

# DexGraspVLA: A Vision-Language-Action Framework Towards General Dexterous Grasping

**Yifan Zhong<sup>1,2\*</sup>, Xuchuan Huang<sup>1,2\*</sup>, Ruochong Li<sup>2,3</sup>, Ceyao Zhang<sup>1,2</sup>**  
**Yitao Liang<sup>1,2</sup>, Yaodong Yang<sup>1,2†</sup>, Yuanpei Chen<sup>1,2†</sup>**

## Abstract

Dexterous grasping remains a fundamental yet challenging problem in robotics. A general-purpose robot must be capable of grasping diverse objects in arbitrary scenarios. However, existing research typically relies on restrictive assumptions, such as single-object settings or limited environments, leading to constrained generalization. We present **DexGraspVLA**, a hierarchical framework for general **dexterous grasping** in cluttered scenes based on RGB image perception and language instructions. It utilizes a pre-trained **Vision-Language** model as the high-level task planner and learns a diffusion-based policy as the low-level Action controller. The key insight to achieve robust generalization lies in iteratively transforming diverse language and visual inputs into domain-invariant representations via foundation models, where imitation learning can be effectively applied due to the alleviation of domain shift. Notably, our method achieves a 90+% success rate under *thousands* of unseen object, lighting, and background combinations in a “zero-shot” environment. Empirical analysis confirms the consistency of internal model behavior across environmental variations, thereby validating our design and explaining its generalization performance. DexGraspVLA also demonstrates free-form long-horizon prompt execution, robustness to adversarial objects and human disturbance, and failure recovery, which are rarely achieved simultaneously in prior work. Extended application to nonprehensile object grasping further proves its generality. Code, model, and video are available at [dexgraspvla.github.io](https://dexgraspvla.github.io).



Figure 1: We propose **DexGraspVLA**, a hierarchical vision-language-action framework that reaches a 90+% dexterous grasping success rate under thousands of unseen object, lighting, and background combinations in a “zero-shot” real-world environment. It robustly handles adversarial objects, human disturbance, failure recovery, and free-form long-horizon grasping prompts.

\*Equal contribution. <sup>1</sup>Institute for Artificial Intelligence, Peking University. <sup>2</sup>PKU-PsiBot Joint Lab.  
<sup>3</sup>Hong Kong University of Science and Technology (Guangzhou). † Corresponding authors: Yuanpei Chen <[yuanpei.chen312@gmail.com](mailto:yuanpei.chen312@gmail.com)>, Yaodong Yang <[yaodong.yang@pku.edu.cn](mailto:yaodong.yang@pku.edu.cn)>.

## 1 Introduction

Dexterous multi-fingered hands, as versatile robotic end-effectors, have demonstrated remarkable capabilities across various manipulation tasks [1, 2, 3, 4, 5, 6, 7, 8, 9]. Among these capabilities, grasping serves as the most fundamental prerequisite, yet remains one of the most challenging problems. Existing dexterous grasping approaches are primarily evaluated on isolated objects or under simplified settings. Nevertheless, real-world applications demand more general grasping capabilities that can function reliably in diverse scenarios such as industrial manufacturing and household environments. However, developing general dexterous grasping capabilities presents multifaceted challenges. At the object level, the policy must generalize across diverse physical properties including geometries, masses, textures, and orientations. Beyond object characteristics, the system must also demonstrate robustness to various environmental factors, such as lighting conditions, background complexities, and potential disturbances. Compounding these challenges, multi-object scenarios introduce additional complexity that demands sophisticated reasoning capabilities. For instance, in cluttered or stacked environments, planning the optimal sequence to grasp all objects becomes a crucial cognitive task that extends beyond simple grasp execution.

One line of research adopts a two-stage pipeline: first predicting a target grasp pose from single-frame perception, then executing open-loop motion planning to reach the pose [10, 11, 12]. However, such methods are heavily constrained by precise camera calibration and mechanical accuracy requirements. By contrast, end-to-end paradigms, such as imitation learning and reinforcement learning, enable closed-loop grasping by continuously adjusting actions based on real-time feedback, offering more robust and adaptive solutions. Reinforcement learning has achieved notable successes in simulation [13, 14, 15, 16], but simulating real-world physical complexity remains challenging, resulting in an inevitable sim-to-real gap. Imitation learning, by learning directly from human demonstrations [17, 18, 19], offers a viable alternative for dexterous grasping. However, such approaches often struggle with generalization beyond the demonstration data. This issue is further compounded by the impracticality of collecting expert trajectories across the full spectrum of objects and environmental variations required for general grasping. As a result, a key challenge is how to effectively leverage limited expert data to achieve broad generalization.

The rapid emergence of vision and language foundation models [20, 21, 22, 23, 24] presents promising opportunities for robotic manipulation. Leveraging internet-scale data in pre-training, these models demonstrate remarkable scene understanding and generalization capabilities for visual and linguistic inputs. To harness these capabilities for decision making, researchers have explored the integration of vision and language foundation models into action generation, leading to the development of vision-language-action (VLA) models [25]. One straightforward approach directly trains vision-language models (VLMs) on robot data in an end-to-end manner [26, 27]. However, this paradigm typically demands an enormous volume of manually collected demonstrations [28, 29] in an attempt to encompass the full range of real-world diversity and complexity. Even so, these models exhibit markedly reduced performance on unseen scenarios and still require further data collection and fine-tuning to handle new conditions. In addition, the substantial disparity between robotics datasets and the massive pre-training corpora leads to catastrophic forgetting, compromising the model’s valuable long-range reasoning capabilities. Other research endeavors have proposed hierarchical VLA architectures [30, 31] to decouple high-level task planning and low-level action control. While they hold promise for long-horizon task completion [32] and more general embodied capabilities, effectively utilizing foundation models to learn generalizable low-level controllers and achieve embodied reasoning remains underexplored.

In this paper, we present **DexGraspVLA**, the first hierarchical Vision-Language-Action framework for general dexterous grasping that integrates the complementary strengths of foundation models and imitation learning. At the high level, it utilizes a pre-trained VLM as a task planner, which interprets and reasons about language instructions, plans the overall grasping task, and generates task affordance signals. Guided by these signals and multimodal inputs, a low-level diffusion-based modularized controller produces closed-loop action sequences. The essence of DexGraspVLA lies in leveraging foundation models to iteratively transform *diverse* vision and language inputs into *domain-invariant* representations, where it then efficiently and effectively applies diffusion-based imitation learning to capture the action distribution in our dexterous grasping dataset. As a result, novel scenarios outside the training set no longer induce failures, because the foundation models translate them into representations resembling those encountered during training — thus remaining within the learned policy’s domain. This approach fuses the extensive world knowledge of foundation models with

the strong action modeling capacity of imitation learning, thereby enabling robust generalization performance in real-world applications.

Notably, DexGraspVLA achieves an unprecedented 90.8% success rate for grasping in cluttered scenes spanning 1,287 unseen object, lighting, and background combinations, all tested in a “zero-shot” environment. It robustly handles adversarial objects, human disturbances, and failure recovery. On a single-object benchmark, it attains 98.6% success, outperforming existing baselines whose controller learns directly from raw visual inputs by at least 48%. Analysis reveals consistent internal representations and attention maps within DexGraspVLA across varying environments, validating framework design and explaining its performance. Beyond single-step tasks, our framework executes free-form, long-horizon language instructions with embodied reasoning, achieving 89.6% average task success rate. We further extend it to nonprehensile object grasping [33, 34], demonstrating its generality for diverse manipulation skills. These results establish DexGraspVLA as a general, instruction-driven framework that learns from limited demonstrations and generalizes reliably to real-world scenarios, marking a promising step toward general dexterous grasping and beyond.

## 2 Related Work

**Dexterous Grasping.** Dexterous grasping typically falls into two categories: two-stage approaches and end-to-end methods. Two-stage approaches first generate a grasp pose and then control the dexterous hand targeting this pose. The main challenge is generating high-quality grasp poses based on visual observation. Current methods employ sample-based [35, 36], optimization-based [11, 12, 37, 38, 39, 40], or regression-based [41, 42] approaches to generate target grasp poses, followed by motion planning for robot execution. While these methods benefit from decoupled perception and control and simulation data generation, they typically suffer from the lack of closed-loop feedback and sensitivity to disturbances and calibration errors. End-to-end methods directly model grasping trajectories using imitation learning or reinforcement learning. Recent works explored training dexterous manipulation using reinforcement learning in simulation environments and transfer to the real-world [43, 1, 44, 2, 45, 46, 13, 14, 15, 16, 47, 3, 4, 5, 6, 7, 48, 49, 50, 51]. [52] and [53] generate affordance hints using computer vision methods and use reinforcement learning to train policies based on these features. [54] and [55] demonstrate some generalization capabilities in the real-world through large-scale parallel training in simulation. However, this reliance on simulation inevitably introduces sim-to-real gaps, while direct training in the real-world suffers from poor sample efficiency. Recently, imitation learning using human demonstration has shown remarkable results in complex tasks [53, 56, 17, 18, 19, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69]. These methods require human teleoperation to collect demonstration data and directly learn the distribution in the dataset. While being easier to train, this approach limits their generalization capabilities.

**Foundation Models for Robotics.** Recent advances in foundation models pre-trained on web-scale data have led to strong generalization in vision [23, 24, 70, 20, 71, 72] and sophisticated multimodal reasoning in VLMs [22, 73]. Applying these models to robotics is a promising direction. A common approach, as in RT-X [29], OpenVLA [26], and Pi0 [27], directly fine-tunes VLMs on robot data. However, this requires massive, diverse demonstrations [28, 29, 27], yet still struggles with unseen scenarios and often degrades original vision-language capabilities due to catastrophic forgetting. An alternative paradigm adopts hierarchical VLA architectures, where a high-level module interprets vision and language to generate intermediate action guidance, in the form of language motion [74], trajectory [75, 31, 76], latent representation [77, 78], etc., conditioned upon which a low-level module executes control. Most relevant to us are methods that also employ affordance as action guidance, such as affordance maps [79] and keypoints [80]. They typically use pre-trained VLMs for affordance prediction and apply open-loop motion planning. In contrast, our high-level planner utilizes bounding boxes as reliable affordance guidance and incorporates explicit reasoning to handle free-form prompts, while the low-level controller is a learned closed-loop policy, enhancing robustness. To learn this policy, we leverage foundation models to extract domain-invariant representations, enabling efficient imitation learning from limited demonstrations and achieving strong zero-shot generalization by offloading real-world complexity to the model’s perceptual backbone.

