# Language-Guided Grasp Detection with Coarse-to-Fine Learning for Robotic Manipulation

Zebin Jiang [2†], Tianle Jin [2†], Xiangtong Yao [2], Alois Knoll [2], *Fellow, IEEE*, Hu Cao [1,2*], *Member, IEEE*

*Abstract*—Grasping is one of the most fundamental challenging capabilities in robotic manipulation, especially in unstructured, cluttered, and semantically diverse environments. Recent researches have increasingly explored language-guided manipulation, where robots not only perceive the scene but also interpret task-relevant natural language instructions. However, existing language-conditioned grasping methods typically rely on shallow fusion strategies, leading to limited semantic grounding and weak alignment between linguistic intent and visual grasp reasoning. In this work, we propose Language-Guided Grasp Detection (LGGD) with a coarse-to-fine learning paradigm for robotic manipulation. LGGD leverages CLIP-based visual and textual embeddings within a hierarchical cross-modal fusion pipeline, progressively injecting linguistic cues into the visual feature reconstruction process. This design enables fine-grained visual-semantic alignment and improves the feasibility of the predicted grasps with respect to task instructions. In addition, we introduce a language-conditioned dynamic convolution head (LDCH) that mixes multiple convolution experts based on sentence-level features, enabling instruction-adaptive coarse mask and grasp predictions. A final refinement module further enhances grasp consistency and robustness in complex scenes. Experiments on the OCID-VLG and Grasp-Anything++ datasets show that LGGD surpasses existing language-guided grasping methods, exhibiting strong generalization to unseen objects and diverse language queries. Moreover, deployment on a real robotic platform demonstrates the practical effectiveness of our approach in executing accurate, instruction-conditioned grasp actions. *The code will be released publicly upon acceptance.*

*Index Terms*—Robotic manipulation, foundation model, vision-language fusion, language-guided grasping.

## I. INTRODUCTION

ROBOTIC grasping serves as a fundamental yet capability for autonomous manipulation, underpinning tasks ranging from industrial logistics and manufacturing to household assistance and service robotics [1], [2], [3], [4]. The ultimate objective of robotic grasping is to endow robots with the ability to reliably detect, localize, and physically grasp objects in an open-world setting, handling diverse objects, occlusions, and ambiguity arising from human intention. Despite extensive research, achieving robust grasping in realistic environments, especially when guided by natural language instructions, remains a significant challenge.

Traditionally, analytical and model-based methods formed the first generation of grasp detection solutions, relying on classical mechanics, friction cone analysis, and 3D geometric modeling [5], [6]. These approaches assume precise object models and compute feasible contact points using force-closure or wrench-space formulations. While theoretically
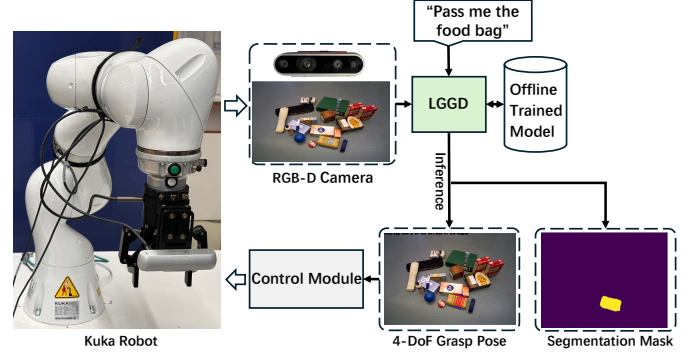
† Equal contribution, ∗ Corresponding author.



Fig. 1: Overview of our system: An RGB-D camera mounted on the robot's wrist captures visual data of objects to be grasped. Our proposed LGGD generates 4-DoF grasp poses based on the RGB image and language query during the inference process. These generated grasp poses are then utilized by the control module to plan and execute robot trajectories for pick-and-place tasks.

grounded, their applicability sharply declines when object models are incomplete or unavailable, or when objects are novel, deformable, or partially occluded-conditions that frequently occur in real domestic and industrial scenes. Moreover, collecting accurate mesh models for every object is impractical, and performance suffers when sensing noise, camera perspective distortion, or cluttered backgrounds are present. The proliferation of deep learning triggered a paradigm shift from geometry-dependent methods to data-driven grasp inference [7], [8], [9], [10]. Discriminative grasping approaches conceptualize grasp detection as a selection problem: sample a pool of grasp hypotheses and evaluate them to select the best one. Two-stage methods, exemplified by Lenz *et al.* [7] and GQCNN [8], often produce high-quality grasp candidates but suffer from several limitations: (1) grasp candidate sampling is a bottleneck, (2) performance degrades when candidate coverage is insufficient, and (3) the separation of sampling and evaluation introduces error propagation. One-stage discriminative detectors [11], [12] merge hypothesis generation and scoring into a single process, improving computational efficiency, yet they still operate within a candidate-based paradigm and do not fundamentally exploit continuous per-pixel grasp reasoning.

This limitation motivated the rise of generative grasp detection, which reconceptualizes grasping as dense prediction over the image domain [13]. Instead of evaluating discrete hypotheses, generative approaches directly output grasp qual-

ity, orientation, and gripper width per pixel. This conceptual leap eliminates the need for sampling, significantly improves inference speed, and harmonizes grasping with pixel-level perception pipelines such as semantic segmentation. Advancements including GR-ConvNet [14], SE-enhanced designs [15], ASPP-based receptive-field scaling [16], and Transformer-based perception [17], [18] have all contributed to improving grasping robustness and accuracy under visual ambiguity. However, generative methods remain fundamentally *purely visual*. They operate solely on image input and do not incorporate semantic intent or contextual constraints from natural language instructions. The inability of visually driven grasping models to integrate task intent forms a crucial barrier to human-robot collaboration. Humans rarely specify grasps purely visually; instead, we express intentions linguistically: *"pick up the red screwdriver" "grab the knife by the handle," "pick up the leftmost apple," "hand me the shallow bowl, not the deep one."* Such statements encode object identity, attributes, relational placement, spatial disambiguation, safety semantics, and implicit operational constraints. A purely visual model has no access to this semantic layer of meaning and may inadvertently grasp a wrong object.

The emergence of large-scale vision-language models (VLMs), such as CLIP [19], BLIP [20], and LLaVA [21], has provided a powerful foundation for connecting linguistic reasoning with visual grounding. These models align language tokens with visual embeddings across broad conceptual spaces and have shown early promise in robot manipulation settings [22], [23], [24]. However, existing language-guided grasping systems primarily attach language to vision at a single stage of processing, either early (shared encoder input) [25], or late (final decision stage) [26], [27]. Such shallow or uni-directional fusion typically fails to maintain semantic context throughout the grasp inference pipeline, resulting in weak grounding of task intent, particularly for fine-grained grasping tasks (e.g., grasping a specific object part).

In this work, we address this gap by introducing **LGGD**, a coarse-to-fine language-guided grasp detection framework that performs multi-level vision-language fusion. Unlike prior approaches that fuse language only once, LGGD injects linguistic intent throughout the perception pipeline. First, the dual cross vision-language fusion module (DCVLF) enables symmetric exchange between visual tokens and word-level language features, allowing visual representations to attend to relevant linguistic cues while refining the language representation based on visual evidence. Second, we introduce a hierarchical language-guided upsampling module (LMAFN) that reconstructs high-resolution features while modulating multi-scale decoder activations using sentence-level semantics via Feature-wise Linear Modulation (FiLM) and attention, preserving instruction context during spatial detail recovery. Third, at the prediction stage, we employ a language-conditioned dynamic convolution head (LDCH) that forms instruction-conditioned decoding kernels through an expert-mixture formulation, improving semantic consistency for dense mask and grasp predictions. Finally, a lightweight residual refinement module corrects local artifacts and stabilizes grasp outputs in cluttered scenes. Another advantage of our approach is the training

data: in addition to OCID-VLG [28], we train on the larger Grasp-Anything++ [23] dataset, which covers more object categories and scene variations with diverse natural language descriptions, improving robustness to novel objects, unusual phrasing, and out-of-distribution scenes. As shown in Figure 1, we deploy LGGD on a real robotic system for pick-and-place execution, consisting of a KUKA robot, a wrist-mounted RGB-D camera, the LGGD model, and a control module. In real-world grasping experiments, LGGD achieves reliable real-robot performance.

The key contributions of this work can be summarized as follows:

- **Framework.** We propose LGGD, an end-to-end language-conditioned grasp detection framework that couples referring segmentation with dense grasp prediction in a coarse-to-fine manner.
- **Fusion + Decoding.** We introduce DCVLF for bidirectional word-level fusion, and a language-guided hierarchical upsampling module to recover spatial details under linguistic conditioning.
- **Instruction-specialized heads + Refinement.** We develop a language-conditioned dynamic convolution head (LDCH) for instruction-aware coarse predictions, together with residual refinement modules (RM) and multi-stage deep supervision to improve robustness in clutter.
- **Evaluation.** We validate our approach on OCID-VLG and Grasp-Anything++ and further demonstrate real-robot interactive pick-and-place performance under varying tabletop layouts.

## II. RELATED WORK

### A. Grasp Detection

Grasp detection is a core research area in robotics, focused on enabling robots to recognize and perform object grasps in complex environments [29], [30], [22], [31], [32].

**Model-based Grasp Detection.** Traditionally, grasp detection relies on analytical approaches [5], [6], which require 3D object models. The models are scanned and stored in a grasp pose database. During grasp pose detection, the objects point cloud or 3D model data collected by the camera are used for object recognition and pose estimation. The object information is then matched in the database to retrieve the optimal grasp pose. This method has a relatively complex process involving object recognition, object pose estimation, and grasp pose generation. Its core is to construct a grasp contact model based on point contacts, analyzing contact forces and torques. It is suitable for simple scenarios with only a small number of objects but requires obtaining object models in advance and annotating grasp poses. It has poor generalization ability for unknown objects and does not fully utilize cameras or other sensors to enhance the robot's perception of the environment. In real robotic operation environments, it is often impossible to obtain precise models of all objects that need to be manipulated, and environmental interference can also affect performance. Therefore, model-based methods have limited applications in real-world scenarios.

**Discriminative Grasp Detection.** With the emergence of deep learning, grasp detection witnessed substantial performance improvements. Lenz *et al.* [7] propose an early two-stage cascade, in which the first network rapidly filters out grasp candidates, while the second network conducts a more refined evaluation based on RGB, depth, and surface-normal features. Following this paradigm, Mahler *et al.* [8] introduce the Grasp Quality Convolutional Neural Network (GQCNN), which employs a Region Proposal Network (RPN) to sample grasp-relevant regions and subsequently evaluates them using a learned grasp-quality metric to determine the optimal grasp pose. To overcome the inherent latency of multi-stage approaches, one-stage methods combine grasp proposal generation and scoring into a single forward pass, offering significantly improved efficiency. Zhang *et al.* [11] modified the Fully Convolutional Network (FCN) [33] by replacing the blue channel of the RGB image with depth information, forming an RG-D input that consistently outperforms its RGB counterpart. Kumra *et al.* [12] further demonstrated the benefit of multimodal fusion by designing a dual-branch ResNet architecture that independently extracts RGB and depth features before combining them for grasp prediction. Building on the one-stage paradigm, subsequent works introduced additional innovations. Xu *et al.* [34] presented the Single-Shot Grasp (SSG) detector inspired by YOLACT [35], integrating instance segmentation with grasp detection to enable instance-aware grasping. Xu *et al.* [36] eliminated reliance on predefined anchor boxes by proposing an *anchor-free* grasp detector. Extending this framework, Wu *et al.* [37] reformulated angle estimation as a classification problem, whereas Chen *et al.* [38] introduced a continuous angle-encoding scheme that resolves the discontinuity of discrete angle bins and yields more accurate angular predictions.

