

# RoboAug: One Annotation to Hundreds of Scenes via Region-Contrastive Data Augmentation for Robotic Manipulation

Xinhua Wang<sup>1,\*</sup>, Kun Wu<sup>1,\*</sup>, Zhen Zhao<sup>1,\*</sup>, Hu Cao<sup>2</sup>, Yinuo Zhao<sup>1,3</sup>, Zhiyuan Xu<sup>1</sup>, Meng Li<sup>1</sup>, Shichao Fan<sup>1,4</sup>, Di Wu<sup>1,5</sup>, Yixue Zhang<sup>1,6</sup>, Ning Liu<sup>1</sup>, Zhengping Che<sup>1,†,✉</sup> and Jian Tang<sup>1,✉</sup>

<sup>1</sup>Beijing Innovation Center of Humanoid Robotics

<sup>2</sup>Computation, Information and Technology, Technical University of Munich

<sup>3</sup>City University of Hong Kong

<sup>4</sup>The School of Mechanical Engineering and Automation, Beihang University

<sup>5</sup>State Key Laboratory of Multimedia Information Processing, School of Computer Science, Peking University

<sup>6</sup>The School of Advanced Manufacturing and Robotics, Peking University

\*Co-first authors; †Project leader; ✉Corresponding authors.



Fig. 1: We introduce RoboAug, a region-contrastive data augmentation framework. RoboAug enables robust robotic generalization in diverse, unseen scenes.

**Abstract**—Enhancing the generalization capability of robotic learning to enable robots to operate effectively in diverse, unseen scenes is a fundamental and challenging problem. Existing approaches often depend on pretraining with large-scale data collection, which is labor-intensive and time-consuming, or on semantic data augmentation techniques that necessitate an impractical assumption of flawless upstream object detection in real-world scenarios. In this work, we propose RoboAug, a novel generative data augmentation framework that significantly minimizes the reliance on large-scale pretraining and the perfect visual recognition assumption by requiring only the bounding

box annotation of a single image during training. Leveraging this minimal information, RoboAug employs pre-trained generative models for precise semantic data augmentation and integrates a plug-and-play region-contrastive loss to help models focus on task-relevant regions, thereby improving generalization and boosting task success rates. We conduct extensive real-world experiments on three robots, namely UR-5e, AgileX, and Tien Kung 2.0, spanning over 35k rollouts. Empirical results demonstrate that RoboAug significantly outperforms state-of-the-art data augmentation baselines. Specifically, when evaluating generalization capabilities in unseen scenes featuring diverse combinations of

**backgrounds, distractors, and lighting conditions, our method achieves substantial gains over the baseline without augmentation. The success rates increase from 0.09 to 0.47 on UR-5e, from 0.16 to 0.60 on AgileX, and from 0.19 to 0.67 on Tien Kung 2.0. These results highlight the superior generalization and effectiveness of RoboAug in real-world manipulation tasks. Our project is available at <https://x-roboaug.github.io/>.**

## I. INTRODUCTION

The deployment of generalist robots in unstructured, real-world environments requires a level of perceptual robustness that extends far beyond controlled training conditions. While end-to-end visuomotor policies have demonstrated impressive capabilities in learning complex skills [5, 80, 77, 41, 4, 18], they remain notoriously brittle to distribution shifts of the observations. When deployed in unseen scenes, the performance of these policies often degrades significantly due to environmental interferences. Addressing this fragility is crucial for realizing reliable robotic systems. In this work, we focus on enhancing policy generalization against three predominant sources of out-of-distribution (OOD) interference: *complex background variations, drastic lighting changes, and the presence of task-irrelevant distractors*.

To tackle the generalization challenge, two primary paradigms have emerged: scaling real-world data collection and leveraging synthetic data augmentation. Inspired by scaling laws in foundation models [43, 68, 49], the first approach advocates for pretraining on massive datasets [45, 64, 27, 66]. However, unlike the passive, internet-scale data acquisition feasible for text and image domains, collecting robotic demonstration data in the real world is prohibitively expensive and labor-intensive. Consequently, Data Augmentation (DA) [33, 25, 72, 35] has become a vital alternative. Traditional “weak” augmentation techniques, such as random cropping and color jittering, modify low-level pixel statistics but fail to introduce the semantic diversity necessary to bridge the gap between training and unstructured deployment environments.

Recent advances [11, 58, 12, 73, 78] have thus pivoted toward “strong” semantic augmentation, utilizing generative models [48, 51] to synthesize novel visual contexts via inpainting. Crucially, these methods rely on the assumption that task-relevant entities, such as the robot and manipulated objects, can be precisely isolated using off-the-shelf segmentation [49] or detection models [44]. However, as highlighted in RoboEngine [73] and corroborated by our empirical analysis, this assumption is often overly optimistic. We collected a dataset, **RoboAug-D**, which contains 7,576 trajectories across 33 tasks for object detection. Then we evaluated state-of-the-art models like GroundingDINO [37] and LLMDet [22] and found substantial failure modes. Imprecise extraction, characterized by missing boundaries or hallucinated regions, propagates to the generative process. For instance, failing to detect a target object causes it to be overwritten by background textures during inpainting, leading the policy to learn incorrect behaviors, such as grasping empty space. This limitation prevents existing pipelines from synthesizing the high-quality data required to immunize policies against real-world interference.

To overcome these deficiencies, we propose **RoboAug**, a novel Region-Contrastive Data Augmentation Framework designed to achieve robust generalization with minimal human intervention. RoboAug synergizes three key technical phases: (1) robust task-relevant region extraction, (2) semantic data augmentation, and (3) region-contrastive policy learning. First, we introduce a task-relevant region extraction phase that generates semantic masks across all trajectory images using annotations from only a single frame. Unlike prior methods requiring labor-intensive frame-by-frame labeling or detector re-training [73], we leverage a one-shot region matching strategy in a training-free manner. By combining GroundingDINO for initial proposals, DINOv2 [44] for category correspondence, and SAM2 [49] for temporal tracking, we ensure precise, pixel-level extraction of task-relevant entities.

Building on these high-quality masks, RoboAug employs a semantic data augmentation phase. Instead of relying on unstable inpainting [60], we directly synthesize diverse full-scene backgrounds [71] and seamlessly composite the foreground regions onto them. To fully exploit the semantic information provided by the masks, we further integrate a region-contrastive policy learning objective. This objective introduces a contrastive loss directly into the visual encoder without architectural modifications, promoting the clustering of feature representations within the same semantic class while repelling disparate classes. This enhances the policy’s ability to attend to task-relevant objects against visual interference.

We validate our framework through a comprehensive experimental campaign comprising over **35k** real-world trials on Tien Kung 2.0, UR-5e, and AgileX robots. Our evaluation rigorously decouples environmental variables, testing background shifts, lighting variations, and distractors both individually and largely in composition. In the most challenging triple-factor variation setting, RoboAug demonstrates superior performance, achieving average success rates of 0.67, 0.47, and 0.60 across the three robots, significantly outperforming the leading baseline (0.42, 0.31, and 0.34). These results confirm that combining precise region extraction with contrastive learning is essential for reliable robot learning. Our main contributions are summarized as follows:

- We propose RoboAug, a region-contrastive data augmentation framework that facilitates the scalable generation of diverse training data with minimal human supervision.
- We introduce a one-shot region matching strategy combined with a region-contrastive loss, significantly improving both the precision of visual extraction and the expressiveness of learned policy features.
- Through extensive real-world experiments exceeding 35k trials, RoboAug exhibits robust generalization against diverse visual perturbations, outperforming baselines by 59.5%, 51.6%, and 76.4% on Tien Kung 2.0, UR-5e, and AgileX, respectively.
- We will open-source the embodied object detection dataset and our multi-task real-world manipulation dataset to facilitate further research.

## II. RELATED WORK

### A. Generalization in Visuomotor Policy Learning

Generalizing visuomotor policies to unstructured environments remains a pivotal challenge in robotic manipulation. While early Imitation Learning (IL) methods struggled with narrow demonstrations [14, 77, 24, 13, 46, 50, 65, 75, 30, 23, 7, 53], recent Vision-Language-Action (VLA) models leverage large-scale data to unlock emergent generalization. Models such as RT-1 [5], RT-2 [80], RT-X [45], and Octo [54] demonstrate that training on diverse cross-embodiment datasets, including BridgeData, Open X-Embodiment, and RoboMIND [16, 57, 31, 45, 64, 6, 29, 66, 27], can significantly enhance robustness. Furthermore, a rapidly expanding family of advanced architectures, ranging from PaLM-E [15] and  $\pi_0$  [4] to recent innovations [8, 59, 38, 62, 3, 34, 19, 39, 74, 76, 69, 47, 63] like HybridVLA [36],  $\pi_{0.5}$  [28], X-VLA [79] and XR-1 [18], has further utilized more internet data to enhance capabilities like multi-task reasoning, spatial understanding, and instruction following. Despite these advancements, acquiring high-quality robotic interaction data remains prohibitively expensive compared to internet-scale NLP or CV resources, leaving existing datasets insufficient to cover the heterogeneous distribution of real-world visual variations. To bridge this gap without the high cost of massive real-world collection, we introduce RoboAug, a data augmentation framework designed to operate at the data level. RoboAug is agnostic to network architecture and training paradigm, enabling seamless integration with diverse visuomotor policies and VLA models to enhance generalization against environmental variants.

### B. Data Augmentation for Robotic Manipulation

Data augmentation serves as a pivotal strategy in robotic learning to circumvent the prohibitive costs of large-scale real-world data collection. While weak augmentation techniques like random cropping and noise injection provide robustness against low-level pixel perturbations, they fail to introduce the semantic diversity required for out-of-distribution generalization. Consequently, the field has witnessed a paradigm shift towards strong generative augmentation [40, 11, 70], which leverages large-scale diffusion models to synthesize high-fidelity, semantically diverse training data. Seminal works like GenAug [11] pioneer this direction by utilizing pre-trained text-to-image models to retarget robot behaviors to unseen situations. By inpainting diverse backgrounds and textures while preserving the robot's pose, GenAug significantly expands the semantic support of the training distribution. Subsequent works [58, 55, 56, 12] have extended this paradigm. ROSIE [72] applies aggressive inpainting to generate distractors, while RoboAgent [2] combines semantic augmentation with action chunking. Similarly, methods like Mirage [9] and RoVi-Aug [10] utilize generative synthesis to bridge domain gaps across distinct robot embodiments and camera viewpoints.