## 3 Problem Formulation

Our goal is to develop a vision-based control policy for language-guided dexterous grasping, formulated as a sequential decision-making problem. Initially, a language instruction  $l$  is given, *e.g.* “grasp the toy”, to directly specify the target object. At each timestep  $t$ , the policy  $\pi$  receives a first-view

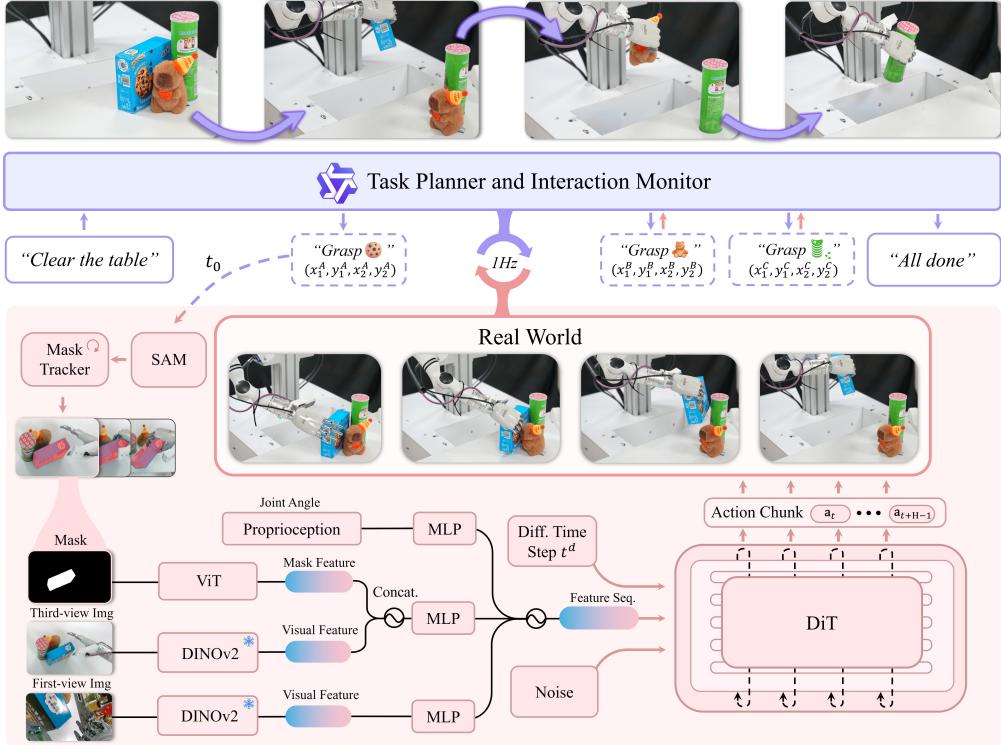


Figure 2: **Overview of DexGraspVLA.** A pre-trained VLM-based high-level **planner** decomposes language prompts into object-level grasping instructions with bounding boxes. The diffusion-based low-level **controller** tracks the target mask, encodes multimodal observations (RGB images, mask, proprioception), and predicts an action chunk via a DiT model. The **planner** monitors execution and continually proposes new instructions based on the updated scene until the task is fully completed.

image  $\mathbf{I}_t^w \in \mathbb{R}^{H \times W \times 3}$  from the wrist camera ( $H$  and  $W$  denote the height and width of the image), a third-view image  $\mathbf{I}_t^h \in \mathbb{R}^{H \times W \times 3}$  from the head camera, and the robot proprioception  $\mathbf{s}_t \in \mathbb{R}^{13}$  consisting of arm and hand joint angles  $\mathbf{s}_t^{\text{arm}} \in \mathbb{R}^7, \mathbf{s}_t^{\text{hand}} \in \mathbb{R}^6$ . Conditioned on these observations, the robot produces an action  $\mathbf{a}_t = (\mathbf{a}_t^{\text{arm}}, \mathbf{a}_t^{\text{hand}}) \in \mathbb{R}^{13}$ , where  $\mathbf{a}_t^{\text{arm}} \in \mathbb{R}^7$  and  $\mathbf{a}_t^{\text{hand}} \in \mathbb{R}^6$  denote the target joint angles for arm and hand respectively, by sampling from the action distribution  $\pi(\cdot | \{\mathbf{I}_j^w\}_{j=0}^t, \{\mathbf{I}_j^h\}_{j=0}^t, \{\mathbf{s}_j\}_{j=0}^t, l)$ . This process continues until a termination condition is reached. The robot receives a binary reward  $r \in \{0, 1\}$  indicating whether it has completed the instruction  $l$  successfully. The goal of the policy  $\pi$  is to maximize the expected reward  $\mathbb{E}_{l, \{(I_j^w, I_j^h, s_j, a_j)\}_{j=0}^T} [r]$ .

More generally, we consider the case where the user prompt  $p$  may be a long-horizon task involving multiple grasping processes, such as “*clear the table*”. This requires the policy  $\pi$  to reason about the prompt, decompose it into individual grasping instructions  $\{l_i\}$ , and complete them sequentially.

## 4 Methods

This section introduces DexGraspVLA, the first hierarchical VLA framework for dexterous grasping. We will first elaborate DexGraspVLA framework (Section 4.1) and then detail our data collection procedure (Section 4.2), which together enable the training of a dexterous grasping policy.

### 4.1 DexGraspVLA Framework

As illustrated in Figure 2, DexGraspVLA adopts a hierarchical and modularized architecture composed of a planner and a controller. Below we explain how each part is designed.

**Planner.** We recognize that to achieve general dexterous grasping, the model must handle multimodal inputs, perform visual grounding, and conduct reasoning about user prompts. Building upon recent advances, we adopt an off-the-shelf pre-trained Qwen VLM [81, 73] as a high-level planner to dynamically plan and monitor the dexterous grasping workflow. Given a user prompt  $p$  (e.g., “clear

the table”), the planner first records the initial head image  $\mathbf{I}_0^h$  and, conditioning on the observation, proposes a grasping instruction  $l$  (e.g., “grasp the cookie”).

For each instruction  $l$ , the planner guides the low-level controller by marking the target object bounding box  $(x_1, y_1, x_2, y_2)$  as task affordance in the head camera image  $\mathbf{I}_{t_0}^h$  at the initial timestep  $t_0$ . While the phrasing and content of language instruction can be diverse and flexible for different users and cases, *i.e.*, showing *domain-variance*, the bounding box is a consistent format for object positioning regardless of the changes in language and visual inputs, *i.e.*, achieving *domain-invariance*. Thus, this transformation alleviates the learning challenge for the controller.

On issuing the bounding box, the planner monitors controller execution by querying cameras at 1Hz. Upon a successful grasp, it triggers a scripted placing motion. After each grasp attempt, the planner resets the robot to the initial state. Based on the initial and current head images as well as the prompt  $p$ , it proposes a new instruction  $l$ , repeating this loop until user prompt  $p$  is completed.

**Controller.** Based on the target bounding box  $(x_1, y_1, x_2, y_2)$ , the controller aims to grasp the intended object in cluttered environments. We feed this bounding box as input to SAM [23] to obtain an initial binary mask  $\mathbf{m}_0 \in \{0, 1\}^{H \times W \times 1}$  of the target object and then use Cutie [82] to continuously track the mask over time, producing  $\mathbf{m}_t$  at each timestep  $t$ . This ensures accurate identification in cluttered scenes throughout the process. The problem is to learn the policy  $\pi$  that effectively models the action distribution  $\pi(\cdot | \mathbf{I}_t^w, \mathbf{I}_t^h, \mathbf{s}_t, \mathbf{m}_t)$ .

To achieve general-purpose dexterous grasping, the system must generalize effectively across diverse real-world scenarios. However, the high variability in raw visual inputs  $\mathbf{I}_t^w, \mathbf{I}_t^h$  poses a fundamental challenge to learning task-critical representations. Traditional imitation learning approaches often fail catastrophically even under minor variations in objects or environmental conditions. To address this issue, our solution is again to convert potentially *domain-varying* inputs into *domain-invariant* representations suitable for imitation learning. We recognize that while pixel-level perception can vary widely, the fine-grained semantic features extracted by large foundation models tend to be more robust and consistent. Thus, we utilize a feature extractor  $\phi$ , such as DINOv2 [20] that has been pre-trained on internet-scale data, to obtain features from raw images. At each timestep  $t$ , we obtain head camera image features  $\mathbf{z}_t^h = \phi^h(\mathbf{I}_t^h) \in \mathbb{R}^{L^h \times D^h}$ , and wrist camera image features  $\mathbf{z}_t^w = \phi^w(\mathbf{I}_t^w) \in \mathbb{R}^{L^w \times D^w}$ , where  $L^h, D^h, L^w, D^w$  denote length and hidden dimension of the feature sequences for head and wrist respectively. As we show in Section 5.4, these extracted features remain comparatively invariant to distracting visual factors.

Up to now, raw language and vision inputs, including instruction  $l$  and images  $\mathbf{I}_t^w, \mathbf{I}_t^h$ , have been iteratively transformed into domain-invariant representations, including mask  $\mathbf{m}_t$  and features  $\mathbf{z}_t^h, \mathbf{z}_t^w$ , by leveraging foundation models. This lays the stage for imitation learning. We now learn the policy  $\pi$  that predicts an action chunk of horizon  $H$  conditioning on these representations.

To fuse the object mask with head camera features, we project  $\mathbf{m}_t$  into the head image feature space using a randomly initialized ViT, producing  $\mathbf{z}_t^m \in \mathbb{R}^{L^h \times D^h}$ , and concatenate it with  $\mathbf{z}_t^h$  patch-wise to obtain  $\tilde{\mathbf{z}}_t^h \in \mathbb{R}^{L^h \times 2D^h}$ . Subsequently, we map  $\tilde{\mathbf{z}}_t^h$ , wrist-camera features  $\mathbf{z}_t^w$ , and robot state  $\mathbf{s}_t$  into a common embedding space with separate MLPs, yielding  $\tilde{\mathbf{z}}_t^h, \tilde{\mathbf{z}}_t^w$ , and  $\tilde{\mathbf{z}}_t^s$ . These embeddings are then concatenated to form the full observation feature sequence  $\tilde{\mathbf{z}}_t^{\text{obs}} \in \mathbb{R}^{(1+L^h+L^w) \times D}$ .

For action prediction, we employ a DiT [83] to generate multi-step actions, following the diffusion policy paradigm [84, 85, 27]. At each timestep  $t$ , we bundle the next  $H$  actions into a chunk  $\mathbf{A}_t = \mathbf{a}_{t:t+H} = [\mathbf{a}_t, \mathbf{a}_{t+1}, \dots, \mathbf{a}_{t+H-1}]$ . During training, a random diffusion step  $t^d = k$  is sampled, and Gaussian noise  $\epsilon$  is added to  $\mathbf{A}_t$ , yielding the noised action tokens  $\mathbf{x}_k = \alpha_k \mathbf{A}_t + \sigma_k \epsilon$ , where  $\alpha_k$  and  $\sigma_k$  are the standard DDPM coefficients. We then feed  $\mathbf{x}_k$  into the DiT alongside the observation feature sequence  $\tilde{\mathbf{z}}_t^{\text{obs}}$ . Each DiT layer performs bidirectional self-attention over the action tokens, cross-attention to  $\tilde{\mathbf{z}}_t^{\text{obs}}$ , and MLP transformations, ultimately predicting the original noise  $\epsilon$ . By minimizing the discrepancy between predicted and true noise, the model learns to reconstruct the ground-truth action chunk  $\mathbf{A}_t$ . At inference time, iterative denoising steps recover the intended multi-step action sequence from the learned distribution, enabling robust imitation of complex, long-horizon behaviors. We also employ the receding horizon control strategy that only executes the first  $H_a$  actions before generating a new action chunk prediction, enhancing real-time responsiveness.

Overall, DexGraspVLA performs imitation learning on *domain-invariant* representations derived from *domain-varying* inputs via foundation models. This approach not only leverages the world

knowledge and generalization capabilities of foundation models, but also effectively captures the mapping from these abstracted representations to the final action output.