**Generative Grasp Detection.** Formulating robotic grasping as a dense prediction problem was first demonstrated by Morrison *et al.* [13] with the Generative Grasping Convolutional Neural Network (GG-CNN). Instead of generating and evaluating discrete grasp hypotheses, GG-CNN directly predicts a grasp pose at every pixel location, analogous to semantic segmentation. Owing to its compact and fully convolutional design, GG-CNN enables near-real-time inference and supports dynamic manipulation scenarios. A major progression in generative grasping followed with the Generative Residual Convolutional Neural Network (GR-ConvNet) proposed by Kumra *et al.* [14]. GR-ConvNet incorporates residual learning and RGBD fusion, compressing feature maps via five residual encoder blocks followed by symmetric transposed-convolution decoding. From the final decoder layer, the network produces four dense grasp maps that represent, respectively, the grasp quality, the orientation angle $\cos 2\theta$ and $\sin 2\theta$, and the required gripper width. This parameterization enables accurate and continuous grasp prediction while maintaining real-time execution speed. To further enhance generative models, attention mechanisms have proven highly effective. Yu *et al.* [15] augmented a U-Net with squeeze-and-excitation blocks, enabling adaptive channel-wise feature reweighting and improved robustness in cluttered environments. To expand receptive fields, Cao *et al.* introduced

atrous spatial pyramid pooling (ASPP) [16], and subsequently proposed gaussian-based grasp representation and the Multi-Dimensional Attention Fusion Network (MDAFN) [1], which jointly applies spatial and channel attention to achieve multi-scale feature aggregation with minimal computational overhead. Beyond CNN-based approaches, Wang *et al.* [17] incorporated hierarchical Swin-UNet [39], leveraging windowed self-attention for improved global context modelling. Zhang *et al.* [18] further advanced this direction by combining a Transformer encoder with a convolutional decoder via a Vision-Mamba architecture, unifying long-range reasoning with localized geometric detail and establishing state-of-the-art (SOTA) performance. Despite these advances, existing generative grasping frameworks remain purely vision-based. They do not interpret or reason over linguistic commands, and thus cannot incorporate task-level semantic intent, such as grasping *"the red mug by the handle"* or *"avoiding the sharp metal blade"*. Without explicit language grounding, models may detect feasible grasps but lack awareness of which object should be grasped and in what manner, limiting their utility in human-interactive and instruction-driven robotic systems.

### B. Language-Guided Grasp Detection

Recent progress in robotic grasping has started to overcome the limitation of analytical and deep learning-based approaches through the incorporation of language understanding, allowing robots to execute grasps based on natural language guidance [40], [41], [42], [43], [23], [25], [24], [44].

Early language-driven systems convert the instruction into a referring object mask and subsequently run a conventional grasp detector on the masked crop. Shridhar *et al.* proposed CLIPort [22], which uses CLIP embeddings to align the pixelwise feature map with the textual query, producing a dense relevance heatmap that guides a Transporter network for planar pick-and-place manipulation. While CLIPort handles compositional commands, its two-stage structure, grounding followed by grasp prediction, means that grounding errors directly degrade the final grasp. To solve the problem of accumulated errors in the two-stage method, Liu *et al.* [45] introduce a hierarchical fusion network that integrates sentence-level, noun-level, and spatial-phrase cues with global, object-level, and spatial visual features extracted by CLIP. Recent researches focus on the fusion of visual and language modalities. Chen et al. [46] used residual networks (ResNet) [47] to process visual information and long short-term memory networks (LSTM) [48] to process text input. However, due to the limitation of internal feature alignment, this method has limitations in processing complex instructions. The emergence of vision-language models (VLMs) provides new technical options for this task. Pre-trained backbones such as BERT [49] and CLIP [19] provide strong language and cross-modal representations. BERT captures contextual token dependencies through masked language modeling, whereas CLIP aligns image and text embeddings via a contrastive objective learned from 400 million image-text pairs.

More recently, GraspMamba [50] introduces a Mamba-based language-driven grasp detection framework with hierarchical feature fusion, aiming to efficiently integrate text and

Fig. 2: Rectangular grasp representation.

multi-scale visual information to improve both accuracy and inference speed. While this linear-time architecture improves long-range dependency modeling and inference efficiency over diffusion- or transformer-based methods, its multimodal fusion mechanism still struggles to preserve fine-grained semantics through the full perception pipeline, making it difficult to accurately localize graspable regions when instructions are complex or object parts are visually similar. Meanwhile, Vo *et al.* [26] propose a mask-guided attention approach to inject language cues into the visual stream by restricting attention to segmentation-derived regions. Although this strategy enhances grounding within localized masks, the fusion is typically applied only at a single stage and relies on pre-extracted masks, limiting semantic propagation and weakening intent consistency. Similarly, segmentation-guided approaches such as GraspSAM [51] convert language or clicking inputs into bounding-box prompts and inject them only at the decoder stage. As a result, all these fusion designs exhibit insufficient bidirectional interaction between modalities, resulting in suboptimal interpretation of fine-grained grasp intentions, particularly in cluttered or ambiguous scenes.

To better align grasp prediction with human intent, our method goes beyond simply using CLIP encoders. We embed language instructions throughout the whole perception pipeline by introducing a hierarchical cross-modal fusion mechanism, enabling bidirectional interaction between image and text tokens. This design preserves fine-grained semantic cues and ensures that visual features are progressively refined by linguistic intent at multiple scales. As a result, the detected grasp poses maintain strong semantic consistency with the user-specified command.

## III. GRASP REPRESENTATION

Given an RGB image and a text prompt specifying the target object, our goal is to detect a grasp pose in the image that corresponds to the provided description. To represent grasps, we adopt the widely used rectangle-based grasp formulation employed in prior works [29], [30], [1], as illustrated in Figure 2. For an $N$-channel input image $\mathbf{I} \in \mathbb{R}^{H \times W \times N}$, a planar grasp is defined by a four-tuple of parameters:

$$g = \{(x, y), \theta, w, q\}. \tag{1}$$

where $(x, y)$ denotes the grasp center in image coordinates, $\theta$ is the rotation angle of the end-effector about the camera $z$-axis, $w$ is the gripper opening width, and $q \in [0, 1]$ represents the estimated probability of grasp quality.

## IV. METHOD

As illustrated in Figure 3, our language-guided grasp detection framework with coarse-to-fine learning (LGGD) consists of several key components: pretrained CLIP-based image and text encoders, a dual cross vision-language fusion (DCVLF) bottleneck, hierarchical language-guided upsampling modules, a coarse mask and grasp prediction head, and subsequent mask and grasp refinement modules.

### A. CLIP-based Encoders

Contrastive Language-Image Pre-training (CLIP) [19] represents a significant advancement in vision-language models that learn visual representations through natural language supervision rather than traditional supervised learning on manually annotated datasets. The CLIP model consists of two encoders: an image encoder and a text encoder that map visual and textual inputs into a shared embedding space. The image encoder processes RGB images through a deep network (commonly ResNet [47] or Vision Transformer [52] architectures), while the text encoder processes natural language descriptions through transformer-based architectures. Both encoders produce normalized feature vectors that enable direct comparison through cosine similarity. The training process utilizes contrastive learning on large-scale image-text pairs. Given a batch of $N$ image-text pairs, CLIP computes embeddings for all images $\{I_1, I_2, \ldots, I_N\}$ and texts $\{T_1, T_2, \ldots, T_N\}$. The image encoder generates normalized embeddings $i_k$ for each image $I_k$, while the text encoder produces normalized embeddings $t_k$ for each text $T_k$.

**Image Encoder.** CLIP has established robust cross-modal feature correlations through large-scale pre-training on approximately 400 million image-text pairs. These correlations persist even when the model parameters are frozen during downstream task training. We employ a CLIP-based image encoder to utilize the pre-established semantic consistency between its visual and textual representations, preparing for subsequent feature fusion. Figure 4 illustrates the structure of the CLIP-based ResNet-50 image encoder. Given a resized RGB input image $I \in \mathbb{R}^{224 \times 224 \times 3}$, the CLIP-based ResNet-50 encoder processes it through a hierarchical sequence of convolutional blocks to extract features. Initially, the image passes through the Stem block, followed by batch normalization and ReLU activations. This process reduces the spatial dimensions by half and outputs an initial feature map $x_{\text{stem}} \in \mathbb{R}^{112 \times 112 \times 64}$.

The output is then processed by four ResNet stages, each built from standard bottleneck blocks but preceded by an anti-alias (AA) down-sampling unit that performs average pooling $(2 \times 2, s = 2)$ followed by a stride-1 convolution. The whole process can be formulated as follows:
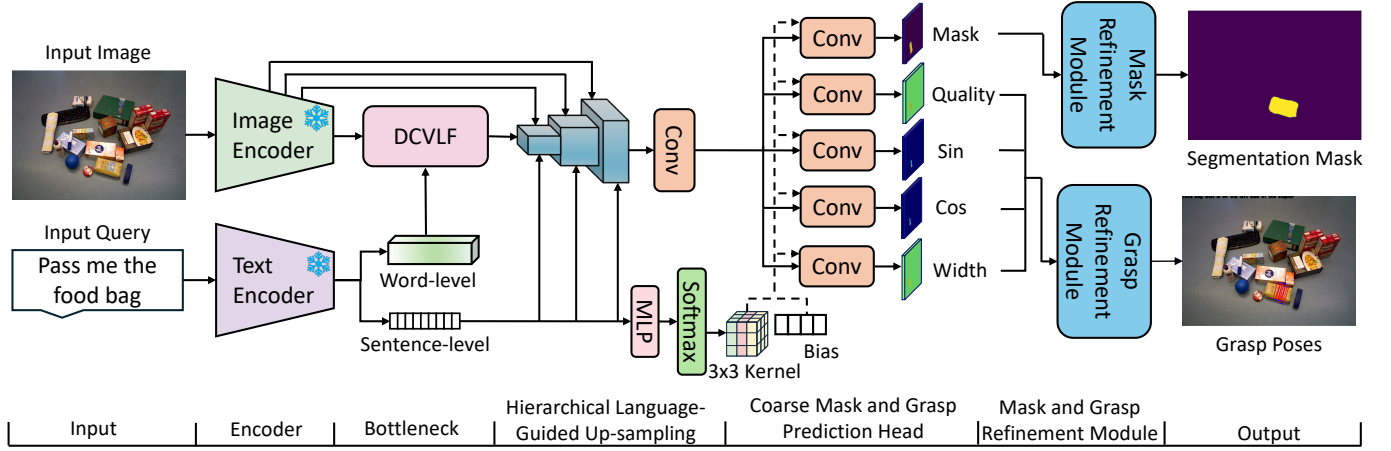
Fig. 3: An overview of our proposed LGGD framework. Given an RGB image and a natural-language command, a CLIP-based image encoder and text encoder extract visual features and word/sentence embeddings. The Dual Cross Vision-Language Fusion (DCVLF) bottleneck aligns the two modalities, after which hierarchical language-guided upsampling progressively refines spatial details according to the textual intent. A coarse mask and grasp prediction head outputs the segmentation mask, grasp quality, angle, and gripper width. Finally, the mask refinement and grasp refinement modules sharpen boundaries and stabilize grasp poses, producing accurate, instruction-consistent grasp poses.