A critical challenge in these generative pipelines is the precise preservation of task-relevant entities like the manipulated objects. Inaccurate masking during generation leads

to semantic corruption, where essential geometric cues are distorted or hallucinated away. To avoid physical constraints like green screens [55], recent works [20, 73, 35] focus on automation segmentation. Notably, RoboEngine [73] combines specialized segmentation with background generation to create physics-aware scenes. EAGLE [78] employs self-supervised control-aware masks. However, the majority of existing methods rely heavily on off-the-shelf generalist vision models (e.g., SAM2 [49]), which often struggle in complex manipulation scenarios involving severe occlusion or intricate object interactions. To address these limitations, RoboAug introduces a task-relevant region extraction phase, which leverages a one-shot region matching strategy to guarantee the structural integrity of visual cues. Furthermore, RoboAug employs a novel region-contrastive loss to enforce representation invariance on manipulated objects while encouraging robustness to background variations, ensuring the policy learns strictly from accurate, task-relevant visual features.

## III. METHODOLOGY

### A. Overview

In this work, we address the challenge of generalization in real-world robotic manipulation through the lens of single-task imitation learning. Formally, given a task specified by a language instruction  $l$ , we collect an expert dataset  $\mathcal{D}_e = \{l, \tau_i\}_{i=1}^N$ , where each trajectory  $\tau = \{(o_t, m_t, a_t)\}_{t=1}^T$  comprises a sequence of camera images  $o_t$ , proprioceptive states  $m_t$ , and robot actions  $a_t$ . To formulate the generalization problem, we decompose the visual observation  $o$  into task-relevant regions  $R_{\text{task}}$  (e.g., the robotic arm and manipulated objects) and task-irrelevant scenario factors  $R_{\text{scen}}$  (e.g., background, distractors, and lighting). Our objective is to learn a visuomotor policy  $\hat{a}_t = \pi(o_t, m_t)$  capable of maximizing success rates in novel environments characterized by unseen scenarios  $R_{\text{scen}}^{\text{new}}$ , all while the task-relevant elements  $R_{\text{task}}$  remain invariant. To achieve this, we introduce **RoboAug**, a region-contrastive data augmentation framework designed to enhance semantic diversity and feature robustness while minimizing annotation costs. As illustrated in Figure 3, our method proceeds in three stages: (1) task-relevant region extraction, (2) semantic data augmentation, and (3) region-contrastive policy learning, which collectively enable the precise identification of task-critical visual features for robust generalization. We provide implementation details in Appendix VI-A and theoretical analysis in Appendix VI-G.

### B. Task-Relevant Region Extraction

The goal in this stage is to obtain pixel-level masks  $M_{\text{task}} \in \{0, 1\}^{H \times W}$  of the task-relevant regions  $R_{\text{task}}$  from demonstration trajectories  $\tau$  without extensive manual annotation and costly detector retraining. We propose a lightweight, two-step extraction pipeline requiring only a single manually labeled reference image per task. First, we employ a training-free, one-shot region matching mechanism to locate key elements in the anchor frame of every trajectory. Second, we propagate these spatial annotations across subsequent frames using

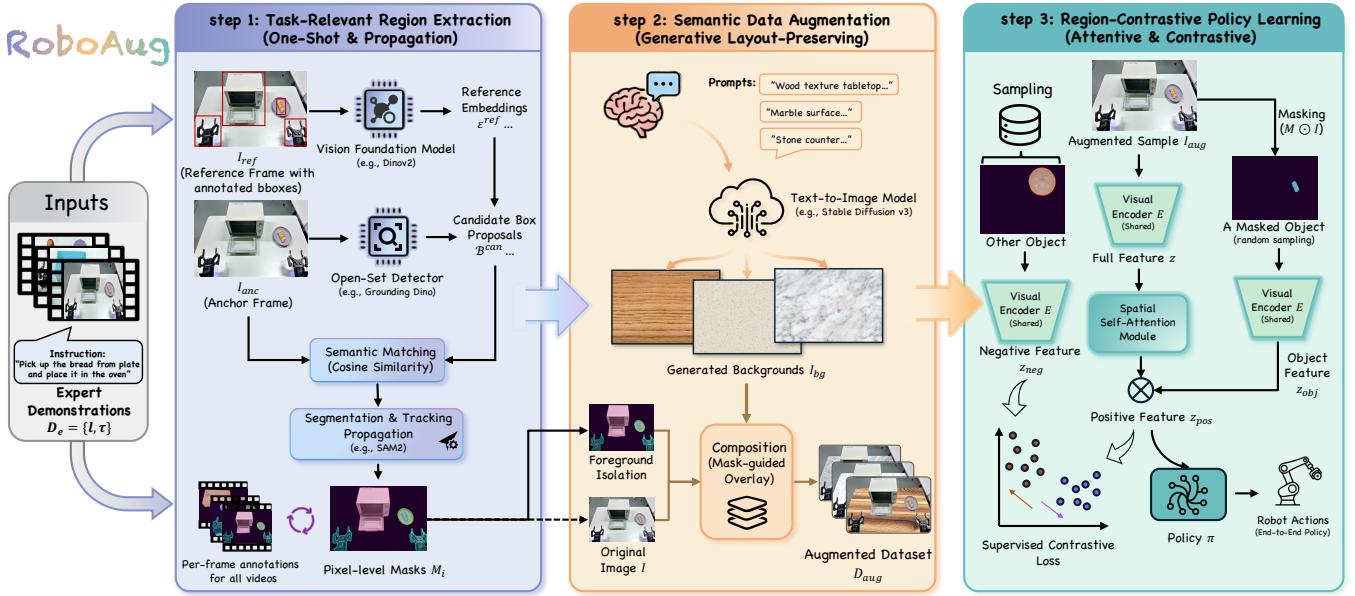


Fig. 3: Overview of RoboAug. RoboAug contains three stages: (1) task-relevant region extraction, (2) semantic data augmentation, and (3) region-contrastive policy learning.

semantic segmentation and tracking, yielding consistent pixel-level masks throughout the dataset.

**One-Shot Region Matching.** We designate the first frame of one trajectory as the *reference frame*, denoted as  $I_{ref}$ , as the first frame typically depicts task-relevant elements clearly without occlusion. We obtain the bounding box annotations  $\mathcal{B}^{ref} = \{B_i^{ref}\}_{i=1}^K$  for  $K$  task-relevant regions (e.g., manipulated objects) via a one-time manual labeling process. These regions are cropped and encoded into a set of reference embeddings  $\mathcal{E}^{ref} = \{e_i^{ref}\}_{i=1}^K$  using a vision foundation model (DINOv2 [44] in our implementation), where  $e_i^{ref} \in \mathbb{R}^d$ .

To transfer these labels to the remaining trajectories, we treat the first frame of each subsequent trajectory  $\tau$  as the *anchor frame*, denoted as  $I_{anc}$ . For each  $I_{anc}$ , we utilize an open-set detector (GroundingDINO [37] in our implementation) to generate candidate bounding box proposals  $\mathcal{B}^{can} = \{B_j^{can}\}_{j=1}^M$ . For each candidate  $B_j^{can}$ , we extract its feature embedding  $e_j^{can}$  and measure its semantic alignment with the reference templates  $\{e_i^{ref}\}_{i=1}^K$  via cosine similarity. The predicted category  $\hat{c}_j$  is determined by the most similar reference embedding:

$$\hat{c}_j = \operatorname{argmax}_{i \in \{1, \dots, K\}} \operatorname{sim}(e_j^{can}, e_i^{ref}). \quad (1)$$

This mechanism ensures robust, training-free alignment of task-relevant regions across diverse demonstrations while filtering out irrelevant background clutter.

**Semantic Mask Propagation.** Upon identifying task-relevant bounding boxes and their corresponding categories within each anchor frame  $I_{anc}$ , our goal is to extend this semantic information to the full trajectories  $\{\tau_i\}_{i=1}^N$ . To achieve this, we leverage a tracking-and-segmentation framework (SAM2 [49] in our implementation), which integrates semantic segmentation with temporal object tracking. This mechanism

allows us to transform sparse bounding box priors  $\mathcal{B}^{can}$  into dense pixel-level masks  $M_{task}$ , and propagate them with spatiotemporal consistency across all frames of each trajectory.

Consequently, every frame in the dataset is equipped with semantic masks  $M_{task}$  corresponding to the task-relevant elements, remarkably requiring manual bounding box annotations for only a single reference frame per task. These semantic masks serve a dual purpose: they not only precisely delineate task-relevant regions for downstream semantic data augmentation via generative models, but also provide fine-grained supervision signals that enhance feature representation in region-contrastive policy learning.

### C. Semantic Data Augmentation

Following the acquisition of precise pixel-level semantic masks  $M_{task}$  for task-relevant regions, RoboAug employs a semantic data augmentation strategy. This process is designed to significantly diversify the training dataset with environmental variations while rigorously preserving the structural integrity of task-critical elements. To automate the creation of diverse environmental contexts, we leverage a Large Language Model (ChatGPT [42] in our implementation) to generate a rich set of descriptive prompts. These prompts guide the image generation model in synthesizing distinct background textures. RoboAug constructs a library of 500 background description templates, systematically categorized into material types, including wood (58%), stone (35%), and composite materials (7%), to ensure a comprehensive coverage of real-world tabletop scenarios.

A key distinction of our approach, compared to prior semantic augmentation methods that rely on inpainting, is the handling of occlusions. We observe that inpainting techniques often introduce visual artifacts and geometric distortions, par-

TABLE I: RoboAug-D Dataset Statistics.

Robot	Task	Traj.	Frame	Obj.	BBox.
Single-Arm Franka	8	2442	20511	19	87426
Single-Arm UR	17	4217	40136	34	197882
Dual-Arm UR	3	669	11077	11	71464
Dual-Arm Agilex	5	248	2025	9	10063
Total	33	7576	73749	46	366835

ticularly when the robotic arm or objects occlude significant portions of the tabletop. To address this, we opt to generate a complete, coherent background image rather than filling in missing regions. This strategy ensures the photorealism and spatial continuity of the background.

Formally, let  $I$  denote the original image and  $M_{\text{task}}$  represent the semantic masks indicating the task-relevant regions (e.g., the robot arm and manipulated objects). We utilize the Stable Diffusion v3 model [17] to synthesize a full-frame background image  $I_{\text{bg}}$ .