## 4.2 Data Collection

To train our dexterous grasping policy, we manually collect a dataset consisting of 2,094 successful demonstrations in cluttered scenes using 36 household objects varying in size, weight, geometry, texture, material, and category. Each episode  $\tau = \{(\mathbf{I}_t^h, \mathbf{I}_t^w, \mathbf{s}_t, \mathbf{m}_t, \mathbf{a}_t)\}_{t=0}^T$  records raw camera images  $\mathbf{I}_t^h, \mathbf{I}_t^w$ , robot proprioception  $\mathbf{s}_t$ , object mask  $\mathbf{m}_t$ , and action  $\mathbf{a}_t$  at each timestep  $t$ . The mask  $\mathbf{m}_t$  is labeled in the same way as in the controller. For each object, we place it randomly and collect multiple grasping demonstrations, with the surrounding objects randomized between episodes. These demonstrations are performed at typical human motion speeds, taking about 3.5 s each. They undergo rigorous manual inspection to ensure quality and reliability. The DexGraspVLA controller is trained on this dataset with imitation learning.

## 5 Experiments

In this section, we comprehensively evaluate DexGraspVLA’s performance. To ensure real-world application, all experiments are conducted in a different environment from the demonstration setup. This “zero-shot” setting is fundamentally more challenging than most prior imitation learning research, which typically requires fine-tuning for high performance. Our experiments seek to address the following questions: (1) **Large-scale Generalization** (Section 5.2): Can DexGraspVLA generalize to thousands of unseen object, lighting, and background combinations? (2) **Baseline Comparison** (Section 5.3): How does it compare to baselines trained directly on raw visual inputs without frozen feature extractors? (3) **Mechanism Analysis** (Section 5.4): Are its internal model behaviors consistent under varying environments? (4) **Long-horizon Task Completion** (Section 5.5): How effectively does DexGraspVLA handle free-form, long-horizon instructions? (5) **Extension to Nonprehensile Grasping** (Section 5.6): Can it be extended to other dexterous manipulation skills beyond grasping?

### 5.1 Experiment Setups

**Hardware Platform.** As illustrated in Figure 3, the robot we use for dexterous grasping is a 7-DoF RealMan RM75-6F arm paired with a 6-DoF PsiBot G0-R hand. A RealSense D405C camera, mounted on the arm’s wrist, provides a first-person viewpoint, while a RealSense D435 camera on the robot’s head offers a third-person perspective. Objects to be grasped are placed on a table in front of the robot. The control frequency of the robot is 20 Hz.

**Baselines.** To the best of our knowledge, no existing method can directly serve as a baseline. Existing dexterous grasping methods either cannot follow language instruction for cluttered scene, rely on optimization unsuited to our linkage-based hand, or consider a different point-cloud based setting, while related VLA frameworks are incompatible with dexterous hands. Therefore, we compare the following methods: (1) *DexGraspVLA (Ours)*: A full implementation of DexGraspVLA. (2) *DexGraspVLA (DINOv2-train)*: Identical to *Ours* but with trainable DINOv2 encoders. (3) *DexGraspVLA (ViT-small)*: Identical to *Ours* but replaces DINOv2 with smaller, trainable ViTs. Empirically, DexGraspVLA (ViT-small) represents an enhanced version of diffusion policy [84], a SOTA imitation learning baseline. The planner is based on Qwen-VL-Chat [81] for all experiments except the long-horizon task (Section 5.5), which uses Qwen2.5-VL-72B-Instruct [73]. Implementation details are in Appendix A. To account for inference randomness, we report *Ours@k* ( $k = 1, 2, 3$ ) in Section 5.2, where up to  $k$  attempts are allowed per test. *Ours@1* is equivalent to *Ours*. Re-grasps performed by the policy after an initial failure within a single attempt are allowed and not counted separately.

### 5.2 Large-Scale Generalization Evaluation

**Tasks.** We curate 360 unseen objects, 6 unseen backgrounds, and 3 unseen lighting conditions. The objects span diverse sizes, weights, geometries, textures, materials, and categories, while remaining graspable by our dexterous hand (measured and visualized in Figure 4). Backgrounds and lighting conditions are selected to be visually distinct. We evaluate generalization through three grasping

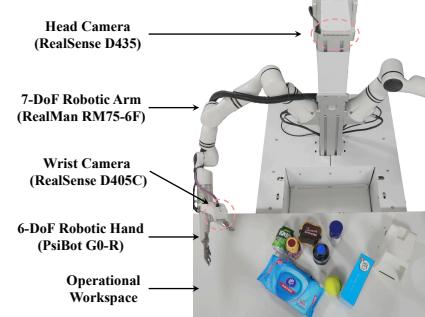


Figure 3: The hardware platform used for dexterous grasping.

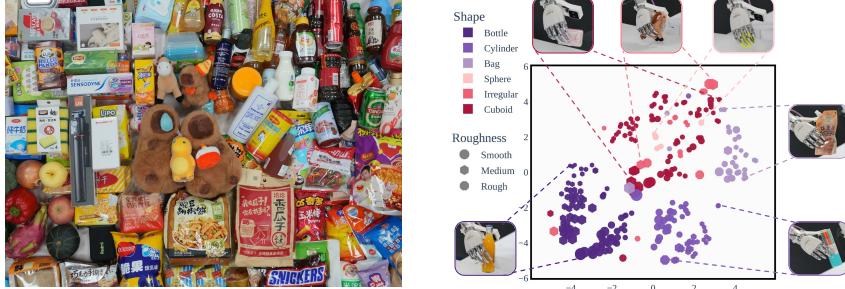


Figure 4: **(Left)** A representative part of all 360 unseen objects used to evaluate DexGraspVLA. **(Right)** A t-SNE projection illustrating the diversity and broad coverage of these objects in length, width, height, mass (denoted by marker size), roughness (marker type), and shape (marker color).

tasks in cluttered scenes (around 6 objects per scene): (1) *Unseen Objects*: Each of the 360 objects is grasped once in a random scene on a white table under white light (360 tests). (2) *Unseen Backgrounds*: A subset of 103 objects  $\mathcal{S}$  is used to create 103 scenes per background under white light, totaling 618 tests. (3) *Unseen Lightings*: The same  $\mathcal{S}$  is used to construct 103 scenes per lighting condition on a white table (309 tests). Details can be found in Appendix B.

**Metric.** A grasp is successful if the object is held 10 cm above the table for 20 s. Success rate is the ratio of successes to total tests; aggregated performance is a weighted average by task proportion.

**Results.** We present the quantitative results in Table 1. From the first row (“Ours@1”), DexGraspVLA achieves a 91.1% single-attempt success rate on 360 unseen objects, 90.5% on 6 unseen backgrounds, and 90.9% under 3 unseen lighting conditions, yielding a 90.8% aggregated success rate. These results demonstrate robust and accurate control of the dexterous hand to grasp specified objects from clutter in diverse unseen conditions, without domain-specific fine-tuning. This highlights strong generalization and suggests that our framework substantially alleviates the longstanding challenge in imitation learning — namely, overfitting to narrow domains. We further analyze the source of this generalization in 5.4 and extend its application in 5.6.

Qualitatively, DexGraspVLA robustly handles challenging cases involving transparent, deformable, reflective, or background-camouflaged objects. It also dexterously adapts its arm and hand to grasp objects with diverse geometries and poses — e.g., grasping a bottle from the side, picking up a small earbud case from the top, or retrieving an awkwardly placed box. The closed-loop policy enables re-grasping after failed attempts and tolerates human-induced perturbations by tracking object motion. Such robustness stems from three factors: first, foundation-model-based perception ensures semantic consistency under appearance variation; second, imitation learning avoids the need for explicit object modeling; and third, diffusion-based action head captures multi-modal action distributions.

From the second and third rows (“Ours@2” and “Ours@3”), we observe that allowing up to three attempts further boosts performance to 96.9%, indicating the capacity to reach even higher success rates. Finally, our model takes around 6 s on average to grasp an object, which is close to that of humans and ensures practical usability in real-world scenarios.

### 5.3 Comparison to Baselines without Frozen Vision Encoders

**Tasks & Metrics.** To compare DexGraspVLA with baselines that learn directly from raw visual inputs without frozen vision encoders, we conduct single-object grasping experiments using 13 seen and 8 unseen objects. Each object is placed at five table locations covering the reachable workspace and camera view. At each location, two independent grasp trials are performed, counted separately to account for inference randomness. This results in 210 tests in total, conducted under white tabletop and white lighting conditions. Success rates are reported in the same way as Section 5.2.

**Results.** Table 2 shows that DexGraspVLA (Ours) consistently achieves over 98% success on both seen and unseen objects, significantly outperforming DINoV2-train and ViT-small variants. Its near-perfect performance in a zero-shot setting indicates strong robustness to domain shift. Interestingly, performance on unseen objects slightly exceeds that on seen ones, suggesting that the model learns the grasping task itself rather than overfitting to training data. In contrast, baselines that map raw inputs to actions fail to generalize, as perceptual changes easily push them out of distribution.

Table 1: Large-scale generalization performance of DexGraspVLA in diverse unseen conditions.

	Unseen Objects (360)	Unseen Backgrounds (6 × 103)	Unseen Lightings (3 × 103)	Aggregated (1287)
Ours@1	91.1%	90.5%	90.9%	90.8%
Ours@2	95.3%	94.2%	95.1%	94.7%
Ours@3	96.7%	96.7%	97.4%	96.9%

#### 5.4 Internal Model Behavior Analysis

To further validate our design, we examine whether internal model behavior remains consistent under varying visual conditions, as shown in Figure 5. We test DexGraspVLA on the same cluttered scene (9 objects, target: “grasp the blue yogurt in the middle”) across four environments: a white table, a calibration board, a colorful tablecloth, and the same tablecloth under disco lighting. For clarity, we display only the tabletop region; full images are in Appendix B. While the head images in the first row of Figure 5 appear to be markedly diverse, the DINOv2 features in the second row look rather consistent. These features are visualized by mapping principal components to RGB channels as done in Oquab et al. [20]. Across environments, the object properties are robustly maintained and matched, which fundamentally allows DexGraspVLA trained on a single data domain to generalize. The third row shows that Cutie accurately tracks the object, providing the correct guidance to the controller. This provides the correct guidance to the controller. The fourth row shows that the attention maps are consistent across environments. The fifth row overlays the attention map on the raw image to confirm the reasonable attention pattern. All visualization details are provided in Appendix B. Therefore, we substantiate that DexGraspVLA indeed transforms perceptually diverse raw inputs into invariant representations, on which it effectively applies imitation learning to model the data distribution, explaining its superior generalization performance. Expectedly, it successfully grasps the yogurt in all four environments.

#### 5.5 Long-Horizon Task Evaluation

**Tasks.** This experiment evaluates DexGraspVLA’s capability to complete complex, long-horizon tasks. We design four types of user prompts: “*Clear the table*”, “*Grasp all bottles*”, “*Grasp all green objects*”, and “*Grasp all food*”. These prompts require common-sense and physical knowledge to identify appropriate grasping targets sequentially. For each prompt, we randomly configure 24 cluttered scenes. The scenes for “*Clear the table*” contain three unseen objects, while the remaining prompts each involve 3 – 4 unseen objects, among which two are relevant targets to be grasped. All tasks are conducted on a white tabletop under white lighting.

