts (**#GA**) counts the average number of grasp attempts required for each scene until the final target object is successfully grasped. If the target is not grasped by the end of the trial, this value is not recorded; Grasp Success Rate (**GSR**) measures the overall success rate of grasping any objects throughout the entire trial. A trial is aborted if any of the following conditions occur four times after the scene is reset to its previous state: i) Grasp failure; ii) More than one object is affected after grasp execution, leading to unintended scene changes or exposure of the target object; iii) A collision of the gripper fingers triggers an emergency stop of the robot.

TABLE II: Average performance results

Method	AFSR (%)	#AGA	AGSR (%)
<i>Breyer's</i>	12.50	4.00	22.5
<i>Top-down</i>	52.50	3.30	56.50
<i>Init view</i>	60.00	3.48	64.08
<i>w/o GF</i>	70.00	3.43	73.19
<i>Ours</i>	87.50	3.10	83.86

B. Experiment evaluation

The experimental results from 8 trials across various scenes are presented in Table. I, which demonstrate the effectiveness of VISO-Grasp (*Ours*), achieving high success in reaching the target with minimal grasp attempts across diverse occlusion scenarios. Compared to baselines with static viewpoints (*Top-down* and *Init view*), VISO-Grasp leverages the synergy between TGV-Planner and RT-UMGF to actively explore viewpoints and refine grasp selection. This facilitate the grasp success on highly occluded or fully unobservable targets. For instance, in Scene (1), where the target object is partially obstructed by a wooden block and positioned beneath a tennis ball, both *Top-down* and *Init view* exhibit limited success due to frequent collisions with the wooden block. In Scene (3), although *Top-down* provides an unobstructed view, it frequently misclassifies objects due to insufficient semantic understanding by the VLM (e.g., recognizing the "Pringle box" as a "Red can"). In contrast, VISO-Grasp mitigates these limitations by leveraging multi-view integration from AMOV3D, refining predictions through diverse observations facilitated by the TGV Planner. As a result, it achieves success in all five trials with the GSR exceeding 80%. In addition, Breyer's NBV planner, which lacks object-centric view planning, struggles to effectively expose the target object, even when predefined BBs are provided around the target. This highlighted the importance of diverse viewpoints generated by TGV-Planner, which enhance comprehensive scene understanding and facilitate more informed grasp planning in general.

In addition, a detailed ablation study on RT-UMGF reveals its crucial role in reducing grasp attempts and improving grasp success by prioritizing high-quality, low-uncertainty grasps over time. In Scene (4), the setting without grasp fusion (w/o GF) achieves a grasp success rate (GSR) of 82.4% with an average of 3.2 grasp attempts. By contrast, VISO-Grasp further increases GSR to 85.7% while reducing grasp attempts to 2.8, confirming that multi-view grasp aggregation refines grasp selection by integrating uncertain and low-quality grasp hypotheses, thereby enhancing selection robustness through the fusion of multiple grasp candidates.

The importance of NBV and grasp fusion becomes more evident when analyzing scenarios where information is severely lacking and semantic information is easily misinterpreted. In Scene (6), where the purple cup is at the bottom of the scene, understanding spatial hierarchy is crucial to successfully grasp the target, as it is one of the most challenging scenes overall. *Breyer's* NBV planner successfully removes cups on top with high GSR but occasionally fails to grasp purple cup due to insufficient object-centric guidance, while *Ours* achieves comparable success with identical GSR. The synergy between NBV and grasp fusion is crucial to the success of *Ours*, since NBV ensures that occluded targets become visible, but grasp predictions from new viewpoints may still exhibit uncertainty.

In general, Table. II presents the overall performance of each method in terms of Average Final Success Rate (**AFSR**), Average Grasp Attempts (**#AGA**) and Average Grasp Success Rate (**AGSR**). *Breyer's* reflects its limits in revealing the target's critical geometry and underscores how the lack of target-focused visibility. *Top-down* and *Init view* achieve moderate performance, which struggles in heavy occlusion and fixed view limit. *Top-down* benefits from its view for some cases and tries to grasp the target in an efficient manner, but results in lower AGSR due to the lack of semantic information. *w/o GF* outperforms all static baselines with 70.00% of AFSR by employing next-best-view planning to reduce occlusion. However, without RT-UMGF, it still produces suboptimal grasps more frequently. In summary, VISO-Grasp (*Ours*) demonstrates the best overall performance by 87.50% combining active NBV exploration with real-time grasp fusion.

VI. CONCLUSION

We develop VISO-Grasp, a novel vision-language-informed system for target-oriented grasping in highly unstructured environments including entire invisibility. By integrating a Vision-Language Model (VLM) with object-centric View planning and real-time uncertainty-driven grasp fusion, our system enhances scene understanding and improves grasp success through continuous velocity fields and semantic spatial reasoning for adaptive grasping in occluded environments with complete invisibility. VISO-Grasp leverages robust multi-view aggregation to refine grasp selection by integrating uncertain grasp hypotheses, ensuring superior stability and accuracy. Experimental results show that VISO-Grasp achieves a final success rate by 87.5% while requiring

the fewest grasp attempts among all baselines, demonstrating its efficiency in occlusion-aware grasping.