Fig. 4: Structure of Image Encoder. The CLIP-based ResNet-50 backbone extracts progressively abstracted feature maps across four residual stages, while attention pooling enhances global semantic perception.

$$x_0 = f_{\text{stage1}}(x_{\text{stem}}), \qquad x_0 \in \mathbb{R}^{56 \times 56 \times 256}, \tag{2}$$
$$x_1 = f_{\text{stage2}}(x_0), \qquad x_1 \in \mathbb{R}^{28 \times 28 \times 512}, \tag{3}$$
$$x_2 = f_{\text{stage3}}(x_1), \qquad x_2 \in \mathbb{R}^{14 \times 14 \times 1024}, \tag{4}$$
$$x_3 = f_{\text{stage4}}(x_2), \qquad x_3 \in \mathbb{R}^{7 \times 7 \times 2048}, \tag{5}$$
$$x_4 = f_{\text{attn}}(x_3), \qquad x_4 \in \mathbb{R}^{7 \times 7 \times 2048}. \tag{6}$$

To enhance global spatial context modeling, we apply an attention pooling module to the stage-4 feature map $x_3$ and obtain an attention-enhanced representation $x_4$. The module leverages multi-head self-attention to aggregate information across the $7 \times 7$ spatial grid, enabling long-range interactions and improving the expressiveness of high-level visual features for subsequent multimodal fusion. In our design, the intermediate feature maps $\{x_0, x_1, x_2\}$ are retained as shallow features for skip connections in the hierarchical upsampling decoder, while $x_4$ is used as the visual input to the DCVLF bottleneck for vision-language interaction.

**Text Encoder.** We employ the language branch of CLIP [19]

as our text encoder. It takes natural language prompts as input and first tokenizes them into 77 tokens using a Byte Pair Encoding (BPE) vocabulary ($\sim$49k tokens). Special tokens `<|startoftext|>` and `<|endoftext|>` are inserted at the beginning and end, respectively. Token IDs are embedded by adding word and learnable positional embeddings, which are then processed by a 12-layer Transformer encoder with a hidden dimension of 512. The encoder produces two different granularities of text representations: $y_{\text{word}} \in \mathbb{R}^{77 \times 512}$ and $y_{\text{sentence}} \in \mathbb{R}^{512}$. $y_{\text{word}}$ corresponds to word-level contextual features and $y_{\text{sentence}}$ is derived from the final `<|endoftext|>` token followed by projection and $\ell_2$-normalization. We employ the two types of features for different roles within our multi-modal design. Word-level features preserve fine-grained semantics such as spatial descriptions or object attributes (e.g. "left," "red," "cup"). They are used in the DCVLF module, where visual features serve as queries and can dynamically attend to the most relevant tokens to achieve pixel-level grounding through bidirectional cross-attention. The sentence-level features provide a global

Fig. 5: Structure of the DCVLF. Visual features are enriched with positional encoding and refined by a self-attention block. Bidirectional cross-attention then allows image features to attend to language cues and textual features to attend to visual regions, achieving full semantic alignment. Residual connections and FFNs stabilize learning, while a subsequent 1×1 Conv and global MHSA-FFN block further integrate local and global context, yielding robust multimodal representations for grasp reasoning.

semantic representation of the entire instruction. We use it as a high-level control signal in later up-sampling and coarse mask and grasp prediction head stages. During up-sampling, we use sentence-level features to modulate visual feature channels via Feature-wise Linear Modulation (FiLM) [53] to enforce task constraints (e.g., emphasizing the "cup" while suppressing irrelevant regions). In prediction head, we use sentence-level features to generate customized convolutional kernels, enabling task-specific behavior depending on the instruction. Through this design, our model can both align visual features with local linguistic concepts and adjust global processing according to the whole intent, leading to precise and instruction-consistent grasp generation.

### B. Dual Cross Vision-language Fusion Bottleneck

To effectively combine visual and language features for grasp pose prediction, we propose a dual cross vision-language fusion (DCVLF) mechanism that integrates image features with word-level features. Specifically, the DCVLF integrates visual and textual representations through multi-head attention mechanisms [52] for cross-modal fusion. It enables the model to accurately identify and localize target regions in images according to natural language instructions. Image features extracted from the visual branch are combined with word-level embeddings. The attention mechanism [52], widely adopted in cross-modal fusion, allows models to selectively focus on salient elements within a sequence to capture intrinsic dependencies. Given an input sequence $X \in \mathbb{R}^{n \times d}$ (with sequence length $n$ and feature dimension $d$), attention first projects the sequence into queries $Q$, keys $K$, and values $V$:

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V. \tag{7}$$

where $W_Q, W_K, W_V$ are learnable projection matrices. Typically, $d_q = d_k = d_v = d$. The attention operation is then defined as:

$$\text{Attention}(Q, K, V) = \text{SoftMax}\left(\frac{QK^\top}{\sqrt{d}}\right) V. \tag{8}$$

where $QK^\top$ forms an $n \times n$ similarity matrix indicating inter-element relevance. The scaling factor $1/\sqrt{d}$ stabilizes gradients, while the row-wise SoftMax converts scores into probability weights. The resulting weighted sum of $V$ emphasizes contextually important elements, allowing explicit modeling of long-range dependencies.

As illustrated in Figure 5, our DCVLF deeply integrates image and text features. First, image features are enriched with sinusoidal positional encoding and processed through a multi-head self-attention (MHSA) block to capture inter-element relevance representations. Because self-attention is permutation-invariant, positional encoding is essential to preserve geometric structure. Embedding spatial coordinates ensures that attention depends on both content and position, maintaining topological sensitivity. This pre-fusion MHSA modules long-range spatial dependencies, globally reweighting convolutional features and reducing intra-modal redundancy. It also enhances the discriminability of visual queries, ensuring that subsequent cross-modal attention focuses on semantic alignment rather than spatial inference. For text, we directly use the word-level outputs from the CLIP encoder, which already incorporate
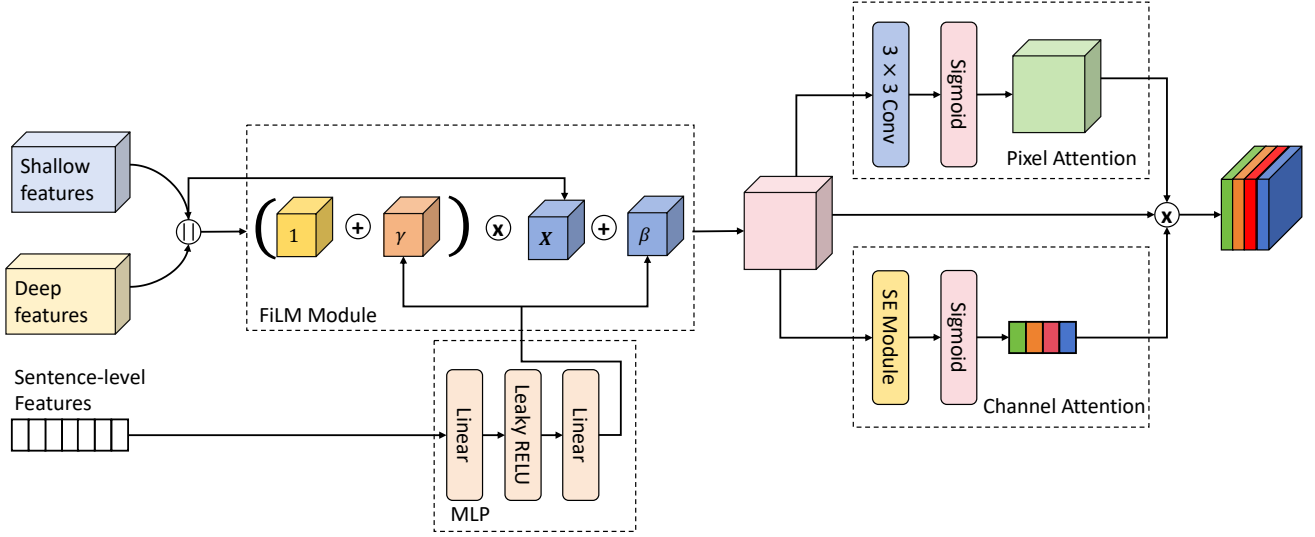
Fig. 6: Structure of the proposed LMAFN. The module integrates sentence-level linguistic context into the decoding pathway via FiLM-based feature modulation and dual attention refinement. Low-level spatial features and high-level semantic features are hierarchically fused, while pixel-wise and channel-wise attentions selectively enhance regions relevant to the instruction. This design ensures progressive semantic alignment and accurate reconstruction of fine-grained structural details during up-sampling.

positional and contextual encoding. The DCVLF then fuses visual and textual information symmetrically:

$$\tilde{I} = \text{SoftMax}\left(\frac{(\hat{I}W_Q)(TW_K)^\top}{\sqrt{d_k}}\right)TW_V, \tag{9}$$

$$\tilde{T} = \text{SoftMax}\left(\frac{(TW_Q)(\hat{I}W_K)^\top}{\sqrt{d_k}}\right)\hat{I}W_V. \tag{10}$$

In the first path, visual features attend to language cues, enriching visual semantics. In the second path, textual features attend to visual regions, enhancing visual grounding. This bidirectional interaction ensures complete multimodal alignment and mitigates information bias inherent in unidirectional fusion. Each path is followed by residual connections, layer normalization, and an independent feed-forward network (FFN) to enhance nonlinearity and training stability. The two enhanced features are concatenated and refined via a $1 \times 1$ convolution and ReLU activation:

$$R = \text{ReLU}(\text{Conv}_{1\times1}(\text{Concat}(\tilde{I}, \tilde{T}))). \tag{11}$$

Finally, MHSA and FFN blocks are used to improve fine-grained representation learning:

$$H = \text{LN}(R + \text{MHSA}(R)), \quad O = \text{LN}(H + \text{FFN}(H)). \tag{12}$$

### C. Hierarchical Language-guided Up-sampling

The decoding process reconstructs full-resolution outputs from fused multimodal features. We integrate language-aware conditioning and multidimensional attention mechanisms into the up-sampling process. Specifically, Feature-wise Linear Modulation (FiLM) injects linguistic context into the visual

domain by adjusting feature activations according to sentence-level semantics. Subsequently, a multidimensional attention module with channel and spatial branches re-weights the feature maps, enhancing regions relevant to the instruction while suppressing background noise. The up-sampling is progressively performed across multiple scales, where language and vision features are fused at each scale. This hierarchical fusion ensures coherent alignment between textual intent and visual details.