**Metric.** For each task, we report the task success rate as the proportion of tests that fully complete all required stages. We further report the average grasping attempts per object in the successful tests, along with success rates for instruction proposal, bounding box prediction, completion check of the planner, and grasp execution of the controller.

**Results.** Table 3 shows that DexGraspVLA achieves an 89.6% aggregated task success rate across four long-horizon prompts, with each target object attempted slightly more than once. The high-level planner grounds prompt semantics on the observation and proposes correct instructions with a 94.3% average success rate. Its bounding box prediction accuracy is consistently above 98%, which we further substantiate with evaluations in distraction conditions in Appendix D. The low-level

Table 2: DexGraspVLA significantly outperforms baselines in both seen and unseen single-object experiments in the zero-shot environment.

	Seen Objects	Unseen Objects	Aggregated
ViT-small	60.0%	35.0%	50.5%
DINOv2-train	30.0%	43.5%	34.8%
Ours	<b>98.5%</b>	<b>98.8%</b>	<b>98.6%</b>

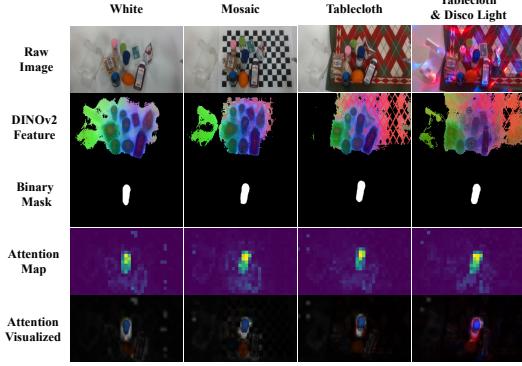


Figure 5: **DexGraspVLA is robust to environmental variations.** The same cluttered scene (1st row) is arranged in four visually different environments (four columns). DINOv2 features (2nd row), masks (3rd row), and attention maps (4th row) are consistent across variations. The 5th row confirms DexGraspVLA is attending to the correct object.

Based on the domain-invariant mask and the DINOv2 features, the DiT action head now predicts the subsequent actions. In the fourth row, we average and normalize all cross-attentions to the head image from DiT. We find that all attention maps exhibit the same behavior of focusing on the target object instead of being distracted by environments. The fifth row overlays the attention map on the raw image to confirm the reasonable attention pattern. All visualization details are provided in Appendix B. Therefore, we substantiate that DexGraspVLA indeed transforms perceptually diverse raw inputs into invariant representations, on which it effectively applies imitation learning to model the data distribution, explaining its superior generalization performance. Expectedly, it successfully grasps the yogurt in all four environments.

Based on the domain-invariant mask and the DINOv2 features, the DiT action head now predicts the subsequent actions. In the fourth row, we average and normalize all cross-attentions to the head image from DiT. We find that all attention maps exhibit the same behavior of focusing on the target object instead of being distracted by environments. The fifth row overlays the attention map on the raw image to confirm the reasonable attention pattern. All visualization details are provided in Appendix B. Therefore, we substantiate that DexGraspVLA indeed transforms perceptually diverse raw inputs into invariant representations, on which it effectively applies imitation learning to model the data distribution, explaining its superior generalization performance. Expectedly, it successfully grasps the yogurt in all four environments.

Table 3: DexGraspVLA on long-horizon prompts.

	Clear Table	Grasp Green	Grasp Bottles	Grasp Food	Aggr.
Task Success Rate	95.8%	87.5%	91.7%	83.3%	89.6%
Avg. Attempts per Grasp	1.09	1.14	1.09	1.19	1.12
Planner: Instruction Proposal	100.0%	92.6%	94.3%	88.1%	94.3%
Planner: BBox Accuracy	98.7%	98.2%	98.1%	98.3%	98.4%
Controller: Grasping	91.0%	92.6%	92.5%	91.5%	92.2%
Planner: Completion Check	98.7%	94.4%	96.2%	94.9%	96.3%

controller, leveraging its robust and generalizable grasping policy, executes individual grasps with over 91% success, enabling reliable multi-step completion. Additionally, the planner detects task completion with over 94% accuracy, preventing redundant actions. These results highlight the synergy between the high-level and low-level modules in DexGraspVLA, showcasing the effectiveness of its hierarchical framework for long-horizon tasks. An example can be found in Appendix C.

### 5.6 Extended Application to Nonprehensile Grasping

**Tasks & Metric.** To demonstrate that DexGraspVLA is applicable to manipulation skills beyond dexterous grasping, we apply the same hierarchical framework to a nonprehensile grasping task. Specifically, we curate 32 flat, wide-surface objects that are difficult to grasp directly—such as plates, boxes, and books—and collect 1,029 human demonstrations in cluttered scenes. In these demonstrations, the robot first performs a pre-grasp manipulation by pushing the object toward the table edge, creating an accessible pose, and then executes a final grasp to lift it. We keep the DexGraspVLA planner unchanged and train the controller on this dataset; details are provided in Appendix A. To evaluate generalization, we curate 18 previously unseen nonprehensile objects across common household categories and design three types of tasks in a zero-shot environment: (1) *Unseen Objects* (36 tests): Each object is placed in two cluttered scenes with varying poses on a white table under white light. (2) *Unseen Lighting* (36 tests): The same protocol is applied on a white table under disco light. (3) *Unseen Backgrounds* (72 tests): The same protocol is applied under white light on a wooden tabletop or a yellow tablecloth. Success rates are reported in the same way as Section 5.2.

**Results.** As shown in Table 4, DexGraspVLA achieves an aggregated generalization performance of 84.7% in the nonprehensile grasping task, demonstrating strong robustness to previously unseen object appearances, shapes, physical properties, as well as novel background and lighting conditions—significantly outperforming baselines. We observe that DexGraspVLA can reliably adapt to object poses, pushing until it extends sufficiently over the edge, followed by a stable grasp. This task is particularly challenging for parallel-jaw grippers, highlighting the dexterity we exhibit. Moreover, DexGraspVLA seamlessly extends to this new task without architectural changes, reflecting three key aspects of generality: (1) the high-level planner’s grounding and reasoning ability; (2) the use of bounding boxes as affordance guidance; and (3) applying imitation learning on domain-invariant representations iteratively obtained from foundation models.

## 6 Limitation and Conclusion

This paper presents DexGraspVLA, a hierarchical VLA framework for general-purpose dexterous grasping. By leveraging a pre-trained VLM as the high-level planner and a diffusion-based low-level controller, the system transforms diverse multimodal inputs into domain-invariant representations and learns robust closed-loop grasping policies via imitation learning. Our large-scale evaluations show over 90% success across thousands of unseen cluttered scenes in a zero-shot setting, with empirical evidence of strong generalization and consistent internal behavior. DexGraspVLA also handles free-form long-horizon prompts, recovers from failures, and extends to nonprehensile grasping, demonstrating broad applicability. While effective, it does not yet address functional grasping and subsequent manipulation, nor does it incorporate tactile sensing. In future work, we aim to extend the high-level planner to generate more fine-grained affordance and learn a task-oriented manipulation controller that also integrates tactile feedback, further broadening the scope of DexGraspVLA.

Table 4: DexGraspVLA significantly outperforms baselines in nonprehensile grasping tasks in diverse unseen conditions.

	Unseen Objects	Unseen Backgrounds	Unseen Lightings	Aggregated
ViT-small	61.1%	37.5%	22.2%	39.6%
DINOv2-train	66.7%	70.8%	55.6%	66.0%
Ours	<b>88.9%</b>	<b>86.1%</b>	<b>77.8%</b>	<b>84.7%</b>



Figure 6: Nonprehensile grasping.

## References

- [1] Tao Chen, Megha Tippur, Siyang Wu, Vikash Kumar, Edward Adelson, and Pulkit Agrawal. Visual dexterity: In-hand dexterous manipulation from depth. In *Icm workshop on new frontiers in learning, control, and dynamical systems*, 2023.
- [2] Haozhi Qi, Brent Yi, Sudharshan Suresh, Mike Lambeta, Yi Ma, Roberto Calandra, and Jitendra Malik. General in-hand object rotation with vision and touch. In *Conference on Robot Learning*, pages 2549–2564. PMLR, 2023.
- [3] Binghao Huang, Yuanpei Chen, Tianyu Wang, Yuzhe Qin, Yaodong Yang, Nikolay Atanasov, and Xiaolong Wang. Dynamic handover: Throw and catch with bimanual hands. In *7th Annual Conference on Robot Learning*, 2023.
- [4] Toru Lin, Zhao-Heng Yin, Haozhi Qi, Pieter Abbeel, and Jitendra Malik. Twisting lids off with two hands. *arXiv preprint arXiv:2403.02338*, 2024.
- [5] Yuanpei Chen, Tianhao Wu, Shengjie Wang, Xidong Feng, Jiechuan Jiang, Zongqing Lu, Stephen McAleer, Hao Dong, Song-Chun Zhu, and Yaodong Yang. Towards human-level bimanual dexterous manipulation with reinforcement learning. *Advances in Neural Information Processing Systems*, 35:5150–5163, 2022.
- [6] Kelvin Xu, Zheyuan Hu, Ria Doshi, Aaron Rovinsky, Vikash Kumar, Abhishek Gupta, and Sergey Levine. Dexterous manipulation from images: Autonomous real-world rl via substep guidance. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5938–5945. IEEE, 2023.
- [7] Yuanpei Chen, Chen Wang, Li Fei-Fei, and C Karen Liu. Sequential dexterity: Chaining dexterous policies for long-horizon manipulation. In *Conference on Robot Learning*, pages 3809–3829, 2023.
- [8] Abhishek Gupta, Justin Yu, Tony Z Zhao, Vikash Kumar, Aaron Rovinsky, Kelvin Xu, Thomas Devlin, and Sergey Levine. Reset-free reinforcement learning via multi-task learning: Learning dexterous manipulation behaviors without human intervention. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6664–6671. IEEE, 2021.
- [9] Kevin Zakka, Laura Smith, Nimrod Gileadi, Taylor Howell, Xue Bin Peng, Sumeet Singh, Yuval Tassa, Pete Florence, Andy Zeng, and Pieter Abbeel. Robopianist: A benchmark for high-dimensional robot control. *arXiv preprint arXiv:2304.04150*, 2023.
- [10] Sirui Chen, Jeannette Bohg, and C Karen Liu. Springgrasp: An optimization pipeline for robust and compliant dexterous pre-grasp synthesis. *arXiv preprint arXiv:2404.13532*, 2024.
- [11] Dylan Turpin, Tao Zhong, Shutong Zhang, Guanglei Zhu, Eric Heiden, Miles Macklin, Stavros Tsogkas, Sven Dickinson, and Animesh Garg. Fast-grasp’d: Dexterous multi-finger grasp generation through differentiable simulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8082–8089. IEEE, 2023.
- [12] Dylan Turpin, Liquan Wang, Eric Heiden, Yun-Chun Chen, Miles Macklin, Stavros Tsogkas, Sven Dickinson, and Animesh Garg. Grasp’d: Differentiable contact-rich grasp synthesis for multi-fingered hands. In *European Conference on Computer Vision*, pages 201–221. Springer, 2022.
- [13] Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- [14] Max Yang, Chenghua Lu, Alex Church, Yijiong Lin, Chris Ford, Haoran Li, Efi Psomopoulou, David AW Barton, and Nathan F Lepora. Anyrotate: Gravity-invariant in-hand object rotation with sim-to-real touch. *arXiv preprint arXiv:2405.07391*, 2024.
- [15] Johannes Pitz, Lennart Röstel, Leon Sievers, and Berthold Bäuml. Dextrous tactile in-hand manipulation using a modular reinforcement learning architecture. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1852–1858. IEEE, 2023.