Subsequently, we superimpose the preserved foreground regions from the original image onto the generated background using linear interpolation based on the mask  $M_{\text{task}}$ :

$$I_{\text{aug}} = M_{\text{task}} \odot I + (1 - M_{\text{task}}) \odot I_{\text{bg}}, \quad (2)$$

where  $\odot$  denotes element-wise multiplication. By iterating this process with randomly sampled prompts for each image, we expand the dataset by orders of magnitude. This results in a final augmented dataset  $\mathcal{D}_{\text{fnl}} = \mathcal{D}_e \cup \mathcal{D}_{\text{aug}}$  enriched with thousands of unique backgrounds, effectively simulating diverse real-world environments while maintaining high fidelity in critical task-relevant regions.

#### D. Region-Contrastive Policy Learning

While prior data augmentation techniques effectively enhance dataset diversity, they often overlook the critical role of task-relevant regional semantics during policy training. To bridge this gap, we introduce a novel Region-Contrastive Learning (RCL) objective that leverages task-relevant regions to refine the policy’s visual representations.

During the training phase, for each image  $I$  sampled from the final dataset  $\mathcal{D}_{\text{fnl}}$  within a batch of size  $B$ , we generate the corresponding masked images  $I_{\text{obj}}$  that isolates task-relevant objects. Formally, given an original image  $I \in \mathcal{D}_{\text{fnl}}$ , a binary mask  $M_{\text{task}}$  delineating the task-relevant object region, and its corresponding category  $c$ , we extract the object-centric image via an element-wise product:  $I_{\text{obj}} = M_{\text{task},c} \odot I$ . These inputs are subsequently processed by a visual encoder  $E(\cdot)$  to extract object feature embeddings  $z_{\text{obj}} = E(I_{\text{obj}})$ .

However, the masked images  $I_{\text{obj}}$  are often dominated by zero-valued (black) regions. This creates feature representations populated by non-informative signals, which can dilute task-critical semantics. To mitigate this, we leverage features from the full image  $z = E(I)$  to accentuate the salient information in  $z_{\text{obj}}$ , because  $z$  share the same visual encoder

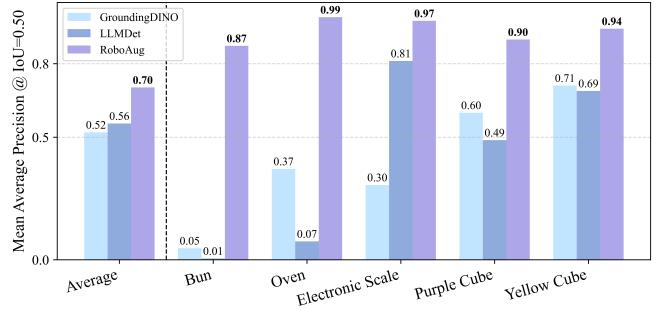


Fig. 4: Comparison of mAP@0.5 across RoboAug-D Dataset. We present the results of 5 representative objects.

and contain all the information. We apply a spatial self-attention mechanism to yield the final attentive features:

$$a_{\text{att}} = \text{sigmoid}(A(z) \odot z), \quad z_{\text{att}} = a_{\text{att}} \odot z_{\text{obj}}. \quad (3)$$

where  $A(\cdot)$  denotes the learnable self-attention module, and the  $\text{sigmoid}(\cdot)$  operation normalizes the attention scores to generate the spatial weight map  $a_{\text{att}}$ .

To align representations within the same category while separating distinct objects, we optimize the visual encoder using a supervised contrastive loss. We construct positive pairs from samples sharing the same object category  $c$ , and negative pairs from images of differing categories. Inspired by [32], we formulated the region-contrastive loss as:

$$\mathcal{L}_{\text{RC}} = \sum_{i \in \mathcal{B}_{\text{obj}}} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(z_{\text{att},i} \cdot z_{\text{att},p}/d)}{\sum_{j \in S(i)} \exp(z_{\text{att},i} \cdot z_{\text{att},j}/d)} \quad (4)$$

where  $i$  is the index of the selected sample within the augmented batch  $\mathcal{B}_{\text{obj}}$ , and  $P(i) = \{p \in \mathcal{B}_{\text{obj}} : c_p = c_i\}$  represents the set of indices for positive samples (i.e., those sharing the class label with sample  $i$ ). The set  $S(i) = \mathcal{B}_{\text{obj}} \setminus \{i\}$  encompasses all indices in the batch excluding the selected itself. The  $d$  denotes the temperature parameter, which regulates the concentration of the feature distribution.

#### E. RoboAug-D Dataset for Object Detection

State-of-the-art vision foundation models, such as GroundingDINO [37] and LLMDet [22], often exhibit performance degradation when applied to robotic manipulation scenes.

To investigate and benchmark model robustness in these domains, we introduce **RoboAug-D**, a large-scale object detection dataset manually annotated from the perspective of robotic manipulators. As shown in Table I, the dataset encompasses 33 distinct tasks, comprising a total of 73,749 keyframes and 366,835 bounding boxes across 46 object categories.

## IV. EXPERIMENTS

In this section, we empirically validate RoboAug from fundamental visual capabilities to complex physical manipulation. We begin by benchmarking the limits of state-of-the-art vision foundation models in Section IV-A. Transitioning to the real world, Sections IV-B–IV-D present comprehensive manipulation experiments, evaluating both broad compositional

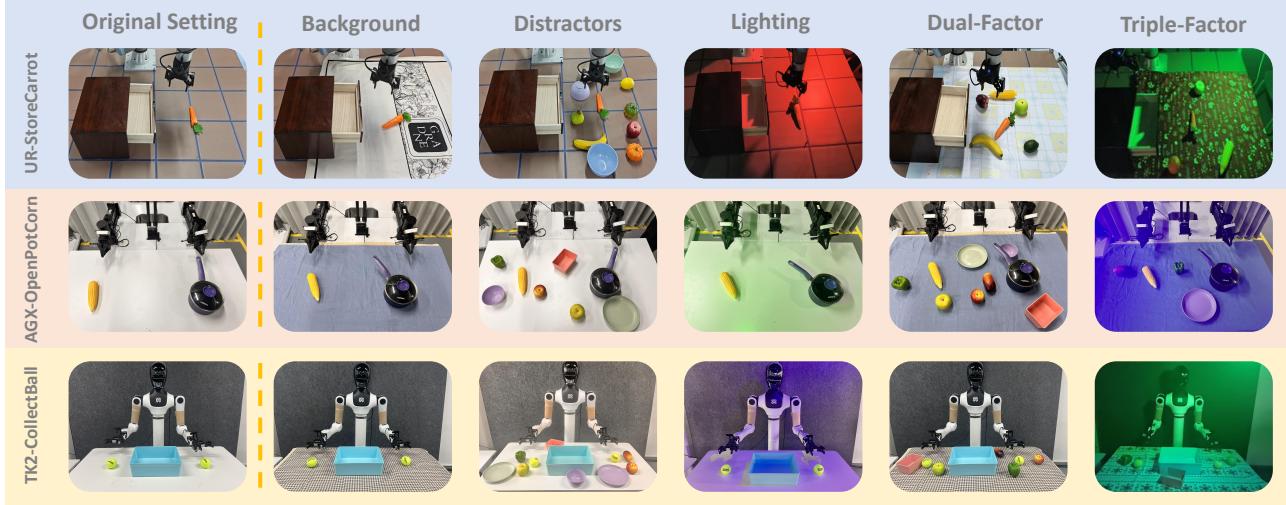


Fig. 5: Overview of the generalization evaluation settings, spanning single-factor variations and compositional dual- and triple-factor scenes involving background, distractors, and lighting.

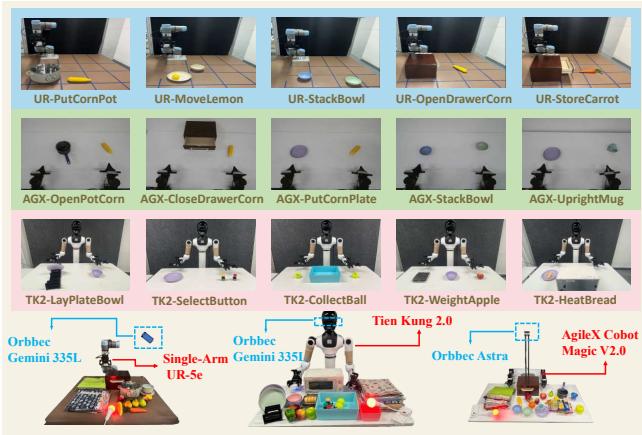


Fig. 6: Experimental Setup. We evaluate RoboAug across three robot embodiments.

generalization and single-factor robustness. In addition, we conduct an ablation study on region-contrastive learning in Section IV-E and analyze scaling laws regarding augmentation magnitude in Section IV-F. Finally, we report complementary simulation results in Section IV-G, with further multidimensional evaluations in Appendices VI-B - VI-F.

#### A. Object Detection on RoboAug-D Dataset

**Experimental Setup.** We evaluated the zero-shot object detection capabilities of Vision Foundation Models (VFs) on the full test set of the RoboAug-D dataset, with a focus on challenges inherent to robotic manipulation scenarios. To ensure a fair assessment of intrinsic generalization, all models were evaluated without fine-tuning. For each of the 46 object categories, we queried the models using five diverse text prompts (e.g., varying in phrasing, synonyms, and functional descriptions). Model performance was compared using mean average precision (mAP@0.5), and our proposed approach was benchmarked against GroundingDINO and LLMDet.

**Experimental Results.** Figure 4 provides the quantitative results, detailing both the overall mAP@0.5 on representative object categories. Our approach demonstrates significant improvements, outperforming the state-of-the-art baselines GroundingDINO and LLMDet by 34.6% and 25.0%, respectively.

For instance, in the “Bun” category, baselines struggle to exceed 0.10 mAP, whereas our method achieves scores of 0.87. These results highlight the limitations of current VFs in handling robotic viewpoints and validate the effectiveness of our one-shot region matching strategy in enhancing detection accuracy for downstream data augmentation.

#### B. Real-World Generalizable Robotic Manipulation

**Hardware Setup.** We validate RoboAug across three diverse robots illustrated in Figure 6: (1) the single-arm UR-5e, (2) the Tien Kung 2.0 humanoid robot, and (3) the AgileX Cobot Magic V2.0 robot.

We collect the dataset via human teleoperation HACTS [67], recording visual observations, robot states, and actions at every frame.