Certain limitations remain for VISO-Grasp that warrant further investigation. First, the system’s reliance on the VLM introduces a computational bottleneck, as VLM inference incurs non-negligible latency, especially for VACL’s context, which demands extensive reasoning over multi-modal inputs. Second, our experiments assume a quasi-static scene, which does not account for highly dynamic objects or external disturbances. Future work could explore adaptive NBV strategies in real-time scene variations.

REFERENCES

- [1] R. Newbury, M. Gu, L. Chumbley, A. Mousavian, C. Eppner, J. Leitner, J. Bohg, A. Morales, T. Asfour, D. Kragic *et al.*, “Deep learning approaches to grasp synthesis: A review,” *TRO*, vol. 39, no. 5, 2023.
- [2] H.-S. Fang, C. Wang, M. Gou, and C. Lu, “Graspnet-1billion: A large-scale benchmark for general object grasping,” in *CVPR*, 2020.
- [3] M. Gilles, Y. Chen, E. Z. Zeng, Y. Wu, K. Furmans, A. Wong, and R. Rayyes, “Metagraspnetv2: All-in-one dataset enabling fast and reliable robotic bin picking via object relationship reasoning and dexterous grasping,” *TASE*, 2023.
- [4] 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,” *arXiv preprint arXiv:1703.09312*, 2017.
- [5] Q. Dai, Y. Zhu, Y. Geng, C. Ruan, J. Zhang, and H. Wang, “Graspnerf: Multiview-based 6-dof grasp detection for transparent and specular objects using generalizable nerf,” in *ICRA*, 2023.
- [6] M. Breyer, J. J. Chung, L. Ott, R. Siegwart, and J. Nieto, “Volumetric grasping network: Real-time 6 dof grasp detection in clutter,” in *CoRL PMLR*, 2021.
- [7] M. Gilles, K. Furmans, and R. Rayyes, “Metamvuc: Active learning for sample-efficient sim-to-real domain adaptation in robotic grasping,” *RAL*, 2025.
- [8] M. Breyer, L. Ott, R. Siegwart, and J. J. Chung, “Closed-loop next-best-view planning for target-driven grasping,” in *IROS*, 2022.
- [9] X. Zhang, D. Wang, S. Han, W. Li, B. Zhao, Z. Wang, X. Duan, C. Fang, X. Li, and J. He, “Affordance-driven next-best-view planning for robotic grasping,” in *CoRL*, 2023.
- [10] H. Ma, M. Shi, B. Gao, and D. Huang, “Active perception for grasp detection via neural graspness field,” in *NeurIPS*, 2024.
- [11] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang *et al.*, “Qwen2.5-vl technical report,” *arXiv preprint arXiv:2502.13923*, 2025.
- [12] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rüdle, C. Rolland, L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár, and C. Feichtenhofer, “Sam 2: Segment anything in images and videos,” *arXiv preprint arXiv:2408.00714*, 2024.
- [13] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, Q. Jiang, C. Li, J. Yang, H. Su *et al.*, “Grounding dino: Marrying dino with grounded pre-training for open-set object detection,” in *ECCV*. Springer, 2024.
- [14] S. Noh, J. Kim, D. Nam, S. Back, R. Kang, and K. Lee, “Graspsam: When segment anything model meets grasp detection,” *arXiv preprint arXiv:2409.12521*, 2024.
- [15] G. Tziafas and H. Kasaei, “Towards open-world grasping with large vision-language models,” in *CoRL*, 2024.
- [16] J. Xu, S. Jin, Y. Lei, Y. Zhang, and L. Zhang, “Rt-grasp: Reasoning tuning robotic grasping via multi-modal large language model,” in *IROS*. IEEE, 2024, pp. 7323–7330.
- [17] Y. Shi, E. Welte, M. Gilles, and R. Rayyes, “vmf-contact: Uncertainty-aware evidential learning for probabilistic contact-grasp in noisy clutter,” *arXiv preprint arXiv:2411.03591*, 2024.
- [18] M. Sundermeyer, A. Mousavian, R. Triebel, and D. Fox, “Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes,” in *ICRA*, 2021.
- [19] H. Liang, X. Ma, S. Li, M. Gerner, S. Tang, B. Fang, F. Sun, and J. Zhang, “Pointnetgpd: Detecting grasp configurations from point sets,” in *ICRA*, 2019.
- [20] 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.
- [21] D. Morrison, P. Corke, and J. Leitner, “Multi-view picking: Next-best-view reaching for improved grasping in clutter,” in *ICRA*, 2019.
- [22] R. Bajcsy, Y. Aloimonos, and J. K. Tsotsos, “Revisiting active perception,” *Autonomous Robots*, vol. 42, pp. 177–196, 2018.