- [16] Ankur Handa, Arthur Allshire, Viktor Makoviychuk, Aleksei Petrenko, Ritvik Singh, Jingzhou Liu, Denys Makoviichuk, Karl Van Wyk, Alexander Zhurkevich, Balakumar Sundaralingam, et al. Dextreme: Transfer of agile in-hand manipulation from simulation to reality. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5977–5984, 2023.
- [17] Chen Wang, Haochen Shi, Weizhuo Wang, Ruohan Zhang, Li Fei-Fei, and C Karen Liu. Dexcap: Scalable and portable mocap data collection system for dexterous manipulation. *arXiv preprint arXiv:2403.07788*, 2024.
- [18] Zoey Qiuyu Chen, Karl Van Wyk, Yu-Wei Chao, Wei Yang, Arsalan Mousavian, Abhishek Gupta, and Dieter Fox. Dextransfer: Real world multi-fingered dexterous grasping with minimal human demonstrations. *arXiv preprint arXiv:2209.14284*, 2022.
- [19] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. In *Robotics: Science and Systems*, 2017.
- [20] Maxime Oquab, Timothée Darcret, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2023.
- [21] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [22] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [23] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [24] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [25] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- [26] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- [27] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al.  $\pi_0$ : A Vision-Language-Action Flow Model for General Robot Control. *arXiv preprint arXiv:2410.24164*, 2024.
- [28] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. In *Robotics: Science and Systems*, 2024.
- [29] Abby O'Neill, Abdul Rehman, Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.

- [30] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [31] Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point trajectory modeling for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- [32] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky.  $\pi_{0.5}$ : a vision-language-action model with open-world generalization, 2025. URL <https://arxiv.org/abs/2504.16054>.
- [33] Wenzhuan Zhou and David Held. Learning to grasp the ungraspable with emergent extrinsic dexterity. In *Conference on Robot Learning*, pages 150–160. PMLR, 2023.
- [34] Yuhao Wang, Yu Li, Yaodong Yang, and Yuanpei Chen. Dexterous non-prehensile manipulation for ungraspable object via extrinsic dexterity. *arXiv preprint arXiv:2503.23120*, 2025.
- [35] Jens Lundell, Francesco Verdoja, and Ville Kyrki. Ddgc: Generative deep dexterous grasping in clutter. *IEEE Robotics and Automation Letters*, 6(4):6899–6906, 2021.
- [36] Andrew T Miller and Peter K Allen. Graspit! a versatile simulator for robotic grasping. *IEEE Robotics & Automation Magazine*, 11(4):110–122, 2004.
- [37] Tengyu Liu, Zeyu Liu, Ziyuan Jiao, Yixin Zhu, and Song-Chun Zhu. Synthesizing diverse and physically stable grasps with arbitrary hand structures using differentiable force closure estimator. *IEEE Robotics and Automation Letters*, 7(1):470–477, 2021.
- [38] Puhao Li, Tengyu Liu, Yuyang Li, Yiran Geng, Yixin Zhu, Yaodong Yang, and Siyuan Huang. Gendexgrasp: Generalizable dexterous grasping. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8068–8074. IEEE, 2023.
- [39] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzhen Xu, Puhao Li, Tengyu Liu, and He Wang. Dexgrasnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11359–11366. IEEE, 2023.
- [40] Jialiang Zhang, Haoran Liu, Danshi Li, XinQiang Yu, Haoran Geng, Yufei Ding, Jiayi Chen, and He Wang. Dexgrasnet 2.0: Learning generative dexterous grasping in large-scale synthetic cluttered scenes. In *8th Annual Conference on Robot Learning*, 2024.
- [41] Yiming Li, Wei Wei, Daheng Li, Peng Wang, Wanyi Li, and Jun Zhong. Hgc-net: Deep anthropomorphic hand grasping in clutter. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 714–720. IEEE, 2022.
- [42] Min Liu, Zherong Pan, Kai Xu, Kanishka Ganguly, and Dinesh Manocha. Deep differentiable grasp planner for high-dof grippers. In *Robotics: Science and Systems*, 2020.
- [43] Tao Chen, Jie Xu, and Pulkit Agrawal. A system for general in-hand object re-orientation. *Conference on Robot Learning*, 2021.
- [44] Zhao-Heng Yin, Binghao Huang, Yuzhe Qin, Qifeng Chen, and Xiaolong Wang. Rotating without seeing: Towards in-hand dexterity through touch. In *Robotics: Science and Systems*, 2023.
- [45] Haozhi Qi, Ashish Kumar, Roberto Calandra, Yi Ma, and Jitendra Malik. In-hand object rotation via rapid motor adaptation. In *Conference on Robot Learning*, pages 1722–1732. PMLR, 2023.
- [46] Sudeep Dasari, Abhinav Gupta, and Vikash Kumar. Learning dexterous manipulation from exemplar object trajectories and pre-grasps. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3889–3896. IEEE, 2023.

- [47] Gagan Khandate, Siqi Shang, Eric T Chang, Tristan Luca Saidi, Yang Liu, Seth Matthew Dennis, Johnson Adams, and Matei Ciocarlie. Sampling-based exploration for reinforcement learning of dexterous manipulation. In *Robotics: Science and Systems*, 2023.
- [48] Yuanpei Chen, Chen Wang, Yaodong Yang, and C Karen Liu. Object-centric dexterous manipulation from human motion data. In *8th Annual Conference on Robot Learning*, 2024.
- [49] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. Unidexgrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3891–3902, 2023.
- [50] Hui Zhang, Sammy Christen, Zicong Fan, Otmar Hilliges, and Jie Song. Graspxl: Generating grasping motions for diverse objects at scale. In *European Conference on Computer Vision*, pages 386–403. Springer, 2024.
- [51] Zhecheng Yuan, Tianming Wei, Shuiqi Cheng, Gu Zhang, Yuanpei Chen, and Huazhe Xu. Learning to manipulate anywhere: A visual generalizable framework for reinforcement learning. *CoRR*, abs/2407.15815, 2024. doi: 10.48550/ARXIV.2407.15815. URL <https://doi.org/10.48550/arXiv.2407.15815>.
- [52] Priyanka Mandikal and Kristen Grauman. Dexvip: Learning dexterous grasping with human hand pose priors from video. In *Conference on Robot Learning*, pages 651–661. PMLR, 2022.
- [53] Priyanka Mandikal and Kristen Grauman. Learning dexterous grasping with object-centric visual affordances. In *2021 IEEE international conference on robotics and automation (ICRA)*, pages 6169–6176. IEEE, 2021.
- [54] Tyler Ga Wei Lum, Martin Matak, Viktor Makoviychuk, Ankur Handa, Arthur Allshire, Tucker Hermans, Nathan D Ratliff, and Karl Van Wyk. Dextrah-g: Pixels-to-action dexterous arm-hand grasping with geometric fabrics. In *8th Annual Conference on Robot Learning*, 2024.
- [55] Ritvik Singh, Arthur Allshire, Ankur Handa, Nathan Ratliff, and Karl Van Wyk. Dextrah-rgb: Visuomotor policies to grasp anything with dexterous hands. *arXiv preprint arXiv:2412.01791*, 2024.
- [56] Zoey Qiuyu Chen, Karl Van Wyk, Yu-Wei Chao, Wei Yang, Arsalan Mousavian, Abhishek Gupta, and Dieter Fox. Learning robust real-world dexterous grasping policies via implicit shape augmentation. In *Conference on Robot Learning*, pages 1222–1232, 2022.
- [57] Ilija Radosavovic, Xiaolong Wang, Lerrel Pinto, and Jitendra Malik. State-only imitation learning for dexterous manipulation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7865–7871. IEEE, 2021.
- [58] Sridhar Pandian Arunachalam, Irmak Güzey, Soumith Chintala, and Lerrel Pinto. Holo-dex: Teaching dexterity with immersive mixed reality. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5962–5969. IEEE, 2023.
- [59] Irmak Guzey, Ben Evans, Soumith Chintala, and Lerrel Pinto. Dexterity from touch: Self-supervised pre-training of tactile representations with robotic play. In *7th Annual Conference on Robot Learning*, 2023.
- [60] Ankur Handa, Karl Van Wyk, Wei Yang, Jacky Liang, Yu-Wei Chao, Qian Wan, Stan Birchfield, Nathan Ratliff, and Dieter Fox. Dexpilot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170. IEEE, 2020.
- [61] Aravind Sivakumar, Kenneth Shaw, and Deepak Pathak. Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube. In *Robotics: Science and Systems*, 2022.
- [62] Yuzhe Qin, Wei Yang, Binghao Huang, Karl Van Wyk, Hao Su, Xiaolong Wang, Yu-Wei Chao, and Dieter Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. *arXiv preprint arXiv:2307.04577*, 2023.

- [63] Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *European Conference on Computer Vision*, pages 570–587. Springer, 2022.
- [64] Zichen Jeff Cui, Yibin Wang, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. From play to policy: Conditional behavior generation from uncurated robot data. In *The International Conference on Learning Representations*, 2023.
- [65] Siddhant Haldar, Jyothish Pari, Anant Rai, and Lerrel Pinto. Teach a robot to fish: Versatile imitation from one minute of demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.
- [66] Yuzhe Qin, Hao Su, and Xiaolong Wang. From one hand to multiple hands: Imitation learning for dexterous manipulation from single-camera teleoperation. *IEEE Robotics and Automation Letters*, 7(4):10873–10881, 2022.
- [67] Sridhar Pandian Arunachalam, Sneha Silwal, Ben Evans, and Lerrel Pinto. Dexterous imitation made easy: A learning-based framework for efficient dexterous manipulation. In *2023 IEEE international conference on robotics and automation (icra)*, pages 5954–5961, 2023.
- [68] Irmak Guzey, Yinlong Dai, Ben Evans, Soumith Chintala, and Lerrel Pinto. See to touch: Learning tactile dexterity through visual incentives. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 13825–13832. IEEE, 2024.
- [69] Toru Lin, Yu Zhang, Qiyang Li, Haozhi Qi, Brent Yi, Sergey Levine, and Jitendra Malik. Learning visuotactile skills with two multifingered hands. *arXiv preprint arXiv:2404.16823*, 2024.
- [70] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- [71] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10371–10381, 2024.
- [72] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2024.
- [73] Qwen Team. Qwen2.5-vl, January 2025. URL <https://qwenlm.github.io/blog/qwen2.5-vl/>.
- [74] Suneel Belkhale, Tianli Ding, Ted Xiao, Pierre Sermanet, Quon Vuong, Jonathan Tompson, Yevgen Chebotar, Debidatta Dwibedi, and Dorsa Sadigh. Rt-h: Action hierarchies using language. In <https://arxiv.org/abs/2403.01823>, 2024.
- [75] Jiayuan Gu, Sean Kirmani, Paul Wohlhart, Yao Lu, Montserrat Gonzalez Arenas, Kanishka Rao, Wenhao Yu, Chuyuan Fu, Keerthana Gopalakrishnan, Zhuo Xu, et al. Rt-trajectory: Robotic task generalization via hindsight trajectory sketches. *arXiv preprint arXiv:2311.01977*, 2023.
- [76] Mengda Xu, Zhenjia Xu, Yinghao Xu, Cheng Chi, Gordon Wetzstein, Manuela Veloso, and Shuran Song. Flow as the cross-domain manipulation interface. *arXiv preprint arXiv:2407.15208*, 2024.
- [77] Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xu Huang, Shu Jiang, et al. Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025.
- [78] Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Se June Joo, Jianwei Yang, Baolin Peng, Ajay Mandlekar, Reuben Tan, Yu-Wei Chao, Bill Yuchen Lin, et al. Latent action pretraining from videos. In *CoRL 2024 Workshop on Whole-body Control and Bimanual Manipulation: Applications in Humanoids and Beyond*.