**Task Design.** As shown in Figure 6, we designed five tasks per embodiment to cover a range of complexities, extending from single-arm pick-and-place to precise dual-arm collaboration. These tasks require diverse skills, including pushing, rotating, and grasping. The UR-5e performs household interactions such as PutCornPot and OpenDrawerCorn. The AgileX and Tien Kung 2.0 robots execute complex bimanual tasks, including UprightMug and LayPlateBowl. For each task, we collected a dataset comprising 50 expert trajectories.

**Generalization Evaluation and Metrics.** We devised two protocols to assess robustness: *Compositional Generalization*, which evaluates adaptability across combined environmental variables, and *Single-Factor Generalization*, which probes stability against intense variations in specific factors. Variables include background textures, lighting conditions, and task-irrelevant distractors. We report the success rate averaged over

TABLE II: Comparative results under triple-factor variations: 3 unseen backgrounds, 4 lighting conditions, and 3 distractors.

Augmentation Method	UR-PutCornPot	UR-MoveLemon	UR-StackBowl	UR-StoreCarrot	UR-OpenDrawerCorn	Average
ACT w/o Aug [77]	0.06	0.06	0.08	0.10	0.15	0.09
RoboEngine-T [73]	0.12	0.16	0.12	0.20	0.32	0.18
RoboEngine-G [73]	0.16	0.24	0.12	0.32	0.40	0.25
GenAug [11]	0.22	0.28	0.16	0.40	0.48	0.31
<b>RoboAug</b>	<b>0.38</b>	<b>0.46</b>	<b>0.28</b>	<b>0.56</b>	<b>0.68</b>	<b>0.47</b>
	AGX-PutCornPlate	AGX-UprightMug	AGX-StackBowl	AGX-OpenPotCorn	AGX-CloseDrawerCorn	
ACT w/o Aug [77]	0.12	0.14	0.14	0.16	0.24	0.16
RoboEngine-T [73]	0.24	0.22	0.28	0.30	0.28	0.26
RoboEngine-G [73]	0.30	0.32	0.32	0.28	0.38	0.32
GenAug [11]	0.32	0.34	0.30	0.32	0.40	0.34
<b>RoboAug</b>	<b>0.51</b>	<b>0.58</b>	<b>0.62</b>	<b>0.55</b>	<b>0.73</b>	<b>0.60</b>
	TK2-WeightApple	TK2-CollectBall	TK2-HeatBread	TK2-SelectButton	TK2-LayPlateBowl	
ACT w/o Aug [77]	0.20	0.16	0.14	0.22	0.24	0.19
RoboEngine-T [73]	0.28	0.25	0.32	0.34	0.30	0.30
RoboEngine-G [73]	0.38	0.36	0.40	0.35	0.42	0.38
GenAug [11]	0.48	0.32	0.42	0.40	0.50	0.42
<b>RoboAug</b>	<b>0.80</b>	<b>0.52</b>	<b>0.65</b>	<b>0.68</b>	<b>0.70</b>	<b>0.67</b>

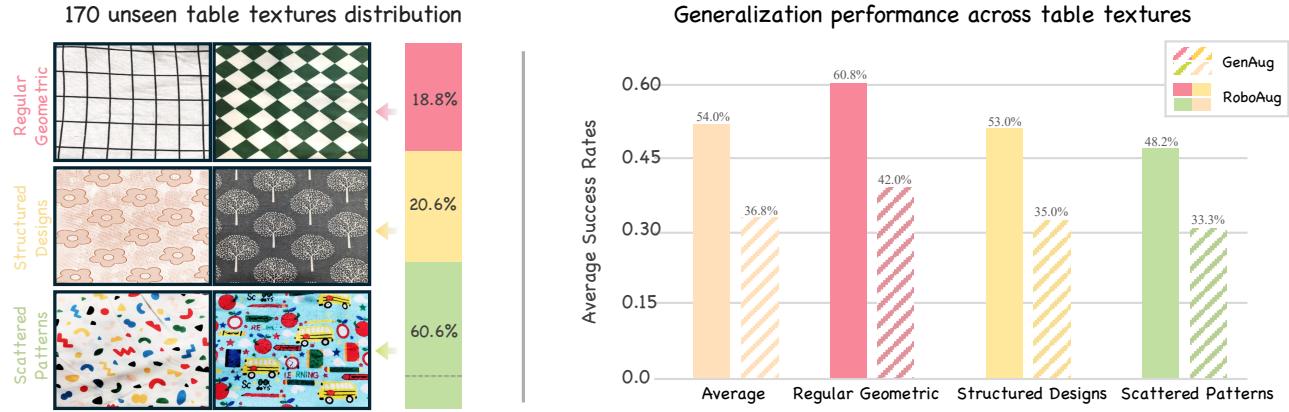


Fig. 7: Background generalization on task UR-PutCornPot across 170 unseen backgrounds.

20 real-world rollouts per configuration.

### C. Compositional Generalization Evaluation

**Evaluation Setup.** We evaluate our policy under a challenging *triple-factor* protocol incorporating 3 unseen backgrounds, 3 task-irrelevant distractors, and 4 distinct lighting conditions. We compared RoboAug against a non-augmented policy (ACT [77]), a texture-replacement method (RoboEngine-T [73]), and two generative baselines (RoboEngine-G [73] and GenAug [11]). All augmentation methods employ a  $5 \times$  data expansion ratio, supplementing 50 real-world expert trajectories with 250 generated trajectories. Additional results of the *Dual-Factor* setting are detailed in Appendix VI-B.

**Results under Triple-Factor Variation.** Table II shows the performance on all the robots. RoboAug consistently outperforms all baselines. Notably, in the AGX-UprightMug task, which requires rotating a mug and coordinating placement, RoboAug achieves a success rate of 0.58, significantly surpassing the strongest baseline GenAug (0.34). Results for UR-5e and Tien Kung 2.0 show similar gains, demonstrating the embodiment-agnostic generalization of RoboAug.

TABLE III: Ablation study on region-contrastive loss.

Method	RCL	UR-Put CornPot	AGX-Put CornPlate	TK2-Weight Apple	Average
ACT w/o Aug [77]	✗	0.06	0.12	0.20	0.13
ACT w/o Aug [77]	✓	0.12	0.20	0.32	0.21
RoboEngine-T [73]	✗	0.12	0.24	0.28	0.21
RoboEngine-T [73]	✓	0.14	0.28	0.30	0.24
RoboEngine-G [73]	✗	0.16	0.30	0.38	0.28
RoboEngine-G [73]	✓	0.16	0.36	0.44	0.32
GenAug [11]	✗	0.22	0.32	0.48	0.34
GenAug [11]	✓	0.28	0.35	0.53	0.39
RoboAug	✗	0.28	0.43	0.68	0.46
RoboAug	✓	0.38	0.51	0.80	0.56

### D. Single-Factor Generalization Evaluation

**Evaluation Setup.** To rigorously assess generalization boundaries, we isolate three environmental factors: (1) *Background Diversity*, where we introduce 170 unseen textures across three complexity levels (Geometric, Structured, Scattered); (2) *Distractor Density*, where we increase workspace clutter to 10 objects; and (3) *Lighting*, which spans 20 distinct conditions including dynamic shifts. We provide results

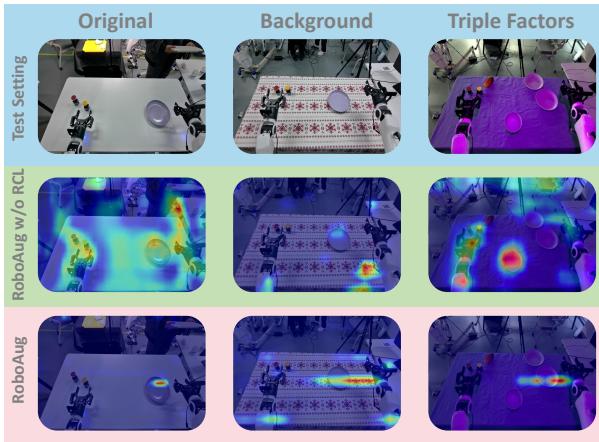


Fig. 8: Feature heatmap comparison of RoboAug with and without region-contrastive loss (RCL).

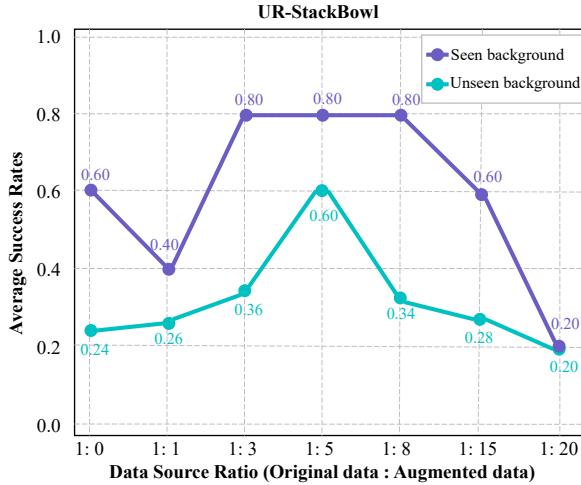


Fig. 9: Ablation study on the effect of data augmentation ratio. We evaluate ratios from 1:0 (only raw data) to 1:20.

regarding distractor and lighting variations in Sections VI-C and VI-D of the Appendix.

**Results on Unseen Background Variation.** Figure 7 illustrates performance on UR-PutCornPot across the 170 unseen backgrounds. RoboAug achieves a mean success rate of 54.0%, significantly outperforming GenAug (36.8%). While success rates for both methods naturally decrease as background patterns become more intricate, GenAug exhibits a sharper decline. These results confirm that RoboAug effectively maintains policy focus on foreground objects despite severe background visual distractions.

#### E. Effectiveness of Region-Contrastive Learning

**Quantitative Results.** Table III presents the impact of the region-contrastive loss on policy performance. We observe a consistent improvement across all baselines when integrating RCL. Even for the standard baseline ACT w/o Aug, RCL doubles the success rate in the UR-PutCornPot task.



Fig. 10: Experimental setup for evaluating generalization on the LIBERO-Plus benchmark.

TABLE IV: Generalization performance on the LIBERO-Plus benchmark.