**Feature-wise Linear Modulation (FiLM).** FiLM [53] introduces lightweight conditioning by dynamically modulating feature channels with linguistic information. The sentence-level text embeddings $s \in \mathbb{R}^{d_{\text{text}}}$ are transformed via a two-layer MLP to generate modulation coefficients:

$$[\gamma, \beta] = W_2 \,\text{LReLU}(W_1 s). \tag{13}$$

where $\gamma, \beta \in \mathbb{R}^{C \times 1 \times 1}$ and LReLU denotes the LeakyReLU activation. The feature maps $x \in \mathbb{R}^{C' \times H' \times W'}$ are modulated as:

$$\text{FiLM}(x, s) = (1 + \gamma) \odot x + \beta. \tag{14}$$

where $\odot$ denotes channel-wise multiplication. The additive offset in $(1 + \gamma)$ preserves information when $\gamma = 0$, while $\beta$ injects language-dependent biases. Unlike spatial concatenation methods, FiLM maintains the tensor structure, providing efficient channel-wise semantic gating that aligns global sentence semantics with visual activations.

As shown in Figure 6, shallow features $x_s \in \mathbb{R}^{C' \times H' \times W'}$ preserve spatial structures, while deep features $x_d \in \mathbb{R}^{C' \times H' \times W'}$ encode high-level semantics. They are concatenated:

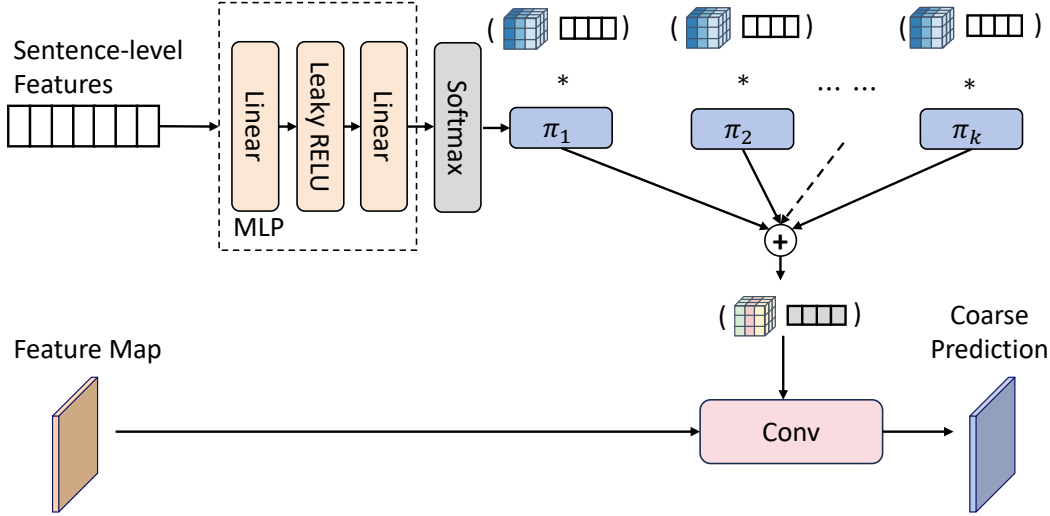$$x_{\text{cat}} = \text{Concat}(x_d, x_s) \in \mathbb{R}^{2C' \times H' \times W'}. \tag{15}$$

Fig. 7: Structure of the Language-conditioned Dynamic Convolution Head(LDCH). The module adopts a kernel-mixture formulation in which sentence-level features guide the weighted combination of multiple convolution experts to generate language-conditioned coarse predictions.

Then, FiLM, conditioned by sentence-level embeddings, refines $x_{\text{cat}}$ to emphasize task-relevant semantics.

$$x_{\text{FiLM}} = \text{FiLM}(x_{\text{cat}}, s). \tag{16}$$

**Language-guided Multidimensional Attention Fusion Network (LMAFN).** Our LMAFN extends channel-wise recalibration [54] with pixel-level and linguistic guidance. Specifically, a $3\times3$ convolution followed by sigmoid activation generates pixel-level attention maps:

$$\mathbf{A}_{\text{pixel}} = \sigma(\text{Conv}_{3\times3}(x_{\text{FiLM}})) \in \mathbb{R}^{C \times H \times W}. \tag{17}$$

Parallel to pixel-level attention, the FiLM-enhanced feature $x_{\text{FiLM}}$ is first processed by global average pooling to obtain $z$. Then it passed through a shared two-layer MLP:

$$z = \text{GAP}(x_{\text{FiLM}}), \tag{18}$$

$$A_{\text{channel}} = \sigma(W_2\, \delta(W_1 z)). \tag{19}$$

where $W_1$ and $W_2$ denote MLP weights. Weight sharing is used to ensure consistent semantic transformations for both modalities. Finally, the outputs integrate both attention maps:

$$x_{\text{out}} = x_{\text{FiLM}} \odot A_{\text{pixel}} \odot A_{\text{channel}}. \tag{20}$$

### D. Coarse Mask and Grasp Prediction Head

To specialize output predictions according to textual instructions, we introduce a language-guided dynamic convolution module for the coarse mask and grasp prediction heads. The key idea is to use sentence-level features to predict a mixture over a small set of learnable convolution experts [55], and to form Language-conditioned kernels via a weighted combination, as illustrated in the Figure 7.

**Task-wise feature preparation.** Given the feature map $x_{\text{out}}$, we first expand its channel dimension using a $1\times1$ convolution:

$$\mathbf{F} = \text{Conv}_{1\times1}(x_{\text{out}}) \in \mathbb{R}^{B \times 5C \times H \times W}. \tag{21}$$

The output is then partitioned into five task-specific feature groups through a fixed channel-wise slicing operation:

$$\mathbf{F} = [F_{\text{mask}}; F_{\text{quality}}; F_{\text{sin}}; F_{\text{cos}}; F_{\text{width}}], \quad F_t \in \mathbb{R}^{B \times C \times H \times W}, \tag{22}$$

corresponding to the segmentation mask, grasp quality, grasp angle (sine and cosine), and gripper width prediction heads. For efficient per-sample dynamic convolution, each $F_t$ is reshaped for grouped convolution:

$$\tilde{F}_t = \text{reshape}(F_t) \in \mathbb{R}^{1 \times (B \cdot C) \times H \times W}. \tag{23}$$

**Language-guided kernel mixing.** For each sample, we extract a sentence-level embedding $s_i \in \mathbb{R}^d$ and predict a distribution over $K$ kernel experts using an MLP followed by a softmax:

$$\pi_i = \text{Softmax}(\text{MLP}(s_i)), \quad \pi_i \in \mathbb{R}^K, \tag{24}$$

where $\text{MLP}(\cdot)$ is a two-layer perceptron with a LeakyReLU nonlinearity. For each prediction head $t$, we maintain a bank of $K$ learnable convolution experts:

$$\{W_t^{(k)}, b_t^{(k)}\}_{k=1}^K, \quad W_t^{(k)} \in \mathbb{R}^{C \times 3 \times 3}, \ b_t^{(k)} \in \mathbb{R}. \tag{25}$$

Given the mixture weights $\pi_i$, the language-conditioned kernel and bias are formed via a weighted sum:

$$W_{i,t} = \sum_{k=1}^K \pi_{i,k} W_t^{(k)}, \qquad b_{i,t} = \sum_{k=1}^K \pi_{i,k} b_t^{(k)}. \tag{26}$$

Stacking all instances yields the dynamic kernel and bias bank:

$$W_{t,\text{d}} = \{W_{1,t}, \ldots, W_{B,t}\} \in \mathbb{R}^{B \times C \times 3 \times 3},$$
$$b_{t,\text{d}} = [b_{1,t}, \ldots, b_{B,t}] \in \mathbb{R}^B. \tag{27}$$

**Coarse prediction via grouped convolution.** We apply the instruction-conditioned parameters through a grouped convolution with the number of groups equal to the batch size:

$$Y_t = \text{GroupConv2D}(\tilde{F}_t, W_{t,\text{d}}, b_{t,\text{d}}, \text{groups} = B), \tag{28}$$
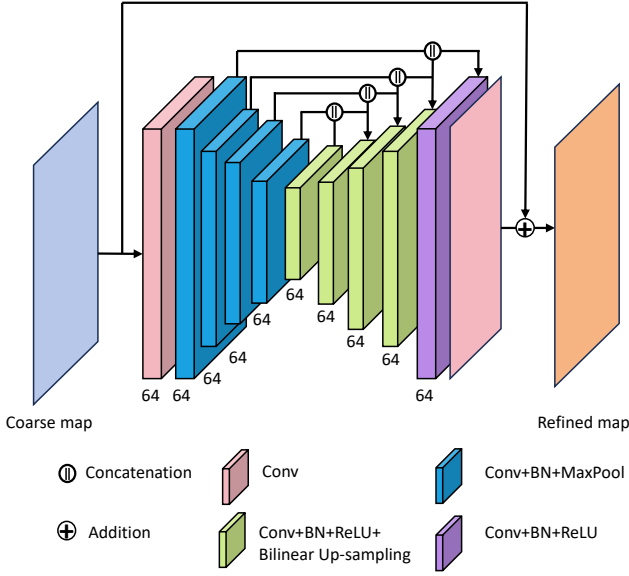
Fig. 8: Structure of the Refinement Module (RM). The RM adopts a lightweight U-Net-style structure to learn residual corrections over coarse predictions.

and reshape the output back to the standard batch format $Y_t \in \mathbb{R}^{B \times 1 \times H \times W}$. The same mechanism is used for all five heads, producing language-conditioned coarse predictions for the mask and grasp parameter maps.

### E. Mask and Grasp Refinement Modules

We employ a lightweight U-Net-style architecture, which learns residual corrections to coarse predictions rather than predicting the final outputs directly. This residual learning design stabilizes training and preserves the information from the preceding prediction stage, as illustrated in Figure 8. Specifically, we employ two refinement branches: the mask refinement module refines the segmentation mask, and the grasp refinement module jointly refines grasp quality, angle (sine and cosine), and gripper width. Given coarse predictions,

$$Y_{\text{mask}}^{\text{coarse}} \in \mathbb{R}^{B \times 1 \times H \times W}, \quad Y_{\text{grasp}}^{\text{coarse}} \in \mathbb{R}^{B \times 4 \times H \times W}, \quad (29)$$

each branch learns a residual correction $\Delta$ that refines the coarse predictions:

$$Y_{\text{mask}} = Y_{\text{mask}}^{\text{coarse}} + \text{Mask}_{\text{refine}}(Y_{\text{mask}}^{\text{coarse}}), \quad (30)$$

$$Y_{\text{grasp}} = Y_{\text{grasp}}^{\text{coarse}} + \text{Grasp}_{\text{refine}}(Y_{\text{grasp}}^{\text{coarse}}). \quad (31)$$

As shown in Figure 8, each refinement module follows a symmetric encoder-decoder structure. The encoder consists of four convolutional blocks with max-pooling for context aggregation, while the decoder entails mirrored convolutional blocks with bilinear up-sampling and skip connections to recover spatial details. The outputs are one-channel (mask) or four-channel (grasp) residual maps aligned with the input resolution. All refinement operations are applied on logits before activation, ensuring gradient flow through both coarse and refined prediction paths for stable end-to-end optimization. This coarse-to-fine framework, coarse estimation followed by

residual refinement, improves precision and semantic coherence, especially in cases involving occluded or irregularly shaped objects, producing sharper masks and more stable grasp orientation predictions.