- [79] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. In *Conference on Robot Learning*, pages 540–562, 2023.
- [80] Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang, and Li Fei-Fei. Rekep: Spatio-temporal reasoning of relational keypoint constraints for robotic manipulation. In *8th Annual Conference on Robot Learning*, 2024.
- [81] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- [82] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Joon-Young Lee, and Alexander Schwing. Putting the object back into video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3151–3161, 2024.
- [83] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- [84] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 2023.
- [85] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024.
- [86] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [87] Yiheng Li, Heyang Jiang, Akio Kodaira, Masayoshi Tomizuka, Kurt Keutzer, and Chenfeng Xu. Immiscible diffusion: Accelerating diffusion training with noise assignment. *arXiv preprint arXiv:2406.12303*, 2024.
- [88] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2020.
- [89] Andreas Peter Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *Transactions on Machine Learning Research*, 2022.

## A Implementation Details

In this section, we present the details of DexGraspVLA implementation (Appendix A.1), base-line implementation (Appendix A.2), dataset collection (Appendix A.3), and assets information (Appendix A.4). Our code, trained model weights, and videos can be found at [dexgraspvla.github.io](https://dexgraspvla.github.io).

### A.1 Details of DexGraspVLA Implementation

**Planner.** The high-level planner operates as described in Section 4.1. By leveraging an off-the-shelf VLM as the planner, our framework gains remarkable flexibility, enabling easy utilization of more advanced models for enhanced performance. Our observations indicate that Qwen2.5-VL-72B-Instruct [73] outperforms Qwen-VL-Chat [81] in reasoning and instruction following, leading to improved long-horizon task completion. Therefore, we base the DexGraspVLA planner on Qwen2.5-VL-72B-Instruct in the long-horizon tasks and provide our prompts below.

These prompts mainly instruct the VLM to function as DexGraspVLA planner via four sub-tasks, including (1) Instruction Proposal: proposing the current grasping instruction  $l$  based on the user prompt  $p$ , (2) Bounding Box Prediction: marking the target object bounding box, (3) Grasp Outcome Verification: checking if the grasp has succeeded, and (4) Prompt Completion Check: evaluating whether the entire user prompt is fully fulfilled. Since instruction proposal, bounding box prediction, and prompt completion check only require information within the operational workspace on the table, we crop the relevant region from the head camera image and fill the remaining area with white pixels. The resulting cropped image is used as the planner’s visual input for these sub-tasks.

To start with, when a user prompt  $p$  is provided, the planner first determines which object in the scene should be grasped next. This step involves interpreting the prompt in context and selecting the best matching object from the current visual input.

You are controlling a robotic arm that needs to complete the following user prompt:  
<user\_prompt>.

I will show you two images. The initial image (before any actions) is: <initial\_head\_image>. The current image (after the latest action) is: <current\_head\_image>.

Your task is to select the **best object to grasp next** from the current image.

To identify objects, **use common sense and everyday knowledge** to infer what each item is. For example, recognize cups, bottles, fruits, snacks, boxes, tools, etc.

When choosing the best object to grasp, follow these principles:

1. Prefer objects on the right, then center, then left.
2. Avoid objects that are blocked or surrounded.
3. Avoid grasping objects that would cause other items to topple.
4. Select objects that best match the user prompt.

Please output ONLY ONE object that the robot should grasp next.

Return format (in English, natural language):

A short sentence precisely describing the target object, including:

- color.
- shape.
- relative position (e.g., "on the right", "in front", "next to the red box").

Example:

Grasp the blue cube on the right side of the table.

After deciding on the next object to grasp, the planner proceeds to locate this object in the image by predicting its bounding box using the following prompts. The generated grasping instruction is used as input to this localization module.

You are a robotic vision assistant. Your task is to locate the object described below in the given image: <current\_head\_image> and return its bounding box.

Grasping instruction: <grasping\_instruction>.

Instructions:

1. Carefully read the grasping instruction and match the target object to the best-fitting visible object in the image.
2. Select EXACTLY ONE object that best matches the description.
3. For the selected object, return the following in strict JSON format:
  - "bbox\_2d": [x1, y1, x2, y2] (integer pixel coordinates, top-left to bottom-right)
  - "label": a short 2-4 word name, (e.g. "blue cup")
  - "description": a complete, natural-language description of the object's appearance and position

Requirements:

- Only return one object.
- Coordinates must be valid and within image boundaries.
- Do not guess if the object is not visible.

During the controller's execution, the planner verifies whether the object has been successfully grasped, using the following prompt.

I will show you two images. The top-down view from the head camera is: <current\_head\_image>. The close-up view from the wrist camera is: <current\_wrist\_image>.

Grasping instruction: <grasping\_instruction>.

Task:

Determine whether the robotic arm has **successfully grasped the target object**.

You should consider:

- Whether the target object is still visible on the table.
- Whether the object is securely held in the robotic hand.

Output format:

A reasoning and a boolean value (True=successfully grasped, False=not grasped).

Keep it short and simple.

After each grasp attempt, the planner checks whether the user prompt has been fulfilled with the following prompt.

The robot is trying to complete the following user prompt: <user\_prompt>.

I will show you two images. The initial image (before any actions) is: <initial\_head\_image>. The current image (after the latest action) is: <current\_head\_image>.

Please compare the two images and determine whether the user prompt has been fully completed.

Instructions:

- Only consider visible 3D objects.
- If all target objects have been removed or grasped, return True.
- If some relevant objects remain, return False.

Output format:

A reasoning and a boolean value (True=completed, False=not completed).

Example:

All blue objects have been removed from the table: True.

In our experiments, we either query the online APIs of these models or host them on an 8-A800 GPU server by ourselves with vLLM [86]. When hosting Qwen2.5-VL-72B-Instruct, we employ Qwen2.5-VL-7B-Instruct for speculative decoding to accelerate inference.

**Controller.** We first elaborate on the implementation details for the controller in the general dexterous grasping experiments. All raw images are produced by head and wrist cameras at a resolution of  $640 \times 480 \times 3$ . Correspondingly, the resolution of mask is  $640 \times 480 \times 1$ . Through preliminary model selection, we decide to use DINOv2 ViT-B/14 as the feature extractor  $\phi^h$  for head camera images and DINOv2 ViT-L/14 as the feature extractor  $\phi^w$  for wrist camera images. Before feeding images into DINOv2, we resize them to  $518 \times 518 \times 3$ . During training, we apply domain randomization via color jittering. Finally, the images are normalized and fed into DINOv2 models. This leads to features  $\mathbf{z}_t^h \in \mathbb{R}^{1369 \times 768}$  and  $\mathbf{z}_t^w \in \mathbb{R}^{1369 \times 1024}$ . By processing the mask  $\mathbf{m}_t$  with a randomly initialized ViT, we extract its features  $\mathbf{z}_t^m \in \mathbb{R}^{1369 \times 768}$ . Patch-wise concatenation of  $\mathbf{z}_t^h$  and  $\mathbf{z}_t^m$  leads to  $\tilde{\mathbf{z}}_t^h \in \mathbb{R}^{1369 \times 1536}$ . We then project  $\tilde{\mathbf{z}}_t^h, \mathbf{z}_t^w, \mathbf{s}_t$  to the same feature space of dimension 1024 with separate MLPs, yielding  $\tilde{\mathbf{z}}_t^h \in \mathbb{R}^{1369 \times 1024}, \tilde{\mathbf{z}}_t^w \in \mathbb{R}^{1369 \times 1024}, \tilde{\mathbf{z}}_t^s \in \mathbb{R}^{1 \times 1024}$ , and concatenate them to form the full observation feature sequence  $\tilde{\mathbf{z}}_t^{\text{obs}} = (\tilde{\mathbf{z}}_t^h, \tilde{\mathbf{z}}_t^w, \tilde{\mathbf{z}}_t^s) \in \mathbb{R}^{2739 \times 1024}$ .

For action modeling, we define an action chunk horizon of  $H = 64$ . When we add noise to the action during training, we employ Immiscible Diffusion [87] to improve data-noise mapping. The noised action chunk  $\mathbf{A}_t$  belongs to  $\mathbb{R}^{64 \times 13}$ .

The DiT implementation is based on the original DiT paper [83], diffusion policy [84], and RDT [85]. It first embeds the diffusion timestep to the same hidden space as  $\tilde{\mathbf{z}}_t^{\text{obs}}$ , yielding  $\tilde{\mathbf{z}}_t^d \in \mathbb{R}^{1 \times 1024}$ , and concatenates it with  $\tilde{\mathbf{z}}_t^{\text{obs}}$  to form the condition sequence  $\tilde{\mathbf{z}}_t = (\tilde{\mathbf{z}}_t^{\text{obs}}, \tilde{\mathbf{z}}_t^d) \in \mathbb{R}^{2740 \times 1024}$ . We project the noised action chunk to the same hidden space, deriving  $\tilde{\mathbf{z}}_t^A \in \mathbb{R}^{64 \times 1024}$ , and feed it into DiT. Each DiT layer performs bi-directional attention within action tokens, cross-attention to the condition sequence, and MLP projections. Finally, the output is projected back to the action space to be the model’s prediction of noise. During training, we compute MSE loss between the noise prediction and ground truth, and back-propagate the gradient to update all trainable parameters. During inference, we start from Gaussian noise and iteratively denoise it using DDIM sampling [88]. At each step, the DiT model predicts the noise given the condition sequence, and we update the action chunk using the DDIM scheduler until we obtain the final action. The controller only executes the first six actions in the predicted action chunk before making a new prediction.

In total, the controller possesses 163M trainable parameters. To accelerate training, we utilize bfloat16 mixed-precision training, reducing memory usage and improving computational efficiency. Additionally, we employ FusedAdamW as the optimizer to further speed up training through optimized memory access and fused kernel execution. With these techniques, we train the controller for 84 epochs over our dataset on an 8-A800 GPU server, which takes less than one day to complete. All hyper-parameters in our implementation are presented in Table 5.

In the nonprehensile grasping experiments, we keep most of the hyper-parameters the same but make the following changes: we use DINOv2 ViT-B/14 as the feature extractors  $\phi^h, \phi^w$  for both head and wrist camera images, and the action horizon is set to 100. This controller has 106M trainable parameters and is trained for 200 epochs on an 8-A800 GPU server, which takes approximately two days to finish.