Method	Background	Distractor	Light	Average
ACT w/o Aug [77]	0.745	0.860	0.789	0.798
RoboEngine-G [73]	0.806	0.942	0.855	0.868
<b>RoboAug</b>	<b>0.913</b>	<b>0.990</b>	<b>0.896</b>	<b>0.933</b>

The addition of RCL boosts RoboAug’s performance on the challenging TK2-WeightApple task from 0.68 to 0.80. This trend indicates that RCL effectively enhances feature robustness and generalization capability, regardless of the underlying data augmentation strategy.

**Visualization Analysis.** To interpret the learned representations, we visualize feature attention heatmaps using Grad-CAM [52] in Figure 8. We compare activations across three increasingly difficult scenarios: original environments, unseen backgrounds, and complex settings with triple factors. The baseline (w/o RCL) exhibits diffuse attention, easily distracted by high-frequency background textures or task-irrelevant objects, particularly under severe lighting changes. In contrast, RoboAug with RCL maintains precise localization on specific grasping points, effectively filtering out environmental noise and distractors. This visual evidence confirms that the region-contrastive objective forces the policy to encode task-relevant semantics invariant to visual perturbations, corroborating the quantitative improvements in Table III.

#### F. Scaling Law Analysis of Data Augmentation

We analyze the scaling laws of data augmentation on the UR-StackBowl task by varying the ratio. As shown in Figure 9, performance follows an inverted U-shaped trend.

Moderate settings (1:3 to 1:8) act as effective regularization and significantly improve success rates. However, excessive augmentation ( $> 1:15$ ) saturates the network’s finite capacity, causing rapid deterioration. These findings confirm that a balanced ratio ( $\approx 1:5$ ), rather than simply maximizing data quantity, is critical for optimal performance.

#### G. Results on Simulation

As shown in Figure 10, we utilize the LIBERO-Plus benchmark [21], which introduces diverse environmental shifts to the standard tasks. Policies were trained on the LIBERO-Object dataset (10 tasks, 50 demonstrations each) and tested under three perturbation types: background, distractors, and lighting. As summarized in Table IV, RoboAug consistently outperforms the baselines across all categories. Our method

achieves an average success rate of 93.3%, surpassing the strongest baseline, RoboEngine-G, by a significant margin of 6.5%. These results highlight the effectiveness of RoboAug in preventing policy degradation, particularly in scenarios with complex visual distractors and background changes.

## V. CONCLUSION

We introduced RoboAug, a data augmentation framework designed to enhance robotic generalization across diverse and unseen scenes. Unlike prior methods that rely on large-scale pre-training or assume perfect object recognition, our approach requires only a single framework annotation. By utilizing generative models for semantic augmentation and integrating a plug-and-play region-contrastive loss, RoboAug effectively guides the model to focus on task-relevant regions. Extensive real-world validation, comprising over 35k trials on UR-5e, AgileX, and Tien Kung 2.0 robots, demonstrates that RoboAug consistently outperforms state-of-the-art baselines. These results highlight the superior effectiveness and robustness of RoboAug in complex real-world manipulation tasks.

## REFERENCES

- [1] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of machine learning research*, 3(Nov):463–482, 2002.
- [2] Homanga Bharadhwaj, Jay Vakil, Mohit Sharma, Abhinav Gupta, Shubham Tulsiani, and Vikash Kumar. Roboagent: Generalization and efficiency in robot manipulation via semantic augmentations and action chunking. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4788–4795. IEEE, 2024.
- [3] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr0ot n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- [4] 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.
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- [6] Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xuan Hu, Xu Huang, et al. Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025.
- [7] Jiahang Cao, Qiang Zhang, Jingkai Sun, Jiaxu Wang, Hao Cheng, Yulin Li, Jun Ma, Kun Wu, Zhiyuan Xu, Yecheng Shao, et al. Mamba policy: Towards efficient 3d diffusion policy with hybrid selective state models. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2025.
- [8] Chi-Lam Cheang, Guangzeng Chen, Ya Jing, Tao Kong, Hang Li, Yifeng Li, Yuxiao Liu, Hongtao Wu, Jiafeng Xu, Yichu Yang, et al. Gr-2: A generative video-language-action model with web-scale knowledge for robot manipulation. *arXiv preprint arXiv:2410.06158*, 2025.
- [9] Lawrence Yunliang Chen, Kush Hari, Karthik Dharmajan, Chenfeng Xu, Quan Vuong, and Ken Goldberg. Mirage: Cross-embodiment zero-shot policy transfer with cross-painting. In *Proceedings of Robotics: Science and Systems*, 2024.
- [10] Lawrence Yunliang Chen, Chenfeng Xu, Karthik Dharmajan, Muhammad Zubair Irshad, Richard Cheng, Kurt Keutzer, Masayoshi Tomizuka, Quan Vuong, and Ken Goldberg. Rovi-aug: Robot and viewpoint augmentation for cross-embodiment robot learning. *arXiv preprint arXiv:2409.03403*, 2024.
- [11] Zoey Chen, Sho Kiami, Abhishek Gupta, and Vikash Kumar. Genaug: Retargeting behaviors to unseen situations via generative augmentation. *arXiv preprint arXiv:2302.06671*, 2023.
- [12] Zoey Chen, Zhao Mandi, Homanga Bharadhwaj, Mohit Sharma, Shuran Song, Abhishek Gupta, and Vikash Kumar. Semantically controllable augmentations for generalizable robot learning. *The International Journal of Robotics Research*, 44(10-11):1705–1726, 2025.
- [13] 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*, page 02783649241273668, 2023.
- [14] Zichen Jeff Cui, Yibin Wang, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. From play to policy: Conditional behavior generation from uncurated robot data. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- [15] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. PaLM-e: An embodied multimodal language model. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 8469–8488, 2023.
- [16] Frederik Ebert, Ynlai Yang, Karl Schmeckpeper, Bernadette Bucher, Georgios Georgakis, Kostas Daniilidis, Chelsea Finn, and Sergey Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets. In *Proceedings of Robotics: Science and Systems (RSS)*, 2022.

- [17] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- [18] Shichao Fan, Kun Wu, Zhengping Che, Xinhua Wang, Di Wu, Fei Liao, Ning Liu, Yixue Zhang, Zhen Zhao, Zhiyuan Xu, et al. Xr-1: Towards versatile vision-language-action models via learning unified vision-motion representations. *arXiv preprint arXiv:2511.02776*, 2025.
- [19] Shichao Fan, Quantao Yang, Yajie Liu, Kun Wu, Zhengping Che, Qingjie Liu, and Min Wan. Diffusion trajectory-guided policy for long-horizon robot manipulation. *IEEE Robotics and Automation Letters*, 10(12):12788–12795, 2025. doi: 10.1109/LRA.2025.3619794.
- [20] Yu Fang, Yue Yang, Xinghao Zhu, Kaiyuan Zheng, Gedas Bertasius, Daniel Szafir, and Mingyu Ding. Rebot: Scaling robot learning with real-to-sim-to-real robotic video synthesis. *arXiv preprint arXiv:2503.14526*, 2025.
- [21] Senyu Fei, Siyin Wang, Junhao Shi, Zihao Dai, Jikun Cai, Pengfang Qian, Li Ji, Xinze He, Shiduo Zhang, Zhaoye Fei, Jinlan Fu, Jingjing Gong, and Xipeng Qiu. Libero-plus: In-depth robustness analysis of vision-language-action models, 2025. URL <https://arxiv.org/abs/2510.13626>.
- [22] Shenghao Fu, Qize Yang, Qijie Mo, Junkai Yan, Xihan Wei, Jingke Meng, Xiaohua Xie, and Wei-Shi Zheng. Llmdet: Learning strong open-vocabulary object detectors under the supervision of large language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 14987–14997, 2025.
- [23] Zipeng Fu, Tony Z Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. In *Conference on Robot Learning (CoRL)*, 2024.
- [24] Theophile Gervet, Zhou Xian, Nikolaos Gkanatsios, and Katerina Fragkiadaki. Act3d: 3d feature field transformers for multi-task robotic manipulation. In *7th Annual Conference on Robot Learning*, 2023.
- [25] Nicklas Hansen and Xiaolong Wang. Generalization in reinforcement learning by soft data augmentation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 13611–13617. IEEE, 2021.
- [26] Alex Hernández-García and Peter König. Data augmentation instead of explicit regularization. *arXiv preprint arXiv:1806.03852*, 2018.
- [27] Chengkai Hou, Kun Wu, Jiaming Liu, Zhengping Che, Di Wu, Fei Liao, Guangrun Li, Jingyang He, Qiuxuan Feng, Zhao Jin, et al. Robomind 2.0: A multimodal, bimanual mobile manipulation dataset for generalizable embodied intelligence. *arXiv preprint arXiv:2512.24653*, 2025.
- [28] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, et al.  $\pi_{0.5}$ : a vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.
- [29] Tao Jiang, Tianyuan Yuan, Yicheng Liu, Chenhao Lu, Jianning Cui, Xiao Liu, Shuiqi Cheng, Jiyang Gao, Huazhe Xu, and Hang Zhao. Galaxea open-world dataset and g0 dual-system vla model. *arXiv preprint arXiv:2509.00576*, 2025.
- [30] Tsung-Wei Ke, Nikolaos Gkanatsios, and Katerina Fragkiadaki. 3d diffuser actor: Policy diffusion with 3d scene representations. In *8th Annual Conference on Robot Learning*, 2024.
- [31] 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. *arXiv preprint arXiv:2403.12945*, 2024.
- [32] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
- [33] Kimin Lee, Kibok Lee, Jinwoo Shin, and Honglak Lee. Network randomization: A simple technique for generalization in deep reinforcement learning. *International Conference on Learning Representations*, 2019.
- [34] Meng Li, Zhen Zhao, Zhengping Che, Fei Liao, Kun Wu, Zhiyuan Xu, Pei Ren, Zhao Jin, Ning Liu, and Jian Tang. Switchvla: Execution-aware task switching for vision-language-action models. *arXiv preprint arXiv:2506.03574*, 2025.
- [35] I Liu, Chun Arthur, Jason Chen, Gaurav Sukhatme, and Daniel Seita. D-coda: Diffusion for coordinated dual-arm data augmentation. *arXiv preprint arXiv:2505.04860*, 2025.
- [36] Jiaming Liu, Hao Chen, Pengju An, Zhuoyang Liu, Renrui Zhang, Chenyang Gu, Xiaoqi Li, Ziyu Guo, Sixiang Chen, Mengzhen Liu, et al. Hybridvla: Collaborative diffusion and autoregression in a unified vision-language-action model. *arXiv preprint arXiv:2503.10631*, 2025.
- [37] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European conference on computer vision*, pages 38–55, 2024.
- [38] 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. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- [39] Zhuoyang Liu, Jiaming Liu, Jiadong Xu, Nuowei Han, Chenyang Gu, Hao Chen, Kaichen Zhou, Renrui Zhang, Kai Chin Hsieh, Kun Wu, et al. Mla: A multisensory language-action model for multimodal understanding and