### F. Loss Function

To facilitate the joint optimization of language-guided segmentation and grasp prediction, we design a comprehensive multi-stage, multi-task loss. The objective supervises both the coarse prediction branch and the refinement branch, and couples binary mask segmentation with continuous grasp regression into a unified training framework.

**Multi-stage Deep Supervision.** Our architecture employs a refinement module with residual connections to progressively enhance the grasp predictions. To fully utilize the learning capacity of intermediate feature stages and stabilize optimization, we adopt deep supervision on both the coarse and refined outputs. The overall training objective is formulated as:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{refine}} + \lambda \cdot \mathcal{L}_{\text{coarse}}. \quad (32)$$

where $\mathcal{L}_{\text{refine}}$ denotes the loss imposed on the final refined predictions, and $\mathcal{L}_{\text{coarse}}$ is an auxiliary loss applied to the coarse predictions. The coefficient $\lambda$ controls the contribution of the coarse supervision. In practice, we set $\lambda = 0.5$ to provide sufficient gradient signals to the coarse stage while preventing it from dominating the optimization of the refinement branch. Each stage-specific loss decomposes into a segmentation term and a grasp regression term:

$$\mathcal{L}_{\text{refine}} = \mathcal{L}_{\text{mask}}^{\text{refine}} + \mathcal{L}_{\text{grasp}}^{\text{refine}}, \quad (33)$$

$$\mathcal{L}_{\text{coarse}} = \mathcal{L}_{\text{mask}}^{\text{coarse}} + \mathcal{L}_{\text{grasp}}^{\text{coarse}}. \quad (34)$$

where $\mathcal{L}_{\text{mask}}^{\text{refine}}$ and $\mathcal{L}_{\text{mask}}^{\text{coarse}}$ supervise the binary segmentation masks at refined and coarse resolutions, respectively, and $\mathcal{L}_{\text{grasp}}^{\text{refine}}$ and $\mathcal{L}_{\text{grasp}}^{\text{coarse}}$ supervise the corresponding dense grasp predictions. This multi-stage deep supervision encourages the network to produce semantically meaningful, language-consistent predictions at early stages, while the refinement module focuses on correcting residual errors and improving prediction accuracy.

**Segmentation Mask Loss.** The segmentation branch predicts a binary mask indicating the language-referred object or region to be grasped. Let $Y_{\text{mask}} \in \mathbb{R}^{H \times W}$ denote the predicted logit map and $M \in \{0,1\}^{H \times W}$ the ground-truth binary mask, where $M_{uv} = 1$ indicates that pixel $(u, v)$ belongs to the target region specified by the language query. We use a weighted binary cross-entropy loss with logits:

$$\mathcal{L}_{\text{mask}} = -\frac{1}{\sum_{u,v} w_{uv}} \sum_{u=1}^{H} \sum_{v=1}^{W} w_{uv} \big[ M_{uv} \log \sigma(Y_{uv}) + (1 - M_{uv}) \log(1 - \sigma(Y_{uv})) \big]. \quad (35)$$

where $\sigma(\cdot)$ is the sigmoid activation that maps logits to probabilities. The normalization by $\sum_{u,v} w_{uv}$ keeps the loss scale consistent across images with different spatial sizes and varying object extents. To alleviate class imbalance between foreground and background pixels, which is particularly severe

when the language-referred object occupies only a small portion of the image, we introduce a simple yet effective foreground reweighting scheme:

$$w_{uv} = \beta \cdot M_{uv} + 1. \tag{36}$$

In this formulation, background pixels receive a weight of $1$, while foreground pixels are weighted by $1 + \beta = 1.5$. This increases the contribution of correctly segmenting the language-targeted region and prevents the model from converging to trivial background-dominant solutions. The hyperparameter $\beta$ can be tuned according to the foreground-to-background ratio of the dataset; higher values impose stronger penalties on misclassified foreground pixels. Note that the segmentation mask serves a dual purpose: it not only evaluates how well the model segments the language-specified object, but also defines the set of valid pixels used for supervising grasp parameters, thereby tightly coupling spatial localization and grasp prediction.

**Grasp Regression Loss.** The grasp head predicts grasp parameters for each spatial location, conditioned jointly on visual features and language-guided fusion features. While a planar grasp can be conceptually defined as $g = \{(x, y), \theta, w, q\}$, in our implementation we adopt a learning-friendly parameterization:

$$G = \{q, \cos(2\theta), \sin(2\theta), w\} \in \mathbb{R}^4. \tag{37}$$

where $q \in [0, 1]$ denotes the grasp quality score, $\theta \in [-\pi/2, \pi/2]$ is the in-plane grasp angle, and $w > 0$ represents the gripper width required to execute the grasp. The angle is encoded by its double-angle trigonometric form $(\cos(2\theta), \sin(2\theta))$ to avoid discontinuities at $\theta = \pm\pi/2$ and to ensure smooth optimization over angular space. Let $\hat{G} \in \mathbb{R}^{H \times W \times 4}$ and $G \in \mathbb{R}^{H \times W \times 4}$ denote the predicted and ground-truth grasp tensors, respectively. We compute the grasp regression loss only over valid grasp locations, which are defined by the ground-truth mask:

$$\mathcal{V} = \{i \mid M_i = 1\}. \tag{38}$$

where $i$ indexes spatial locations in the image. The loss is given by:

$$\mathcal{L}_{\text{grasp}} = \frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} \sum_{m \in \mathcal{M}} \text{SmoothL1}\left(\hat{G}_i^m - G_i^m\right). \tag{39}$$

where $\mathcal{M} = \{q, \cos(2\theta), \sin(2\theta), w\}$ indexes the grasp components. This masking strategy ensures that supervision focuses on language-relevant object regions and avoids penalizing predictions on irrelevant background areas. We use the Smooth L1 loss due to its robustness to outliers and stable gradients:

$$\text{SmoothL1}(x) = \begin{cases} \frac{1}{2}(\sigma x)^2, & \text{if } |x| < \frac{1}{\sigma^2} \\ |x| - \frac{1}{2\sigma^2}, & \text{otherwise.} \end{cases} \tag{40}$$

where $\sigma$ controls the transition point between the quadratic and linear regimes. We set $\sigma = 1$, which results in a quadratic penalty for small errors ($|x| < 1$), encouraging precise grasp predictions, and a linear penalty for larger deviations, limiting

the influence of noisy or ambiguous annotations. This is particularly important in language-guided grasping, where fine-grained object parts (e.g., handle vs. body) may be annotated with some degree of uncertainty.

Overall, the proposed loss formulation tightly couples segmentation and grasp prediction under language guidance: the mask loss drives the network to localize the language-referred object, while the regression loss encourages accurate and semantically consistent grasp parameters within that region. The multi-stage supervision further ensures that both coarse and refined predictions benefit from strong training signals, leading to robust and stable convergence.

## V. EXPERIMENTAL RESULTS

### A. Dataset

To evaluate our model's performance, we use the OCID-VLG [28] and Grasp-Anything++ [23] datasets, which are the most comprehensive benchmarks available for language-guided grasp detection.

**OCID-VLG.** OCID-VLG is constructed on top of OCID [56], [29], extending it to support language-driven grasping tasks. It is derived from 1,763 unique OCID scenes and comprises nearly 90k fully annotated tuples. Each tuple includes an RGB image, a natural language instruction specifying the target object, a pixel-level segmentation mask of the target, and a set of valid grasp rectangles. To assess our models accuracy, generalization, and zero-shot capabilities, we conduct experiments using the following three dataset splits:

- **Multiple-Split.** This is a standard random scene split in which both the object classes and instances appearing in the test set are also present during training. We use this split to evaluate the models performance under conventional conditions.
- **Novel-Instances Split.** In this split, the test set includes object instances that are unseen during training, although their corresponding classes are observed. This setup is used to evaluate the models instance-level generalization-its ability to transfer knowledge from seen instances to novel ones within the same class.
- **Novel-Classes Split.** This is the most challenging split, where the test set contains object classes that never appear during training. We use this split to evaluate the models zero-shot capability-its ability to reason about and localize novel classes based solely on language descriptions and prior visual knowledge, without having seen any corresponding examples.

**Grasp-Anything++.** We further train and evaluate our model on the synthetic Grasp-Anything++ dataset [23], an extension of the original Grasp-Anything dataset [31]. Grasp-Anything is generated using foundation models and contains 1 million samples featuring over 3 million diverse objects, along with 10 million grasping instructions paired with corresponding ground truth. Experiments on Grasp-Anything++ enable us to assess whether the proposed architecture can perform effectively in large-scale synthetic environments with procedurally generated scenes and full-sentence prompts.

## B. Evaluation Metrics

To evaluate our model's performance across different tasks, we adopt a set of established metrics that measure its effectiveness in two domains: referring object segmentation and referring grasp detection.

**Grasp Detection Metrics.** For referring grasp detection, where the goal is to identify valid grasp candidates in an image, we adopt the rectangle-based evaluation protocol used in prior works [14], [16], [28]. A predicted grasp is considered correct if it satisfies both of the following conditions:

- **Angle Difference:** The absolute difference between the predicted grasp rectangles orientation and that of the ground-truth rectangle is less than $30°$.
- **Jaccard Index:** The Intersection over Union (IoU) between the predicted grasp rectangle $g_p$ and a ground-truth rectangle $g_t$ exceeds 0.25. The IoU, also known as the Jaccard index, is defined as

$$J(g_p, g_t) = \frac{|g_p \cap g_t|}{|g_p \cup g_t|} > 0.25. \tag{41}$$

To evaluate the models ability to predict a valid grasp corresponding to a given referring expression, we report the $J@N$ metric. Following the protocol in [28], a prediction is considered successful if the top-N ranked predicted grasp:

- Refers to the correct object as specified by the language instruction.
- Satisfies both the angle criterion (orientation difference less than $30°$) and the IoU (Jaccard index) criterion (greater than 0.25) with respect to the ground-truth grasp on the referred object.

These evaluations capture the models ability to integrate visual grounding with grasp estimation in language-conditioned tasks.

**Referring Segmentation Metrics.** To assess the models ability to segment the object specified by a textual instruction, we follow the evaluation protocol in [57], [28]. Two metrics are employed:

- **Intersection over Union(IoU):** The average IoU between the predicted and ground-truth segmentation masks.
- **Precision@X (Pr@X):** The percentage of test instances for which the IoU exceeds a threshold $X$. Results are reported for $X \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$.

## C. Training Setup

Our experiments are implemented in PyTorch 2.4.0 [58] with CUDA 12.4 support. All models are trained end-to-end on a single NVIDIA H100 GPU with 80 GB memory. We use the AdamW optimizer [59] with a base learning rate of $1 \times 10^{-4}$ and a weight decay of 0.06. The learning rate follows a cosine annealing schedule [60] with a linear warm-up during the first 2,000 iterations [47], starting from $1 \times 10^{-5}$ to improve early-stage training stability. Models are trained for 50 epochs with a batch size of 72.