## A.2 Details of Baseline Implementation

In both general dexterous grasping and nonprehensile grasping experiments, the baseline DexGraspVLA (DINOv2-train) is the same as DexGraspVLA (Ours) described in Appendix A.1 except that the two DINOv2 models are trainable instead of frozen. The baseline DexGraspVLA (ViT-small) is the same as DexGraspVLA (Ours) except that the two DINOv2 models are replaced with two small trainable pre-trained ViTs (the R26-S-32 ResNet-ViT hybrid from Steiner et al. [89]). Correspondingly, we resize the images to  $224 \times 224 \times 3$  to feed them into ViT-small. Each image is split into 49 patches, and the feature dimension is 384.

Table 5: Hyper-parameters of DexGraspVLA in general dexterous grasping experiments.

Hyper-parameter	Value	Hyper-parameter	Value
epoch	84	attention dropout	0.1
learning rate	0.0001	noise scheduler	DDIMScheduler
learning rate scheduler	cosine	num_train_timesteps	50
learning rate warmup steps	2000	beta_start	0.0001
weight decay	0.0001	beta_end	0.02
AdamW betas	[0.95, 0.999]	beta_schedule	squaredcos_cap_v2
seed	42	clip_sample	True
batch size per GPU	48	set_alpha_to_one	True
action horizon	64	steps_offset	0
number of DiT layers	12	prediction_type	epsilon
number of DiT head	8	num_inference_steps	16

### A.3 Details of Data Collection

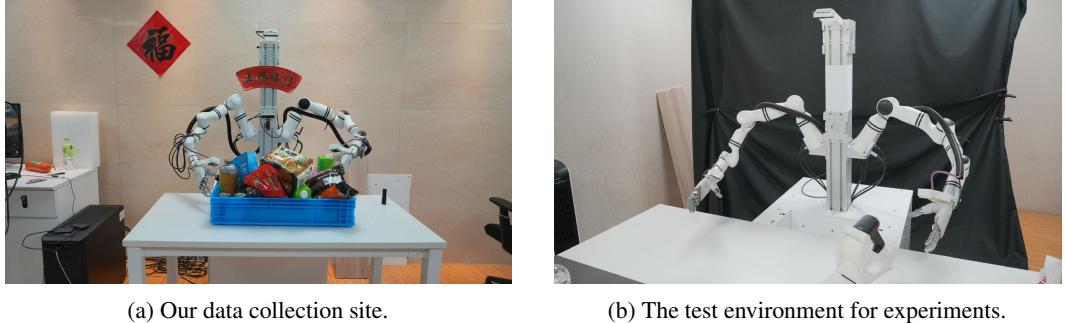
We collect demonstrations through kinesthetic teaching. At the beginning, the robot is set to teaching mode, allowing manual guidance to grasp target objects. The operator then physically guides the robot to the target position and performs the grasping motion. Subsequently, we reset the environment and execute PD control using the recorded joint angles as target. At the same frequency, these target joint angles serve as actions, while images and current joint angles are collected as states. Following the same approach as the low-level controller, we post-process the collected data to generate masks, completing one demonstration sequence. In the general dexterous grasping experiments, each episode has a fixed duration of 75 timesteps, while in nonprehensile grasping, demonstrations have variable lengths, depending on the amount of manipulation required to push the object toward the table edge and complete the grasp. The control frequency is 20Hz.

We hire external contractors, provide them with training, and engage them to assist with data collection. All contractors were compensated with fair wages.

### A.4 Information of Assets

We list the information of assets as below:

1. DINOv2 ViT-B/14 and ViT-L/14 [20]
  - License: Apache License 2.0
  - URL: <https://github.com/facebookresearch/dinov2>
2. Qwen-VL-Chat [81]
  - License: Tongyi Qianwen LICENSE AGREEMENT (<https://github.com/QwenLM/Qwen-VL/blob/master/LICENSE>)
  - URL: <https://huggingface.co/Qwen/Qwen-VL-Chat>
3. Qwen2.5-VL-72B-Instruct and Qwen2.5-VL-7B-Instruct [73]
  - License: Apache-2.0 license
  - URL: <https://github.com/QwenLM/Qwen2.5-VL>
4. vit\_small\_r26\_s32\_224 [89]
  - License: Apache-2.0 license
  - URL: <https://github.com/huggingface/pytorch-image-models>



(a) Our data collection site.

(b) The test environment for experiments.

Figure 7: A comparison of the data collection and test environments, located in different rooms. The visual scenes captured by the robot’s cameras differ significantly, especially for the wrist camera.

## B Experiment Details

### B.1 The “Zero-Shot” Evaluation Environment

Figure 7 contrasts our data collection site and the test site, which are located in separate rooms. We gather all human demonstrations at the data collection site (Figure 7a), whereas the experiments in Section 5 are conducted at the test site (Figure 7b). Because these sites differ in layout and background, both the head camera and the wrist camera encounter scenes not present in the training data during evaluation — particularly the wrist camera, which observes a notably altered environment, capturing a variety of front and peripheral views during operation. Despite these environmental discrepancies, we do not collect any data from the test site to fine-tune the models. Instead, the models are deployed and evaluated directly, resulting in a genuinely “zero-shot” testing environment. Even under these conditions, DexGraspVLA achieves an over 90% success rate in grasping tasks in cluttered scenes across thousands of unseen object, lighting, and background combinations, clearly demonstrating its strong generalization capability.

### B.2 Additional Details of Objects, Lightings, and Backgrounds in General Dextrous Grasping

We collect a total of 360 unseen objects, from which 103 items are randomly selected as the *object subset  $S$* . In the main paper, the *Unseen Objects* experiment is conducted on all 360 objects, while the *Unseen Lightings* and *Unseen Backgrounds* experiments use only the objects in  $S$ . The three unseen lighting conditions comprise disco light, lamp light, and dark light. Meanwhile, the six unseen backgrounds include a black mouse pad, a pink towel, a colorful tablecloth, a black-and-white mouse pad, a wooden board, and a calibration board. These conditions are illustrated in Figure 8.

### B.3 Additional Details of Objects, Lightings, and Backgrounds in Nonprehensile Grasping

In Figure 9, we present the 32 objects curated for collecting nonprehensile grasping demonstrations and 18 unseen objects used for evaluation, covering a wide range of appearances, geometries, sizes, and categories. In Figure 8, we show the unseen background and lighting conditions used in the generalization evaluation. DexGraspVLA demonstrates robust performance on challenging cases, including fully white or irregularly shaped objects. In these scenarios, it successfully pushes the objects toward the table edge to enable stable grasping, even under complex and unseen lighting and background conditions.

### B.4 Details of Visualization

In this part, we explain how we visualize the internal model behavior shown in Figure 5. Due to space constraints, Figure 5 only presents the relevant portion of images containing the tabletop workspace. The full version is shown in Figure 10. The first row is raw images from the head camera resized to  $518 \times 518 \times 3$ . The second row illustrates the DINOv2 ViT-B/14 features following the practice introduced in DINOv2 paper [20]. To make the resulting feature map recognizable for visualization purpose, we enlarge both the height and weight of images by a factor of six before feeding them into DINOv2. After obtaining the feature sequences for all four images, we combine these features, perform a PCA between all patches, and set a threshold to remove background regions. We then apply PCA again, this time to the remaining foreground features, map the top three principal components to the RGB channels, and normalize the result. This yields the visualization shown in the second row. The third row showcases the binary masks  $\mathbf{m}_t \in \mathbb{R}^{518 \times 518 \times 1}$  tracked by Cutie. The fourth row

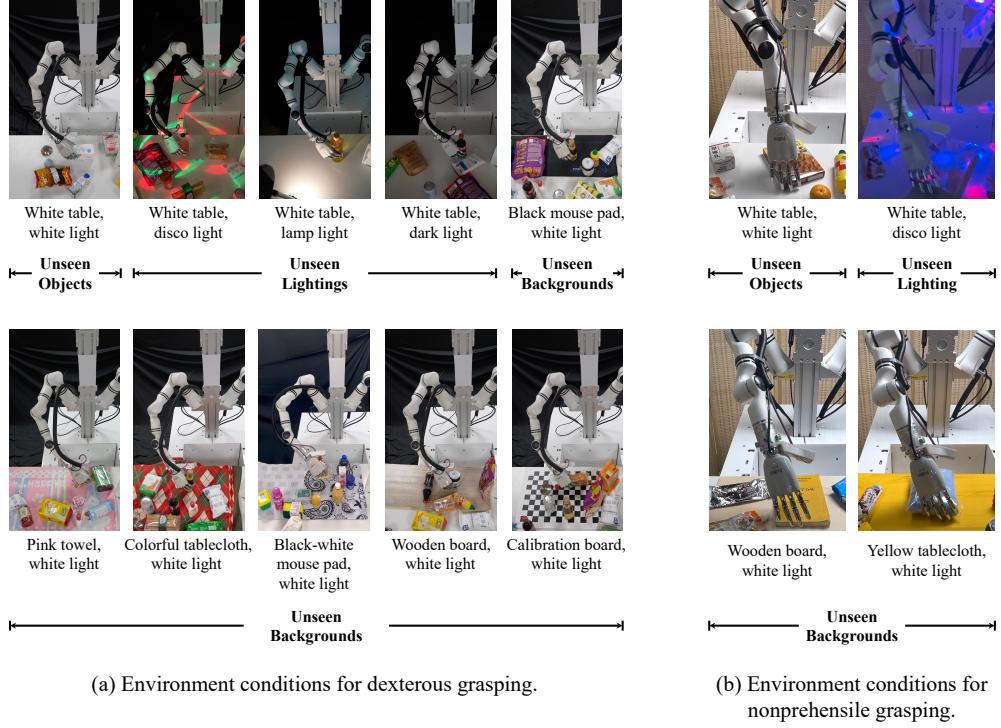


Figure 8: Environment conditions used in our generalization evaluations of dexterous grasping (Section 5.2) and nonprehensile grasping (Section 5.6).



Figure 9: Objects used to train and test methods in nonprehensile grasping. DexGraspVLA achieves robust generalization performance on diverse unseen objects.

displays the averaged DiT attention maps over the head image features. This is computed by summing attention weights to each head image patch across all diffusion steps, DiT layers, DiT heads, and action tokens, and normalize the sum to one. The shape of the averaged attention map is  $37 \times 37 \times 1$ . Finally, we upsample the attention map to  $518 \times 518 \times 1$ , multiply it by 2 to increase brightness, and use it to scale the value channel of head images in HSV space, resulting in the visualization shown in the fifth row.

## C Additional Results

This section provides additional results for the experiments in the main paper. In Table 6, we report the detailed success rates for our large-scale generalization evaluation under each environment condition, corresponding to Table 1 in Section 5.2. From the first row (“Ours@1”), it is evident that DexGraspVLA maintains consistently high success rates across various unseen object, lighting, and background combinations. Many observed failures stem from randomness in policy inference; allowing additional attempts often recovers these failed cases. Accordingly, the second and third rows

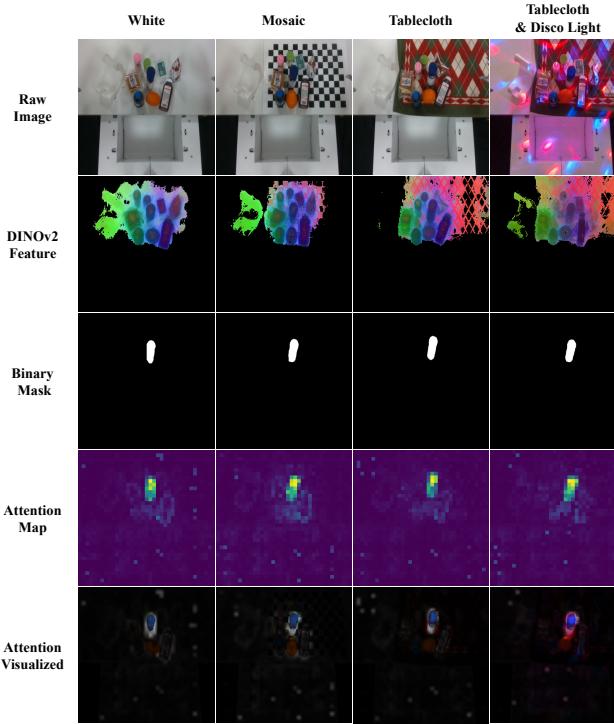


Figure 10: The complete, uncropped version of Figure 5.