- [23] H. H. González-Banos and J.-C. Latombe, “Navigation strategies for exploring indoor environments,” *IJRR*, 2002.
- [24] L. Jin, X. Zhong, Y. Pan, J. Behley, C. Stachniss, and M. Popović, “Activegs: Active scene reconstruction using gaussian splatting,” *arXiv preprint arXiv:2412.17769*, 2024.
- [25] L. Jin, X. Chen, J. Rückin, and M. Popović, “Neu-nbv: Next best view planning using uncertainty estimation in image-based neural rendering,” in *2023 IROS*, 2023.
- [26] R. Zeng, W. Zhao, and Y.-J. Liu, “Pc-nbv: A point cloud based deep network for efficient next best view planning,” in *IROS*, 2020.
- [27] M. Liu, Y. Zhu, H. Cai, S. Han, Z. Ling, F. Porikli, and H. Su, “Partslip: Low-shot part segmentation for 3d point clouds via pretrained image-language models,” in *CVPR*, 2023.
- [28] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei, “Voxposer: Composable 3d value maps for robotic manipulation with language models,” *arXiv preprint arXiv:2307.05973*, 2023.
- [29] P. Sharma, A. Torralba, and J. Andreas, “Skill induction and planning with latent language,” *arXiv preprint arXiv:2110.01517*, 2021.
- [30] S. Jin, J. Xu, Y. Lei, and L. Zhang, “Reasoning grasping via multi-modal large language model,” *arXiv preprint arXiv:2402.06798*, 2024.
- [31] Y. Sun, S. Ma, R. Madaan, R. Bonatti, F. Huang, and A. Kapoor, “Smart: Self-supervised multi-task pretraining with control transformers,” *arXiv preprint arXiv:2301.09816*, 2023.
- [32] I. Radosavovic, T. Xiao, S. James, P. Abbeel, J. Malik, and T. Darrell, “Real-world robot learning with masked visual pre-training,” in *CoRL PMLR*, 2023, pp. 416–426.
- [33] R. Firoozj, J. Tucker, S. Tian, A. Majumdar, J. Sun, W. Liu, Y. Zhu, S. Song, A. Kapoor, K. Hausman *et al.*, “Foundation models in robotics: Applications, challenges, and the future,” *CoRR*, 2023.
- [34] Y. Hu, Q. Xie, V. Jain, J. Francis, J. Patrikar, N. V. Keetha, S. Kim, Y. Xie, T. Zhang, S. Zhao *et al.*, “Toward general-purpose robots via foundation models: A survey and meta-analysis,” *CoRR*, 2023.
- [35] N. P. Bhatt, Y. Yang, R. Siva, D. Milan, U. Topcu, and Z. Wang, “Know where you’re uncertain when planning with multimodal foundation models: A formal framework,” *arXiv preprint arXiv:2411.01639*, 2024.
- [36] C. Tang, D. Huang, W. Dong, R. Xu, and H. Zhang, “Foundationgrasp: Generalizable task-oriented grasping with foundation models,” *TASE*, 2025.
- [37] Y. Qian, X. Zhu, O. Biza, S. Jiang, L. Zhao, H. Huang, Y. Qi, and R. Platt, “Thinkgrasp: A vision-language system for strategic part grasping in clutter,” in *2nd CoRL Workshop on Learning Effective Abstractions for Planning*.
- [38] V. Sripada, S. Carter, F. Guerin, and A. Ghalamzan, “Ap-vlm: Active perception enabled by vision-language models,” *arXiv preprint arXiv:2409.17641*, 2024.
- [39] D. Dimitrov, C. Knauer, K. Kriegel, and G. Rote, “On the bounding boxes obtained by principal component analysis,” in *EuroCG*, 2006.
- [40] C. Si, W. Shi, C. Zhao, L. Zettlemoyer, and J. Boyd-Graber, “Getting more out of mixture of language model reasoning experts,” *arXiv preprint arXiv:2305.14628*, 2023.
- [41] J. Zhang, N. Gireesh, J. Wang, X. Fang, C. Xu, W. Chen, L. Dai, and H. Wang, “Gamma: Graspability-aware mobile manipulation policy learning based on online grasping pose fusion,” in *ICRA*, 2024.
- [42] K. V. Mardia and S. El-Atoum, “Bayesian inference for the von mises-fisher distribution,” *Biometrika*, vol. 63, no. 1, 1976.
- [43] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, “A density-based algorithm for discovering clusters in large spatial databases with noise,” in *kdd*, vol. 96, no. 34, 1996.
- [44] B. Charpentier, O. Borchert, D. Zügner, S. Geisler, and S. Günnemann, “Natural Posterior Network: Deep Bayesian Predictive Uncertainty for Exponential Family Distributions,” in *ICLR*, 2022.
- [45] G. Qian, Y. Li, H. Peng, J. Mai, H. Hammoud, M. Elhoseiny, and B. Ghanem, “Pointnext: Revisiting pointnet++ with improved training and scaling strategies,” in *NeurIPS*, 2022.