- forecasting in robotic manipulation. *arXiv preprint arXiv:2509.26642*, 2025.
- [40] Zhao Mandi, Homanga Bharadhwaj, Vincent Moens, Shuran Song, Aravind Rajeswaran, and Vikash Kumar. Cacti: A framework for scalable multi-task multi-scene visual imitation learning. *arXiv preprint arXiv:2212.05711*, 2022.
- [41] 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 : Open x-embodiment collaboration0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903, 2024.
- [42] OpenAI. chatgpt. <https://chatgpt.com/>, 2025.
- [43] OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2025.
- [44] Maxime Oquab, Timothée Darcet, 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 Journal*, 2024.
- [45] Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024.
- [46] Tim Pearce, Tabish Rashid, Anssi Kanervisto, Dave Bignell, Mingfei Sun, Raluca Georgescu, Sergio Valcarcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, et al. Imitating human behaviour with diffusion models. *ICLR*, 2023.
- [47] Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Zhao, Dong Wang, et al. Spatialavia: Exploring spatial representations for visual-language-action model. In *Proceedings of Robotics: Science and Systems (RSS)*, 2025.
- [48] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831, 2021.
- [49] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [50] Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal conditioned imitation learning using score-based diffusion policies. In *Robotics: Science and Systems*, 2023.
- [51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [52] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [53] Yue Su, Xinyu Zhan, Hongjie Fang, Han Xue, Hao-Shu Fang, Yong-Lu Li, Cewu Lu, and Lixin Yang. Dense policy: Bidirectional autoregressive learning of actions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2025.
- [54] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, 2024.
- [55] Eugene Teoh, Sumit Patidar, Xiao Ma, and Stephen James. Green screen augmentation enables scene generalisation in robotic manipulation. *arXiv preprint arXiv:2407.07868*, 2024.
- [56] Stephen Tian, Blake Wulfe, Kyle Sargent, Katherine Liu, Sergey Zakharov, Vitor Guizilini, and Jiajun Wu. View-invariant policy learning via zero-shot novel view synthesis. *arXiv preprint arXiv:2409.03685*, 2024.
- [57] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning*, pages 1723–1736. PMLR, 2023.
- [58] Jun Wang, Yuzhe Qin, Kaiming Kuang, Yigit Korkmaz, Akhilan Gurumoorthy, Hao Su, and Xiaolong Wang. Cyberdemo: Augmenting simulated human demonstration for real-world dexterous manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17952–17963, 2024.
- [59] Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaiming He. Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers. *Advances in neural information processing systems*, 37:124420–124450, 2024.
- [60] Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricu, et al. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18359–18369, 2023.
- [61] Wei Wang. Advanced auto labeling solution with added features. <https://github.com/CVHub520/X-AnyLabeling>,

2023.

- [62] Junjie Wen, Yichen Zhu, Jinming Li, Minjie Zhu, Zhibin Tang, Kun Wu, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin Shen, et al. Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation. *IEEE Robotics and Automation Letters*, 2025.
- [63] Junjie Wen, Yichen Zhu, Minjie Zhu, Zhibin Tang, Jinming Li, Zhongyi Zhou, Xiaoyu Liu, Chaomin Shen, Yixin Peng, and Feifei Feng. Diffusionvla: Scaling robot foundation models via unified diffusion and autoregression. In *Forty-second International Conference on Machine Learning*, 2025.
- [64] Kun Wu, Chengkai Hou, Jiaming Liu, Zhengping Che, Xiaozhu Ju, Zhiqin Yang, Meng Li, Yinuo Zhao, Zhiyuan Xu, Guang Yang, et al. Robomind: Benchmark on multi-embodiment intelligence normative data for robot manipulation. In *Proceedings of Robotics: Science and Systems (RSS)*, 2025.
- [65] Kun Wu, Yichen Zhu, Jinming Li, Junjie Wen, Ning Liu, Zhiyuan Xu, and Jian Tang. Discrete policy: Learning disentangled action space for multi-task robotic manipulation. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025.
- [66] Shihan Wu, Xuecheng Liu, Shaoxuan Xie, Pengwei Wang, Xinghang Li, Bowen Yang, Zhe Li, Kai Zhu, Hongyu Wu, Yiheng Liu, et al. Robocoin: An open-sourced bimanual robotic data collection for integrated manipulation. *arXiv preprint arXiv:2511.17441*, 2025.
- [67] Zhiyuan Xu, Yinuo Zhao, Kun Wu, Ning Liu, Junjie Ji, Zhengping Che, Chi Harold Liu, and Jian Tang. Hacts: a human-as-copilot teleoperation system for robot learning. *arXiv preprint arXiv:2503.24070*, 2025.
- [68] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chen-gen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- [69] Shuai Yang, Hao Li, Yilun Chen, Bin Wang, Yang Tian, Tai Wang, Hanqing Wang, Feng Zhao, Yiyi Liao, and Jiangmiao Pang. Instructvla: Vision-language-action instruction tuning from understanding to manipulation. *arXiv preprint arXiv:2507.17520*, 2025.
- [70] Sizhe Yang, Wenyu Yu, Jia Zeng, Jun Lv, Kerui Ren, Cewu Lu, Dahua Lin, and Jiangmiao Pang. Novel demonstration generation with gaussian splatting enables robust one-shot manipulation. *arXiv preprint arXiv:2504.13175*, 2025.
- [71] Raymond A Yeh, Chen Chen, Teck Yian Lim, Alexander G Schwing, Mark Hasegawa-Johnson, and Minh N Do. Semantic image inpainting with deep generative models. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5485–5493, 2017.
- [72] Tianhe Yu, Ted Xiao, Austin Stone, Jonathan Tompson, Anthony Brohan, Su Wang, Jaspiar Singh, Clayton Tan, Jodilyn Peralta, Brian Ichter, et al. Scaling robot learning with semantically imagined experience. *arXiv preprint arXiv:2302.11550*, 2023.
- [73] Chengbo Yuan, Suraj Joshi, Shaoting Zhu, Hang Su, Hang Zhao, and Yang Gao. Roboengine: Plug-and-play robot data augmentation with semantic robot segmentation and background generation. *arXiv preprint arXiv:2503.18738*, 2025.
- [74] Yifu Yuan, Haiqin Cui, Yibin Chen, Zibin Dong, Fei Ni, Longxin Kou, Jinyi Liu, Pengyi Li, Yan Zheng, and Jianye Hao. From seeing to doing: Bridging reasoning and decision for robotic manipulation. *arXiv preprint arXiv:2505.08548*, 2025.
- [75] Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, Muhan Wang, and Huazhe Xu. 3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- [76] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han, Chelsea Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1702–1713, 2025.
- [77] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- [78] Yinuo Zhao, Kun Wu, Tianjiao Yi, Zhiyuan Xu, Zhengping Che, Chi Harold Liu, and Jian Tang. Efficient training of generalizable visuomotor policies via control-aware augmentation. In *Proceedings of the 24th International Conference on Autonomous Agents and Multiagent Systems*, pages 2832–2834, 2025.
- [79] Jinliang Zheng, Jianxiong Li, Zhihao Wang, Dongxiu Liu, Xirui Kang, Yuchun Feng, Yinan Zheng, Jiayin Zou, Yilun Chen, Jia Zeng, et al. X-vla: Soft-prompted transformer as scalable cross-embodiment vision-language-action model. *arXiv preprint arXiv:2510.10274*, 2025.
- [80] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023.

## VI. APPENDIX

### A. Implementation Details

In this section, we provide a comprehensive description of the RoboAug framework, detailing the pipeline from task-relevant region extraction to region-contrastive policy learning.

**Object Detection Details.** To obtain accurate bounding boxes for task-relevant elements  $\mathcal{B}^{\text{ref}}$ , we manually annotated the initial dataset using the X-AnyLabeling tool [61]. These regions were cropped and resized to  $224 \times 224$  pixels to align with the input specifications of DINOv2 (86M parameters) [44] and encoded into a set of reference embeddings  $\mathcal{E}^{\text{ref}}$ . During inference, we employ the open-set detector GroundingDINO [37] to generate candidate bounding boxes  $B_j^{\text{can}}$  with both box and text thresholds set to 0.15. For each candidate box, Grounding DINO outputs a confidence score  $\delta$ ; if  $\delta > 0.7$ , the candidate  $B_j^{\text{can}}$  is assigned to the corresponding category directly. Otherwise, the predicted category  $\hat{c}_j$  is determined by selecting the category with the highest cosine similarity to the reference embeddings.

**Semantic Data Augmentation Details.** We use Stable Diffusion 3 Medium [17] for text-to-image generation. The inference steps (`num_inference_steps`) and guidance scale (`guidance_scale`) are set to 30 and 10.0, respectively. Image augmentation is performed with a batch size of 12, and generating a trajectory of 200 frames takes approximately 0.2 GPU-hours on an NVIDIA A100 GPU.

**Region-Contrastive Policy Learning Details.** We apply the proposed Region-Contrastive Loss to two policy architectures: ACT [77] and  $\pi_0$  [4]. While ACT relies solely on third-view RGB images and robot states as input,  $\pi_0$  additionally incorporates language instructions. Further training hyperparameter details are summarized in Table V.

TABLE V: Implementation Details.

	Hyperparameter	Value		Hyperparameter	Value
ACT	Batch Size	24	$\pi_0$	Batch Size	256
	Learning Rate	1e-4		Learning Rate	5e-5
	Optimizer	AdamW		Optimizer	AdamW
	Vision Encoder	ResNet50		Vision Encoder	SigLip
	Action Loss	L2 + RCL		Action Loss	Flow Matching + RCL
	Training Step temperature	50K 0.07		Training Step temperature	30K 0.07

**Real-world Task Setup.** The real-world evaluation is conducted on three robot embodiments: UR-5e (UR), AgileX (AGX), and TienKung2 (TK2). The evaluated tasks are detailed below.