Table 6: The detailed performance of DexGraspVLA under different unseen conditions, which indicates that our approach consistently achieves high success rates across various objects, lightings, and backgrounds. The second and third rows highlight its potential to reach even higher success rates given more chances.

Tasks	Unseen Objects (360)		Unseen Lightings (3 × 103)			Unseen Backgrounds (6 × 103)						Aggregated (1287)
	White Light	Disco Light	Lamp Light	Dark Light	White Light	White Light	White Light	White Light	White Light	White Light	White Light	
Lighting Conditions	White Light	Disco Light	Lamp Light	Dark Light	White Light	White Light	White Light	White Light	White Light	White Light	White Light	
Background Conditions	White Table	White Table	White Table	White Table	Black Mouse Pad	Pink Towel	Colorful Tablecloth	Black-White Mouse Pad	Wooden Board	Calibration Board		
Ours@1	91.1%	92.2%	89.3%	91.2%	94.2%	84.5%	90.3%	92.2%	93.2%	88.3%	90.8%	
Ours@2	95.3%	97.0%	95.1%	93.2%	97.1%	90.3%	91.3%	95.1%	98.1%	93.2%	94.7%	
Ours@3	96.7%	98.1%	98.1%	96.1%	98.1%	91.3%	94.2%	98.1%	100.0%	98.1%	96.9%	

(“Ours@2” and “Ours@3”) show further improvements in performance, highlighting the potential for DexGraspVLA to reach even higher success rates.

We next present the generated responses of DexGraspVLA planner in an exemplar long-horizon task, which demonstrates visual understanding and embodied reasoning. The cluttered scene is arranged as shown in Figure 11, with the user prompt being “Grasp all edible objects, including food and drinks”. Below is the generated responses of DexGraspVLA planner. We omit our system prompts as they have been elaborated in Appendix A.1.

```
<System Prompt of Prompt Completion Check>
.....
<Planner Response>
To determine if the user prompt has been fully completed, let's identify all the edible objects in the initial image and check if they have been removed or grasped in the current image.
```



Figure 11: The head camera images observed by DexGraspVLA planner and generated masks for a long-horizon task: "Grasp all edible objects, including food and drinks."

### Edible Objects in the Initial Image:

1. A bottle (drink).
2. A green packet with an orange label (food).
3. A small white object with a heart symbol (possibly food).
4. A packaged item with a brown appearance (food).

### Analysis of Current Image:

- The bottle is still present on the table.
- The green packet with an orange label is still present on the table.
- The small white object with a heart symbol is still present on the table.
- The packaged item with a brown appearance is still present on the table.

Since all the edible objects are still visible on the table in the current image, none of them have been removed or grasped.

**Reasoning:** All target objects (edible items) remain on the table, indicating that the user prompt to grasp all edible objects has not been completed.

**Output:** False

<System Prompt of Instruction Proposal>

.....

<Planner Response>

The green and orange snack bag on the right side of the table.

<System Prompt of Bounding Box Prediction>

.....

<Planner Response>

```
{
    "bbox_2d": [357, 104, 486, 249],
    "label": "green orange snack",
    "description": "The green and orange snack bag on the right
    side of the table."
}
```

<System Prompt of Grasp Outcome Verification>

.....

<Planner Response>

The target object is no longer visible on the table and appears to be securely held in the robotic hand.

True

<System Prompt of Prompt Completion Check>

.....

<Planner Response>

To determine if the user prompt has been fully completed, let's identify all the edible objects in the initial image and check if they have been removed or grasped in the current image.

### Initial Image Edible Objects:

1. A bottle (drink).
2. A green packet with a brown label (food).
3. A small packaged food item with a white label (food).

### Current Image Analysis:

- The bottle is still present on the table.
- The green packet with a brown label is no longer visible.
- The small packaged food item with a white label is still present on the table.

### Reasoning: The user prompt was to grasp all edible objects. However, the bottle and the small packaged food item are still on the table, indicating that not all edible objects have been removed or grasped.

### Conclusion: Since some relevant objects remain on the table, the user prompt has not been fully completed.

\*\*Output:\*\* All edible objects have been removed or grasped: False.

<System Prompt of Instruction Proposal>

.....

<Planner Response>

The brown snack packet on the right side of the table.

<System Prompt of Bounding Box Prediction>

.....

<Planner Response>

```
{  
    "bbox_2d": [450, 168, 519, 263],  
    "label": "brown snack packet",  
    "description": "The brown snack packet on the right side of the  
    table"  
}
```

<System Prompt of Grasp Outcome Verification>

.....

<Planner Response>

The target object is not visible on the table, and it appears to be securely held in the robotic hand based on the close-up view.

Reasoning: The object is no longer on the table and is being held by the robotic arm.

True

<System Prompt of Prompt Completion Check>

.....

<Planner Response>

The user prompt is to grasp all edible objects, including food and drinks. In the initial image, there are several items that could be considered edible:

1. A bottle of liquid (drink).
2. A bag of snacks.
3. A packaged food item.

In the current image, the bottle of liquid remains on the table, while the other two items (the bag of snacks and the packaged food) are no longer visible. This suggests that only some of the edible objects have been grasped or removed.

Reasoning: Since the bottle of liquid is still present on the table, not all edible objects have been grasped.

Output: False

<System Prompt of Instruction Proposal>  
.....

<Planner Response>  
The white bottle with a brown label next to the yellow toy on the right side of the table.

<System Prompt of Bounding Box Prediction>  
.....

<Planner Response>

```
{
    "bbox_2d": [304, 150, 368, 227],
    "label": "white bottle",
    "description": "The white bottle with a brown label is located next to the yellow toy on the right side of the table."
}
```

<System Prompt of Grasp Outcome Verification>  
.....

<Planner Response>  
The target object is no longer visible on the table and appears to be securely held in the robotic hand.

True

<System Prompt of Prompt Completion Check>  
.....

<Planner Response>  
The user prompt is to grasp all edible objects, including food and drinks. In the initial image, there are several items that could be considered edible:

1. A bottle of liquid (possibly a drink).
2. A packaged snack with visible text.
3. Another packaged item that appears to be food.

In the current image, these items are no longer visible on the table. This suggests they have been removed or grasped by the robot.

Reasoning: Since all the potential edible objects from the initial image are no longer present in the current image, it can be concluded that the user prompt has been fully completed.

Output: True

## D Additional Experiments

This section presents an additional experiment that aims to separately evaluate the bounding-box prediction accuracy of the DexGraspVLA planner.

**Tasks.** The bounding-box prediction accuracy of the planner is crucial to the success of grasping, as it determines the target for the controller. To evaluate this accuracy, we design three types of tasks featuring different environmental distractions: (1) *No Distraction* (1 scenario): The cluttered scene is arranged on a white table under white light; (2) *Background Distraction* (2 scenarios): The cluttered scene is placed on either a calibration board or a brightly colored tablecloth, both under white light; (3) *Lighting Distraction* (2 scenarios): The scene is set up in a dark room illuminated by either a desk lamp or a disco light. Scenarios with distractions are shown in Figure 12. For each scenario, we randomly arrange five cluttered scenes, each containing six randomly selected objects, and then record head-camera images. For each object, we provide a textual prompt describing its appearance

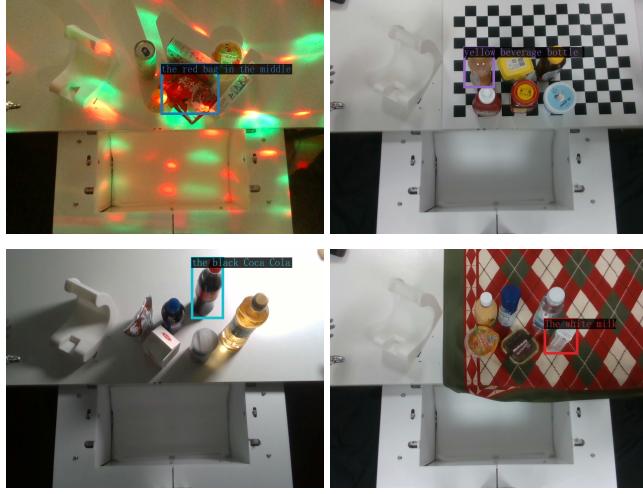


Figure 12: Bounding-box predictions made by DexGraspVLA planner. Across diverse lighting and background conditions, it accurately grounds the language instruction to the target object in cluttered scenes and marks the correct bounding box.

and location, and check whether the planner’s bounding-box prediction accurately marks the target. In total, *No Distraction* accounts for 30 tests, while *Background Distraction* and *Lighting Distraction* both have 60 tests, amounting to 150 tests overall.

**Metric.** We define a bounding box as accurate if it tightly encloses the target object. Accuracy is then measured as the proportion of accurate bounding boxes over all tested objects.

**Results.** The accuracy is reported in Table 7. For 150 prompts, the planner only mislabels one bounding box while succeeding in the other 149 tests, resulting in an aggregated accuracy exceeding 99%. In Figure 12, we present examples of bounding-box predictions produced by the DexGraspVLA planner. Despite substantial variation in environmental conditions, the planner consistently grounds grasping instructions in cluttered scenes and provides the correct bounding boxes. Notably, we can identify objects by names such as “Coca Cola” or “milk,” reflecting the system’s extensive common sense and world knowledge. By drawing on the broad knowledge embedded in each of its foundation models, DexGraspVLA achieves robust generalization across diverse scenarios.

## E Broader Impacts

DexGraspVLA demonstrates strong generalization for dexterous grasping in unseen cluttered scenes, achieving high success rates. It robustly handles adversarial objects, human disturbances, failure recovery, and free-form long-horizon grasping prompts. Furthermore, we show that the framework can extend to additional tasks such as nonprehensile grasping, highlighting its potential for broader application. In many real-world robotics scenarios, such capabilities can serve as a foundation for downstream manipulation, thereby improving productivity and benefiting society. The generality of both the method and the framework also supports broader use in task learning and deployment, contributing to the advancement of the field.

The main potential negative impact relates to safety concerns. As with any autonomous manipulation system, erroneous model predictions could damage objects or, in the worst case, pose a risk to nearby humans. However, we have thoroughly validated the high success rate of our approach, which mitigates such risks. In practice, safety can be further ensured through reduced execution speed, emergency stopping mechanisms, and power-cut safeguards during deployment.

Table 7: Planner accuracy in bounding-box prediction under different environment conditions.

	No Distraction	Background Distraction	Lighting Distraction	Aggregated
Planner	96.7%	100.0%	100.0%	99.3%