## D. Results on OCID-VLG

**Multi-Split.** The evaluation on the OCID-VLG dataset under the Multi-Split setting provides a comprehensive benchmark for model performance, simulating a scenario in which both object classes and instances in the test set are seen during training. As shown in Table I, our proposed method outperforms the recent method, CROG [28], achieving an IoU of **83.1%** for segmentation and a J@1 of **85.4%** for top-1 grasp prediction, compared to CROG's 81.1% IoU and 77.2% J@1. These results validate the effectiveness of our approach.

**Novel-Instances Split.** To assess the models ability to generalize to new instances of object classes seen during training, we adopt the Novel-Instances Split. Table II reports the performance of each model under this setting. The recent method, CROG [28], exhibits a substantial performance decline, whereas our model demonstrates stronger generalization capabilities. This degradation in CROG can be attributed to its strategy of fine-tuning the entire CLIP [19] backbone on a task-specific dataset rather than freezing it. While such full fine-tuning may improve in-domain performance, it increases the risk of overfitting and can induce catastrophic forgetting of the broad visual-semantic knowledge acquired during large-scale pre-training. As a result, CROG struggles when encountering novel instances, leading to the observed performance drop.

**Novel-Classes Split.** The Novel-Classes Split provides a rigorous evaluation of a model's language understanding and reasoning capabilities, as it requires handling object classes that are entirely unseen during training. The experimental results are presented in Table III. CROG exhibits a substantial decline in performance under this setting, with its J@1 and IoU dropping sharply. As discussed earlier, fully fine-tuning the CLIP [19] backbone tends to contaminate the open-world knowledge acquired during pre-training with task-specific biases, thereby severely weakening the models zero-shot generalization ability. In contrast, our proposed method achieves a J@1 of 46.0% and an IoU of 63.1%, demonstrating strong robustness when confronted with novel object classes.

## E. Results on Grasp-Anything++

To assess whether our approach is effective on synthetic data, we further train and evaluate our model on the Grasp-Anything++ dataset using full-sentence prompts. This dataset features procedurally generated tabletop scenes with diverse object configurations and natural language instructions, and it provides separate splits for scenes containing seen and unseen objects. The quantitative results are summarized in Table IV. Our method outperforms prior approaches on both splits. Compared to the second-best method, LGD [23], our performance improves to 0.59 and 0.31 on the seen and unseen splits, respectively. These results demonstrate that our language-guided grasping framework generalizes well not only to real-world datasets but also to large-scale synthetic environments with full-sentence prompts.

TABLE I: Evaluation results of different methods on the OCID-VLG dataset with the Multi-Split setting. Bold and underlined numbers denote the best and second-best results, respectively.

| Model | J@1 | J@5 | IoU | Pr@50 | Pr@60 | Pr@70 | Pr@80 | Pr@90 |
|---|---|---|---|---|---|---|---|---|
| GT-Grounding | 28.7 | - | - | - | - | - | - | - |
| GT-Masks + CLIP [19] | 11.9 | - | 35.0 | 35.0 | 35.0 | 35.0 | 35.0 | 35.0 |
| GR-ConvNet [14] + CLIP [19] | 9.7 | 15.4 | 31.3 | 21.0 | 11.6 | 5.5 | 2.4 | 0.5 |
| SAM [61] + CLIP [19] | 7.2 | - | 25.7 | 29.3 | 28.5 | 27.4 | 22.7 | 9.1 |
| GLIP [62] + SAM [61] | 10.7 | - | 30.3 | 34.7 | 34.1 | 33.5 | 28.6 | 11.7 |
| Det-Seg [29] + CLIP [19] | 28.1 | - | 29.0 | 27.2 | 20.9 | 17.5 | 17.2 | 16.0 |
| SSG [34] + CLIP [19] | 33.5 | - | 33.6 | 35.6 | 35.6 | 35.5 | 35.5 | **32.8** |
| CLIP [19] + MHSA [52] | 82.5 | 87.7 | 71.3 | 93.2 | 84.4 | 54.5 | 27.7 | 0.1 |
| EfficientGrasp [1] + CLIP [19] | 79.2 | 82.4 | 80.2 | 87.4 | 71.8 | 43.1 | 32.1 | 15.6 |
| CROG [28] | 77.2 | - | 81.1 | 96.9 | **94.8** | 87.2 | 64.1 | 16.4 |
| Ours | **85.4** | **90.2** | **83.1** | **98.0** | <u>93.7</u> | **91.3** | **66.7** | <u>26.1</u> |

TABLE II: Evaluation results of different methods on the OCID-VLG dataset with Novel-Instances Split setting.

| Model | J@1 | J@5 | IoU | Pr@50 | Pr@70 | Pr@90 |
|---|---|---|---|---|---|---|
| CROG [28] | 48.3 | 52.0 | 61.3 | 74.8 | 52.2 | 1.7 |
| Ours | **57.6** | **64.1** | **67.6** | **76.4** | **61.9** | **17.1** |

TABLE III: Evaluation results of different methods on the OCID-VLG dataset with Novel-Classes Split setting.

| Model | J@1 | J@5 | IoU | Pr@50 | Pr@70 | Pr@90 |
|---|---|---|---|---|---|---|
| CROG [28] | 14.6 | 15.2 | 40.4 | 46.5 | 23.9 | 0.1 |
| Ours | **46.0** | **54.4** | **63.1** | **72.5** | **54.8** | **19.4** |

TABLE IV: Evaluation results of different methods on the Grasp-Anything++ dataset (full prompt sentence).

| Model | Seen | Unseen |
|---|---|---|
| GR-ConvNet [14] + CLIP [19] | 0.21 | 0.12 |
| GGCNN [13] + CLIP [19] | 0.11 | 0.07 |
| EfficientGrasp [1] + CLIP [19] | 0.33 | 0.18 |
| CLIPort [22] | 0.29 | 0.19 |
| LGD [23] | 0.41 | – |
| Ours | **0.59** | **0.31** |



Fig. 9: Results of reasoning grasping with key words.

### F. Qualitative Evaluations

The complexity of natural language instructions also affects model performance. To analyze this effect, we compare how different models behave under varying instruction types, including reasoning grasping with key words, reasoning grasping with location instructions, reasoning grasping with direct descriptions, and reasoning grasping with novel instances. The visualizations of the CROG and our model under three different instruction types are shown in Figure 9, Figure 10, Figure 11, and Figure 12.

**Reasoning Grasping with Key Words.** For simple instructions containing only object keywords, all models successfully localize the target and generate correct grasp poses, as shown in Figure 9. This demonstrates that both CROG and our model exhibit basic vision-language alignment when the semantic cues are straightforward.

**Reasoning Grasping with Location Instructions.** When the instruction includes spatial location descriptions (as shown in Figure 10), CROG's performance begins to deteriorate, often misidentifying multiple objects as potential grasp targets. In contrast, our model performs reliably across different positional terms. Even when the instruction requires more complex spatial reasoning, our approach consistently interprets and resolves composite spatial relationships correctly.

**Reasoning Grasping with Direct Descriptions.** For instructions that require understanding specific attributes or brand names, our model also outperforms CROG, as illustrated in Figure 11.

**Reasoning Grasping on the Novel Instances.** The proposed model further demonstrates strong zero-shot understanding by correctly identifying objects described using non-obvious brand names. The visualizations in Figure 12 highlight this difference in generalization ability. For the instruction *"Pick the brown kleenex tissues"*, the target refers to a specific tissue box variant not seen during training. CROG fails completely in this scenario, incorrectly localizing its grasp predictions on unrelated objects-an indication of its limited generalization. In contrast, our model is able to correctly identify and approxi-
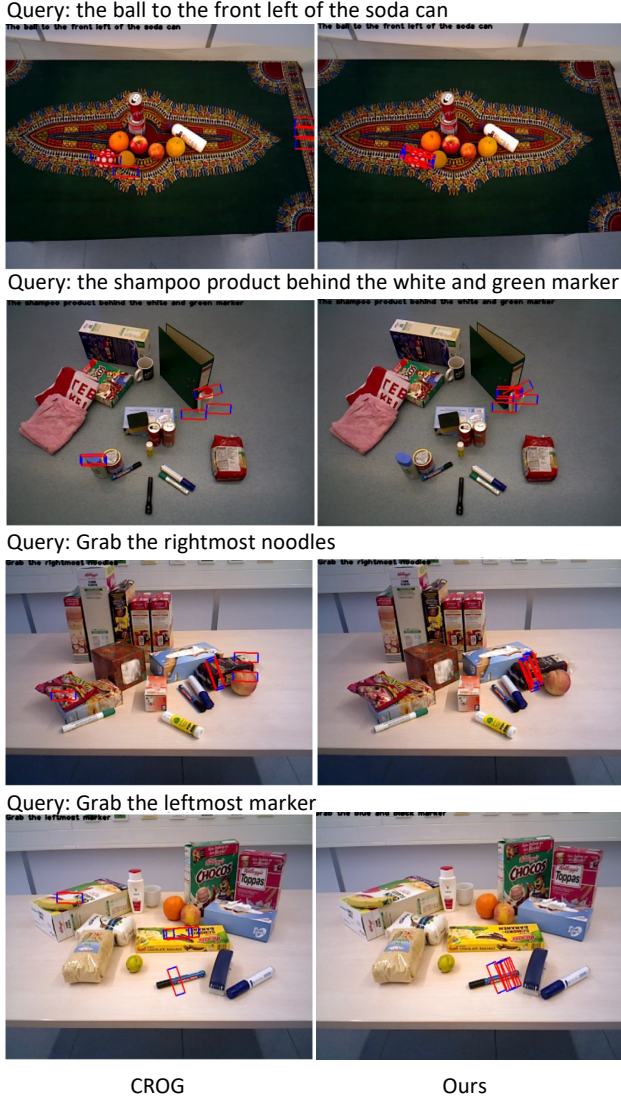
Query: the ball to the front left of the soda can

Query: the shampoo product behind the white and green marker

Query: Grab the rightmost noodles

Query: Grab the leftmost marker

CROG                                            Ours



Fig. 10: Results of reasoning grasping with location instructions.

Query: Pick the Alnatura noodles

Query: Grab the blue and black marker

CROG                                            Ours



Fig. 11: Results of reasoning grasping with direct descriptions.

Query: white food bag product

Query: Pick the brown kleenex tissues

Query: The blue and beige food box

CROG                                            Ours



Fig. 12: Results of reasoning grasping on the novel instances.

mately localize the intended tissue box.

**Attention Map Visualization.** To further illustrate how our model leverages language to guide visual reasoning, we visualize its attention maps for several representative instruction-scene pairs, as shown in Figure 13. The visualizations show that the model consistently attends to task-relevant regions aligned with the textual instructions. For the command *"Pick the lime"* the attention map focuses precisely on the small lime despite numerous distractors. When instructed to *"Pick the keyboard"*, the attention spreads along the keyboards elongated structure, accurately covering its graspable area. For the more fine-grained instruction *"Get the Vichy shampoo bottle"*, the model successfully isolates the target bottle from several visually similar containers.