- **UR-PutCornPot:** Transporting a piece of corn into a cooking pot. The corn is randomly placed within a rectangular region of  $20 \text{ cm} \times 60 \text{ cm}$ .
- **UR-MoveLemon:** Relocating a lemon from a plate to a bowl. The bowl is placed stochastically within the region of  $20 \text{ cm} \times 20 \text{ cm}$  grid.
- **UR-StackBowl:** Stacking bowls in a controlled manner. The position of one bowl is fixed, whereas the second bowl is uniformly sampled along a straight line of length 60 cm to introduce spatial variation.

- **UR-OpenDrawerCorn:** Opening a drawer and placing corn inside. A random orientation is assigned to the corn, with the rotation angle sampled uniformly from the interval  $[-\pi/4, \pi/4]$  radians.
- **UR-StoreCarrot:** Placing a carrot into a drawer and closing it. The carrot is initialized with a random rotation angle sampled uniformly from  $[-\pi/4, \pi/4]$  radians.
- **AGX-OpenPotCorn:** Opening a pot lid and placing corn inside. The pot is placed at a fixed location, while a random orientation is assigned to the corn, with the rotation angle sampled uniformly from the interval  $[-\pi/4, \pi/4]$  radians.
- **AGX-CloseDrawerCorn:** Picking up a corn and securely closing a drawer. The initial orientation of the corn is drawn from a uniform distribution over the interval  $[-\pi/4, \pi/4]$  radians.
- **AGX-PutCornPlate:** Placing corn pieces on a plate. The plate is placed at a fixed location, while the corn is initialized with a random rotation angle sampled uniformly from  $[-\pi/4, \pi/4]$  radians.
- **AGX-StackBowl:** Stacking bowls into a stable configuration. The position of the blue bowl is fixed, while the green bowl is randomly sampled from a  $15 \text{ cm} \times 15 \text{ cm}$  grid region.
- **AGX-UprightMug:** Restoring a tilted mug to an upright position and placing it on the plate. And the plate remains stationary, whereas the mug is sampled from a uniform distribution over a  $20 \text{ cm} \times 20 \text{ cm}$  region.
- **TK2-LayPlateBowl:** Taking a plate from the rack and laying a bowl on the plate. While the plate is fixed at the same position on the rack, the bowl is randomly sampled from a  $15 \text{ cm} \times 15 \text{ cm}$  grid region.
- **TK2-SelectButton:** Selecting a yellow button and placing it on the plate. The position of the plate is fixed, while the yellow and red button are randomly sampled from a  $15 \text{ cm} \times 15 \text{ cm}$  grid region.
- **TK2-CollectBall:** Collecting tennis balls into the box. Two tennis balls are positioned on opposite sides of a box, with each ball randomly located within a designated  $10 \text{ cm} \times 10 \text{ cm}$  area.
- **TK2-WeighApple:** Placing the apple in the bowl and weighing them together. The electronic scale and bowl are fixed, whereas the apple is randomly placed within a  $10 \text{ cm} \times 10 \text{ cm}$  area.
- **TK2-HeatBread:** Putting the bread into the oven and closing the door. The bread is initialized at a random position on the plate

**Dataset Collection.** We collected RoboAug-D dataset using the HACTS teleoperation system [67] on five robot embodiments: single-arm Franka, single-arm UR-5e, dual-arm UR-5e, AgileX and TienKung2. To accommodate task-specific temporal structures, we tailored the keyframe extraction strategy to each task:

- For basic, short-horizon tasks (e.g., single-arm UR-5e), we annotated semantic events (initial, gripper-close, and

TABLE VI: **Quantitative results under the Dual-Factor Variation setting.** We report the average success rates (%) on both UR-5e and AgileX robots. The evaluation involves 5 unseen background textures combined with 10 task-irrelevant distractors, totaling 100 trials for each task.

Augmentation Method	UR-PutCornPlot	UR-MoveLemon	UR-StackBowl	UR-StoreCarrot	UR-OpenDrawerCorn	Average
No Aug	0.20	0.26	0.12	0.28	0.36	0.24
RoboEngine-T	0.25	0.38	0.20	0.36	0.55	0.35
RoboEngine-G	0.28	0.42	0.22	0.40	0.62	0.39
GenAug	0.30	0.44	0.20	0.46	0.68	0.42
<b>RoboAug</b>	<b>0.36</b>	<b>0.52</b>	<b>0.28</b>	<b>0.58</b>	<b>0.84</b>	<b>0.52</b>
	AGX-PutCornPlate	AGX-UprightMug	AGX-StackBowl	AGX-OpenPotCorn	AGX-CloseDrawerCorn	
No Aug	0.20	0.12	0.18	0.20	0.20	0.18
RoboEngine-T	0.26	0.18	0.22	0.25	0.22	0.23
RoboEngine-G	0.30	0.22	0.26	0.30	0.26	0.27
GenAug	0.28	0.20	0.28	0.26	0.32	0.27
<b>RoboAug</b>	<b>0.56</b>	<b>0.36</b>	<b>0.40</b>	<b>0.40</b>	<b>0.56</b>	<b>0.46</b>

TABLE VII: Performance Comparison of different methods under a set of 10 distinct distractors.

Augmentation Method	UR-PutCornPlot	UR-MoveLemon	UR-StackBowl	UR-StoreCarrot	UR-OpenDrawerCorn	Average
No Aug	0.25	0.20	0.30	0.10	0.25	0.22
RoboEngine-T	0.50	0.25	0.50	0.20	0.35	0.36
RoboEngine-G	0.55	0.30	0.50	0.20	0.40	0.39
GenAug	0.60	0.30	0.55	0.20	0.45	0.42
<b>RoboAug</b>	<b>0.90</b>	<b>0.45</b>	<b>0.60</b>	<b>0.40</b>	<b>0.50</b>	<b>0.57</b>
	AGX-PutCornPlate	AGX-UprightMug	AGX-StackBowl	AGX-OpenPotCorn	AGX-CloseDrawerCorn	
No Aug	0.15	0.25	0.25	0.15	0.00	0.16
RoboEngine-T	0.20	0.30	0.40	0.30	0.15	0.27
RoboEngine-G	0.25	0.35	0.45	0.30	0.20	0.31
GenAug	0.30	0.35	0.55	0.25	0.10	0.31
<b>RoboAug</b>	<b>0.45</b>	<b>0.50</b>	<b>0.90</b>	<b>0.40</b>	<b>0.30</b>	<b>0.51</b>

gripper-open frames).

- For complex, long-horizon tasks (e.g., dual-arm UR-5e and AgileX), we employed uniform sampling at 50-frame intervals.

All keyframes feature manual annotations of task-relevant entities, including manipulated objects and robot end-effectors.

### B. Dual-Factor Generalization Evaluation

**Evaluation Setup.** To rigorously assess the robustness of visual policies under complex environmental shifts, we introduce the *Dual-Factor Variation* protocol. This setting challenges the agent with a combination of two distinct perturbations: unseen background textures and object clutter. Formally, we evaluate the model using five background textures that were not present in the training set. For each background, we introduce 10 task-irrelevant distractors placed randomly across the workspace. This configuration aims to verify whether the policy can effectively decouple task-essential features from compounded visual distractions. To validate the cross-embodiment stability of RoboAug, we conduct these experiments on two distinct robot embodiments, the UR-5e (UR) and the AgileX (AGX). For each background-clutter configuration, we perform 20 evaluation trials, resulting in a total of 100 trials, and report the average success rate.

**Results under Dual-Factor Variation.** Quantitative results are summarized in Table VI. The complexity of this setting poses a substantial hurdle for existing methods. The baseline

model, ACT w/o Aug, fails to generalize, yielding average success rates of only 0.24 and 0.18 on the UR and AGX suites, respectively. While state-of-the-art augmentation methods such as RoboEngine-G and GenAug offer moderate improvements by addressing individual visual factors, their performance degrades significantly when facing simultaneous background and object shifts.

In contrast, our proposed RoboAug exhibits superior robustness across all evaluated benchmarks. On the UR-series tasks, RoboAug achieves an average success rate of 0.52, surpassing the strongest baseline GenAug by a relative margin of 10%. A consistent trend is observed in the AGX-series, where our method attains an average score of 0.46. These results demonstrate that RoboAug effectively synthesizes a diverse training distribution that captures the underlying visual logic, enabling the model to maintain high robustness even under high-variance dual-factor perturbations.

### C. Single-Factor Generalization: Robustness to Task-Irrelevant Distractors.

Complementing the background generalization experiments presented in the main manuscript, we further established a single-factor generalization evaluation focused on task-irrelevant distractors. For each task, we randomly placed 10 task-irrelevant objects on the tabletop and conducted 20 evaluation rollouts. Table VII reports performance under this extreme clutter setting. In this challenging scenario, baseline

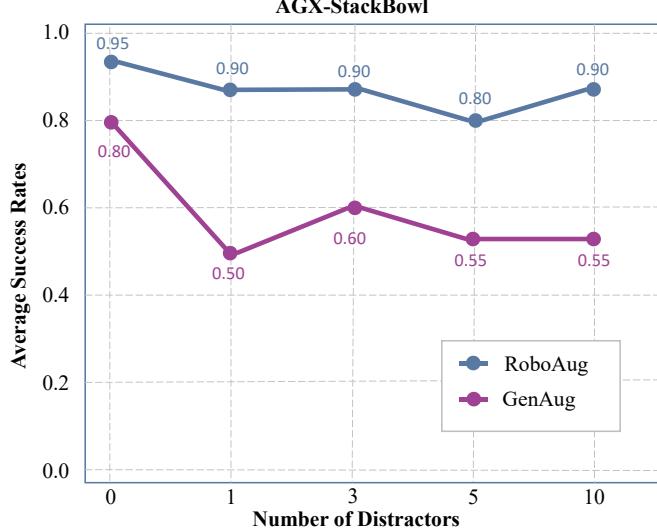


Fig. 11: Comparison of task success rates between RoboAug and the best baseline method under varying numbers of distractors.

methods exhibit substantial performance degradation. Specifically, RoboEngine-T and ACT w/o Aug achieve success rates of only 0.15 and 0.20 on the AGX-PutCornPlate task, whereas RoboAug attains a markedly higher success rate of 0.45.

Failure inspection indicates that baseline methods often struggle to semantically distinguish the target object from dense background clutter, leading to incorrect object selection or unstable execution. In contrast, by leveraging the proposed region-contrastive loss, RoboAug learns more discriminative region-level representations, enabling the policy to consistently attend to the task-relevant object and maintain robust performance under heavy visual distraction.