## G. Efficiency Analysis

In robotic applications, efficiency is as important as accuracy. Table V summarizes the computational costs of the evaluated models, including parameter size, FLOPs, MACs,
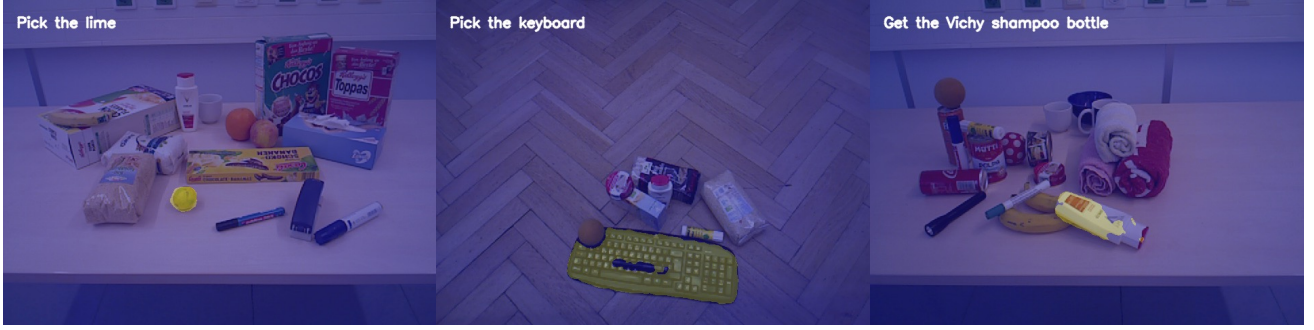
Fig. 13: Attention map visualization of our model.

TABLE V: Comparison of model complexity and efficiency.

| Model | Parameters (M) | FLOPs (G) | MACs (G) | Inference Time (ms) | Trainable Parameters (M) |
|---|---|---|---|---|---|
| GR-ConvNet [14] + MHSA [52] | 55.77 | 111.53 | 64.81 | 21.00 | 52.21 |
| CROG [28] | 147.11 | 223.50 | 111.72 | 39.11 | 87.11 |
| Ours | 171.76 | 289.21 | 144.61 | 43.00 | 67.70 |

TABLE VI: Ablation study on the key components.

| DCVLF | LMAFN | RM | LDCH | IoU | J@1 |
|---|---|---|---|---|---|
| | | | | 31.30 | 9.73 |
| ✓ | | | | 71.34 | 79.36 |
| ✓ | ✓ | | | 80.16 | 82.27 |
| ✓ | ✓ | ✓ | | 82.29 | 84.53 |
| ✓ | ✓ | ✓ | ✓ | **83.14** | **85.36** |

TABLE VII: Results of Real Robot Experiment.

| Setup | Ball | Fruit | Can | Glue | Average |
|---|---|---|---|---|---|
| Isolated | 95% | 85% | 95% | 100% | 93.75% |
| Scattered | 100% | 95% | 95% | 100% | 97.5% |
| Cluttered | 90% | 90% | 85% | 90% | 88.75% |

and single-sample inference time. When considered alongside the accuracy results in Table I, these metrics highlight the trade-off between model performance and computational cost. The results demonstrate that our proposed method strikes an effective balance between efficiency and accuracy, achieving an inference time of 43.00 ms with a parameter budget of 171.76 M, while delivering superior grasping performance.

### H. Ablation Study

Table VI reports a step-by-step ablation where we progressively add each key component to the baseline. The baseline achieves only 31.30% IoU and 9.73% J@1, suggesting that simple feature concatenation fails to capture the fine-grained alignment required by language-guided grasping. After introducing Dual Cross Vision-language Fusion Bottleneck (DCVLF), performance increases to 71.34% IoU and 79.36% J@1 (+40.04/+69.63), indicating that DCVLF is effective establishing vision-language fusion. Building on this, adding Hierarchical Language-guided Up-sampling with Language-guided Multidimensional Attention Fusion Network (LMAFN) further improves results to 80.16% IoU and 82.27% J@1 (+2.57/+0.85), demonstrating that it benefits dense prediction by recovering spatial details and improving boundary fidelity with text guidance. With the refinement module (RM), the model reaches 82.29% IoU and 84.53% J@1 (+2.13/+2.26), showing that refinement further enhances grasp reliability and segmentation consistency. Finally, replacing the static convolution heads with the proposed LDCH results in the best overall

performance: 83.14% IoU and 85.36% J@1 (+0.85/+0.83), validating that instruction-conditioned kernel mixing provides more robust and task-adaptive predictions.

### I. Real-Robot Interactive Grasp Experiments

We evaluate the proposed system in an interactive grasping setting, where a user issues an instruction specifying a target object and the robot must localize the correct target and execute a pick-and-place operation. This experiment is designed to assess end-to-end language-conditioned manipulation performance under varying tabletop layouts. We consider four target objects, Ball, fruit, Can, and Glue, as representative items with different shapes. Our experiments are conducted on a real robotic platform using a KUKA LBR iiwa 14 R820 manipulator equipped with a Robotiq 2F-85 adaptive gripper. Model inference and the control pipeline are executed on a workstation with an NVIDIA RTX 4090 (24GB) GPU and an AMD Ryzen Threadripper Pro 7955WX CPU. To assess sensitivity to object arrangement and occlusion, trials are performed under three tabletop configurations:

1) **Isolated**, where the target is presented without nearby distractors;
2) **Scattered**, where 4-5 objects are distributed across the workspace with clear separation;
3) **Cluttered**, where objects are densely packed, increasing both occlusion and physical interference.

Representative scenes are shown in Figure 14.

Figure 15 further illustrates the procedure of a representative trial. Specifically, the workspace is partitioned into a predefined pick region and a destination (place) region. For each
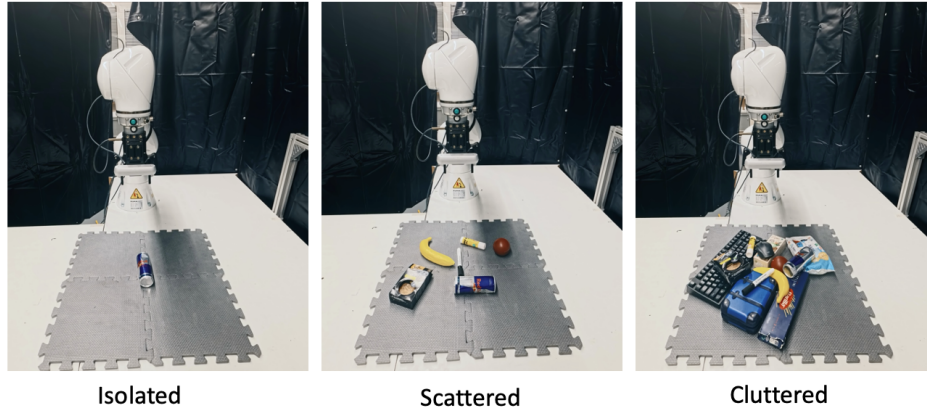
Fig. 14: Interactive real robot experiment in 3 Setups: Isolated (left), Scattered (middle), Cluttered (right).
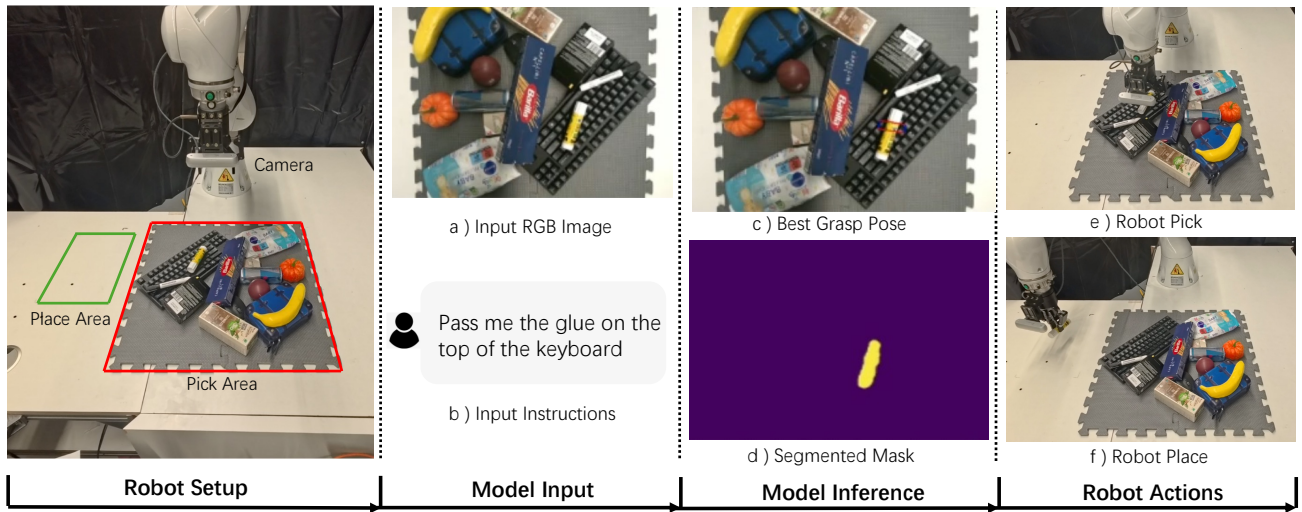


Fig. 15: Overview of a real-robot interactive pick-and-place trial. Left: experimental workspace with a wrist-mounted RGB-D camera and predefined pick and place regions. Given an input RGB image (a) and a language instruction (b), the model predicts the best grasp pose (c) and the target segmentation mask (d). The robot then executes the pick (e) and places the object at the predefined destination (f).

trial, the system takes an RGB observation together with a user instruction (We use only RGB for inference; depth is used only for execution/calibration.), and outputs an instruction-conditioned target mask and a corresponding grasp pose. These predictions are then used to execute a pick-and-place behavior on the physical robot. Performance is summarized using the task success rate, computed over repeated executions for each target object and scene configuration. A trial is counted as successful if the robot grasps the instructed target object, achieves a stable pickup, and places it at a predefined destination location. For each target object in each scene configuration, 20 trials are conducted to ensure consistency. Quantitative results are reported in Table VII. In the Isolated setting, success rates are 95% for Ball, 85% for Fruit, 95% for Can, and 100% for Glue. Under the Scattered configuration, success rates remain high across all targets. In the Cluttered setting, performance decreases to 90% (Ball), 90% (fruit), 85% (Can), and 90% (Glue). Across the four targets, the averaged success rate is 93.75% for Isolated scene, 97.5% for Scattered

scene, and 88.75% for Cluttered scene. In our real-robot trials, failures are most frequently observed when the arm must execute longer motions across the workspace and when specular highlights caused by lighting introduce reflective artifacts on object surfaces, which can degrade perception and subsequently affect grasp execution. Overall, the system produces stable qualitative behavior on the physical platform, and the observed success rates are consistent with the trends reported earlier in our dataset-based evaluation.

## VI. CONCLUSION

This work investigates the language-guided robotic grasping in cluttered and ambiguous environments, where robots must accurately localize and grasp objects based on both visual perception and natural language instructions. To this end, we introduced LGGD, an end-to-end grasp detection framework that performs deep cross-modal integration between RGB images and textual queries. Our model distributes linguistic information throughout the entire grasp

prediction pipeline, rather than fusing modalities only at an early stage. Through a feature fusion module with dual-path cross-attention, a language-conditioned upsampling mechanism, and a text-guided decoder, our model progressively refines spatial-semantic alignment between language intent and visual features. Furthermore, the residual refinement stage enhances grasp robustness and reduces prediction noise in cluttered scenes. Experimental evaluations on the OCID-VLG and Grasp-Anything++ datasets demonstrate that LGGD significantly outperforms previous language-guided grasping baselines in both segmentation accuracy and grasp success metrics. Real-robot experiments further confirm its practical effectiveness, showing that the model can reliably interpret human instructions and generate feasible grasp poses even in visually complex scenarios. Overall, this research highlights the potential of deeply integrated vision-language reasoning for robotic manipulation. The proposed framework represents a step toward more semantically grounded, instruction-aware robotic systems that can flexibly collaborate with humans in real-world environments.

## REFERENCES

[1] H. Cao, G. Chen, Z. Li, Q. Feng, J. Lin, and A. Knoll, "Efficient grasp detection network with gaussian-based grasp representation for robotic manipulation," *IEEE/ASME Transactions on Mechatronics*, vol. 28, no. 3, pp. 1384–1394, 2022. 1, 3, 4, 12

[2] N. Correll, K. E. Bekris, D. Berenson, O. Brock, A. Causo, K. Hauser, K. Okada, A. Rodriguez, J. M. Romano, and P. R. Wurman, "Analysis and observations from the first amazon picking challenge," *IEEE Transactions on Automation Science and Engineering*, vol. 15, no. 1, pp. 172–188, 2018. 1

[3] F. Surez-Ruiz, X. Zhou, and Q.-C. Pham, "Can robots assemble an ikea chair?" *Science Robotics*, vol. 3, no. 17, 2018. 1

[4] Z. Fu, T. Z. Zhao, and C. Finn, "Mobile aloha: Learning bimanual mobile manipulation using low-cost whole-body teleoperation," in *CoRL*, 2024. 1

[5] C. Borst, M. Fischer, and G. Hirzinger, "Grasp planning: how to choose a suitable task wrench space," in *ICRA*, 2004. 1, 2

[6] A. Bicchi and V. Kumar, "Robotic grasping and contact: a review," in *ICRA*, 2000. 1, 2

[7] I. Lenz, H. Lee, and A. Saxena, "Deep learning for detecting robotic grasps," *The International Journal of Robotics Research*, vol. 34, no. 4-5, pp. 705–724, 2015. 1, 3

[8] 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." in *RSS*, 2017. 1, 3

[9] H. Cao, G. Chen, Z. Li, Y. Hu, and A. Knoll, "Neurograsp: Multimodal neural network with euler region regression for neuromorphic vision-based grasp pose estimation," *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1–11, 2022. 1

[10] B. Li, H. Cao, Z. Qu, Y. Hu, Z. Wang, and Z. Liang, "Event-based robotic grasping detection with neuromorphic vision sensor and event-grasping dataset," *Frontiers in neurorobotics*, vol. 14, p. 51, 2020. 1

[11] H. Zhang, X. Zhou, X. Lan, J. Li, Z. Tian, and N. Zheng, "A real-time robotic grasping approach with oriented anchor box," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 51, no. 5, pp. 3014–3025, 2019. 1, 3

[12] S. Kumra and C. Kanan, "Robotic grasp detection using deep convolutional neural networks," in *IROS*, 2017. 1, 3

[13] D. Morrison, P. Corke, and J. Leitner, "Closing the Loop for Robotic Grasping: A Real-time, Generative Grasp Synthesis Approach," in *RSS*, 2018. 1, 3, 12

[14] S. Kumra, S. Joshi, and F. Sahin, "Antipodal robotic grasping using generative residual convolutional neural network," in *IROS*, 2020. 2, 3, 11, 12, 14

[15] S. Yu, D.-H. Zhai, Y. Xia, H. Wu, and J. Liao, "Se-resunet: A novel robotic grasp detection method," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 5238–5245, 2022. 2, 3

[16] H. Cao, G. Chen, Z. Li, J. Lin, and A. Knoll, "Residual squeeze-and-excitation network with multi-scale spatial pyramid module for fast robotic grasping detection," in *ICRA*, 2021. 2, 3, 11

[17] S. Wang, Z. Zhou, and Z. Kan, "When transformer meets robotic grasping: Exploits context for efficient grasp detection," *IEEE robotics and automation letters*, vol. 7, no. 3, pp. 8170–8177, 2022. 2, 3

[18] Q. Zhang, J. Zhu, X. Sun, and M. Liu, "Htc-grasp: A hybrid transformer-cnn architecture for robotic grasp detection," *Electronics*, vol. 12, no. 6, p. 1505, 2023. 2, 3

[19] 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. 2, 3, 4, 5, 11, 12

[20] J. Li, D. Li, C. Xiong, and S. Hoi, "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation," in *ICML*, 2022. 2

[21] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *NeurIPS*, 2023. 2

[22] M. Shridhar, L. Manuelli, and D. Fox, "Cliport: What and where pathways for robotic manipulation," in *CoRL*, 2022. 2, 3, 12

[23] A. D. Vuong, M. N. Vu, B. Huang, N. Nguyen, H. Le, T. Vo, and A. Nguyen, "Language-driven grasp detection," in *CVPR*, 2024. 2, 3, 10, 11, 12

[24] K. Xu, S. Zhao, Z. Zhou, Z. Li, H. Pi, Y. Zhu, Y. Wang, and R. Xiong, "A joint modeling of vision-language-action for target-oriented grasping in clutter," in *ICRA*, 2023. 2, 3

[25] N. Nguyen, M. N. Vu, B. Huang, A. Vuong, N. Le, T. Vo, and A. Nguyen, "Lightweight language-driven grasp detection using conditional consistency model," in *IROS*, 2024. 2, 3

[26] V. Vo, M. Vu, B. Huang, A. Vuong, N. Le, T. Vo, and A. Nguyen, "Language-driven grasp detection with mask-guided attention," in *IROS*, 2024. 2, 4

[27] S. Jin, J. Xu, Y. Lei, and L. Zhang, "Reasoning grasping via multimodal large language model," *CoRL*, 2024. 2

[28] G. Tziafas, X. Yucheng, A. Goel, M. Kasaei, Z. Li, and H. Kasaei, "Language-guided robot grasping: Clip-based referring grasp synthesis in clutter," in *CoRL*, 2023. 2, 10, 11, 12, 14

[29] S. Ainetter and F. Fraundorfer, "End-to-end trainable deep neural network for robotic grasp detection and semantic segmentation from rgb," in *ICRA*, 2021. 2, 4, 10, 12

[30] A. Depierre, E. Dellandréa, and L. Chen, "Jacquard: A large scale dataset for robotic grasp detection," in *IROS*, 2018. 2, 4

[31] A. D. Vuong, M. N. Vu, H. Le, B. Huang, H. T. T. Binh, T. Vo, A. Kugi, and A. Nguyen, "Grasp-anything: Large-scale grasp dataset from foundation models," in *ICRA*, 2024. 2, 10

[32] 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. 2

[33] 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. 3

[34] Y. Xu, M. Kasaei, H. Kasaei, and Z. Li, "Instance-wise grasp synthesis for robotic grasping," in *ICRA*, 2023. 3, 12

[35] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "Yolact: Real-time instance segmentation," in *ICCV*, 2019. 3

[36] R. Xu, F.-J. Chu, and P. A. Vela, "Gknet: Grasp keypoint network for grasp candidates detection," *The International Journal of Robotics Research*, vol. 41, no. 4, pp. 361–389, 2022. 3

[37] Y. Wu, F. Zhang, and Y. Fu, "Real-time robotic multigrasp detection using anchor-free fully convolutional grasp detector," *IEEE Transactions on Industrial Electronics*, vol. 69, no. 12, pp. 13 171–13 181, 2021. 3

[38] H. Cheng, Y. Wang, and M. Q.-H. Meng, "Anchor-based multi-scale deep grasp pose detector with encoded angle regression," *IEEE Transactions on Automation Science and Engineering*, 2023. 3

[39] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, "Swin-unet: Unet-like pure transformer for medical image segmentation," in *ECCVW*, 2022. 3

[40] Y. Lu, Y. Fan, B. Deng, F. Liu, Y. Li, and S. Wang, "Vl-grasp: a 6-dof interactive grasp policy for language-oriented objects in cluttered indoor scenes," in *IROS*, 2023. 3

[41] Q. Sun, H. Lin, Y. Fu, Y. Fu, and X. Xue, "Language guided robotic grasping with fine-grained instructions," in *IROS*, 2023. 3

[42] T. Nguyen, M. N. Vu, B. Huang, A. Vuong, Q. Vuong, N. Le, T. Vo, and A. Nguyen, "Language-driven 6-dof grasp detection using negative prompt guidance," in *ECCV*, 2024. 3

[43] C. Cheang, H. Lin, Y. Fu, and X. Xue, "Learning 6-dof object poses to grasp category-level objects by language instructions," in *ICRA*, 2022. 3

[44] Y. Yang, X. Lou, and C. Choi, "Interactive robotic grasping with attribute-guided disambiguation," in *ICRA*, 2022. 3

[45] J. Liu, J. Xie, L. Xiao, C. Wang, and F. Zhou, "Hierarchical multi-modal fusion for language-conditioned robotic grasping detection in clutter," *IEEE Robotics and Automation Letters*, 2024. 3

[46] Y. Chen, R. Xu, Y. Lin, and P. A. Vela, "A joint network for grasp detection conditioned on natural language commands," in *ICRA*, 2021. 3

[47] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016. 3, 4, 11

[48] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, p. 17351780, Nov. 1997. 3

[49] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 2019, pp. 4171–4186. 3

[50] H. H. Nguyen, A. Vuong, A. Nguyen, I. Reid, and M. N. Vu, "Graspmamba: A mamba-based language-driven grasp detection framework with hierarchical feature learning," *IROS*, 2025. 3

[51] S. Noh, J. Kim, D. Nam, S. Back, R. Kang, and K. Lee, "Graspsam: When segment anything model meets grasp detection," in *ICRA*, 2025. 4

[52] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *NeurIPS*, 2017. 4, 6, 12, 14

[53] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "Film: Visual reasoning with a general conditioning layer," in *AAAI*, 2018. 6, 7

[54] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *CVPR*, 2018. 8

[55] Y. Chen, X. Dai, M. Liu, D. Chen, L. Yuan, and Z. Liu, "Dynamic convolution: Attention over convolution kernels," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 030–11 039. 8

[56] M. Suchi, T. Patten, D. Fischinger, and M. Vincze, "Easylabel: A semi-automatic pixel-wise object annotation tool for creating robotic rgb-d datasets," in *ICRA*, 2019. 10

[57] Z. Wang, Y. Lu, Q. Li, X. Tao, Y. Guo, M. Gong, and T. Liu, "Cris: Clip-driven referring image segmentation," in *CVPR*, 2022. 11

[58] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *NeurIPS*, 2019. 11

[59] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *ICLR*, 2019. 11

[60] L. Ilya and H. Frank, "SGDR: Stochastic gradient descent with warm restarts," in *ICLR*, 2017. 11

[61] 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. 12

[62] L. H. Li, P. Zhang, H. Zhang, J. Yang, C. Li, Y. Zhong, L. Wang, L. Yuan, L. Zhang, J.-N. Hwang *et al.*, "Grounded language-image pre-training," in *CVPR*, 2022. 12