To investigate the impact of clutter density, Figure 11 illustrates performance trends as the number of distractors increases from 0 to 10. On the AGX-StackBowl task, RoboAug consistently outperforms the strongest baseline, GenAug, across all clutter levels. Notably, RoboAug preserves a high success rate of 0.90 even with 10 distractors, while GenAug experiences a pronounced drop from 0.80 to 0.55. These results demonstrate that RoboAug is substantially more robust to task-irrelevant visual perturbations.

#### D. Single-Factor Generalization: Robustness to Illumination Variation.

As illustrated in Figure 12, we further evaluate policy generalization across 10 tasks on the UR-5e and AgileX robots under 20 unseen lighting conditions, including 4 types of dynamic lighting changes. Under these conditions, the baseline method GenAug exhibits notable performance degradation, particularly in scenarios involving drastic color temperature shifts, strong cast shadows, and dynamic high-contrast illumination. Such lighting variations alter the apparent shape

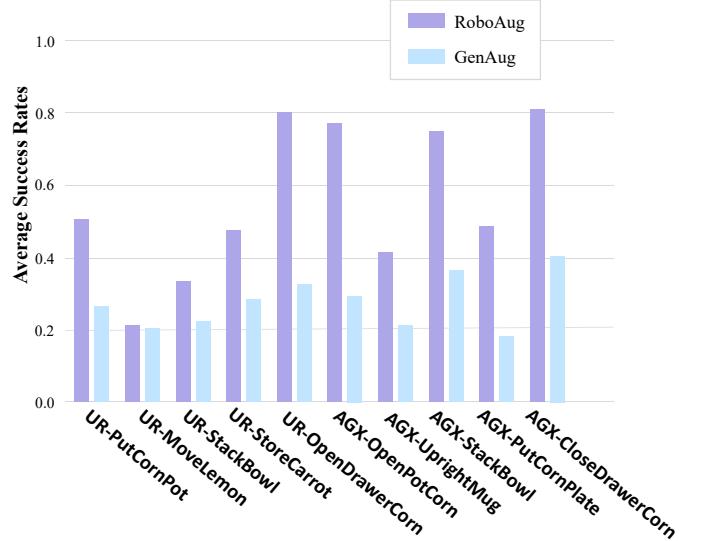


Fig. 12: Comparison of RoboAug and best baseline method on 20 lighting conditions, evaluated using the task success rates.

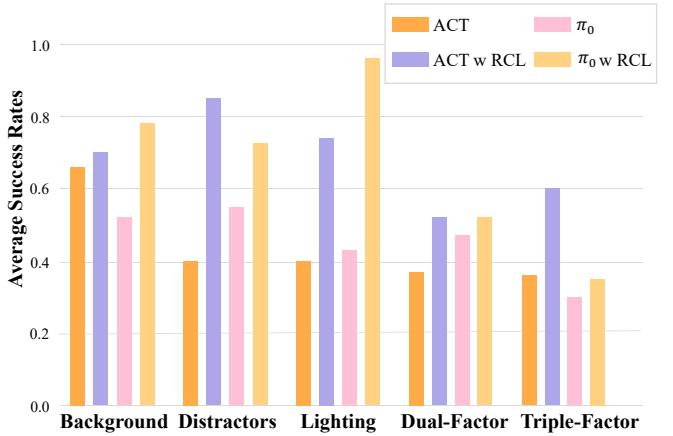


Fig. 13: Region-contrastive loss based on ACT [77] and  $\pi_0$  [4].

and texture of objects, frequently leading to perception and recognition failures.

In contrast, RoboAug demonstrates substantially more stable performance across the full range of lighting conditions. For example, on the AGX-OpenPotCorn task, RoboAug achieves a success rate of 0.75, significantly outperforming GenAug, which attains only 0.30. We attribute this improvement to the proposed generative augmentation pipeline combined with region-contrastive policy learning, which systematically exposes the policy to diverse lighting variations during training. As a result, the learned visual representations are less sensitive to illumination changes, enabling more reliable task execution in unseen lighting environments.

#### E. Effectiveness of Region-Contrastive Loss across Different Policies

To verify the universality of our proposed method, we evaluate the impact of the Region-Contrastive Loss (RCL) when integrated into different policy backbones. Specifically, we apply RCL to two representative architectures, ACT [77]

TABLE VIII: Tasks Success Rates on Novel Backgrounds. We evaluated the success rate of augmentation methods across 10 unseen backgrounds.

Method	UR-PutCornPlot	UR-MoveLemon	UR-StackBowl	UR-StoreCarrot	UR-OpenDrawerCorn	Average
ACT w/o Aug	0.24	0.15	0.12	0.22	0.34	0.21
RoboEngine-T	0.24	0.26	0.22	0.20	0.40	0.26
RoboEngine-G	0.46	0.28	0.32	0.50	0.42	0.40
GenAug	0.48	0.36	0.28	0.46	0.42	0.40
<b>RoboAug</b>	<b>0.90</b>	<b>0.60</b>	<b>0.60</b>	<b>0.65</b>	<b>0.65</b>	<b>0.68</b>
	AGX-PutCornPlate	AGX-UprightMug	AGX-StackBowl	AGX-OpenPotCorn	AGX-CloseDrawerCorn	
ACT w/o Aug	0.16	0.18	0.12	0.28	0.12	0.17
RoboEngine-T	0.22	0.20	0.38	0.20	0.30	0.26
RoboEngine-G	0.32	0.35	0.46	0.30	0.32	0.35
GenAug	0.36	0.30	0.47	0.40	0.42	0.39
<b>RoboAug</b>	<b>0.74</b>	<b>0.84</b>	<b>0.70</b>	<b>0.52</b>	<b>0.72</b>	<b>0.70</b>

and  $\pi_0$  [4], within the AGX-StackBowl task. As illustrated in Figure 13, the experimental results demonstrate that incorporating RCL consistently improves performance for both policies. These findings validate the effectiveness of RCL and suggest that it is compatible with diverse policy formulations.

#### F. Background Generalization across Multiple Embodiments

To further validate the effectiveness of RoboAug across different robot configurations and multiple tasks, we conducted extensive experiments focusing on single-factor generalization with respect to background variations. Specifically, we evaluated our method on five distinct tasks for both the UR-5e and AgileX robotic arms. For each task, we tested the policy on 10 different unseen backgrounds and reported the average success rate.

Table VIII presents the quantitative results for these experiments. As indicated by the data, our proposed method, RoboAug, demonstrates superior generalization capabilities compared to all baselines. In the UR-5e robot tasks, RoboAug achieves an average success rate of 68%, significantly surpassing the strongest baseline, GenAug, which achieves 40%. A similar trend is observed in the AgileX robot tasks, where our method reaches a 70% success rate, consistently outperforming other augmentation strategies.

#### G. Theoretical Analysis of RoboAug

To provide a theoretical foundation for RoboAug, we analyze the generalization error bound using Rademacher complexity [1, 26]. We demonstrate that our method improves generalization through two synergistic mechanisms: increasing the effective sample size via semantic augmentation and reducing the hypothesis space complexity via region-contrastive learning.

1) *Preliminaries and Definitions:* Let  $\mathcal{X}$  and  $\mathcal{Y}$  denote the input observation space and action space, respectively. We assume the data is drawn from an underlying distribution  $\mathcal{D}$ . A policy is a function  $\pi : \mathcal{X} \rightarrow \mathcal{Y}$  chosen from a hypothesis class  $\mathcal{H}$ . The goal is to minimize the expected risk  $\mathcal{R}(\pi) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\ell(\pi(\mathbf{x}), \mathbf{y})]$ , where  $\ell$  is a bounded continuous loss function.

In RoboAug, we suppose that the observation  $\mathbf{x}$  can be decomposed into task-relevant regions  $\mathbf{x}_{\text{task}}$  and task-irrelevant scenario factors  $\mathbf{x}_{\text{scen}}$  (e.g., background, lighting). The ideal expert policy  $\pi^*$  depends solely on the task-relevant regions, such that  $\pi^*(\mathbf{x}_{\text{task}}, \mathbf{x}_{\text{scen}}) = \pi^*(\mathbf{x}_{\text{task}}, \mathbf{x}_{\text{scen}}^{\text{new}})$  for any variations in  $\mathbf{x}_{\text{scen}}$ .

2) *Analysis of Semantic Data Augmentation:* Standard generalization bounds depend heavily on the number of training samples  $N$ . We first recall the classical generalization bound based on Rademacher complexity [1].

*Theorem 6.1 (Generalization Bound for Loss Functions):* Let  $\mathcal{H}$  be the policy hypothesis class. Assume the loss function  $\ell$  is Lipschitz continuous with respect to its first argument with constant  $L_\ell$  and is bounded by  $c$ . For any  $\delta > 0$ , with probability at least  $1 - \delta$  over the draw of a dataset  $S$  of size  $N$ , the following inequality holds for all  $\pi \in \mathcal{H}$ :

$$\mathcal{R}(\pi) \leq \hat{\mathcal{R}}_N(\pi) + 2L_\ell \mathfrak{R}_N(\mathcal{H}) + c\sqrt{\frac{\log(1/\delta)}{2N}}, \quad (5)$$

where  $\hat{\mathcal{R}}_S(\pi)$  is the empirical risk on the dataset  $S$ , and  $\mathfrak{R}_N(\mathcal{H})$  is the Rademacher complexity of  $\mathcal{H}$  given  $N$  samples.

RoboAug expands the original expert dataset of size  $N$  to a significantly larger augmented dataset of size  $N_{\text{total}} = N + N_{\text{aug}}$  by generating diverse  $\mathbf{x}_{\text{scen}}$  while preserving  $\mathbf{x}_{\text{task}}$ . Assuming the augmented samples are valid (i.e., they share the correct action labels  $\mathbf{y}$  derived from the expert trajectories), this expansion leads to a tighter generalization bound.

*Theorem 6.2 (Generalization Bound with Semantic Augmentation):* Let  $S_{\text{total}}$  be the augmented dataset of size  $N_{\text{total}}$ . Under the assumptions of Theorem 6.1, with probability at least  $1 - \delta$ , for all  $\pi \in \mathcal{H}$ :

$$\mathcal{R}(\pi) \leq \hat{\mathcal{R}}_{S_{\text{total}}}(\pi) + 2L_\ell \mathfrak{R}_{N_{\text{total}}}(\mathcal{H}) + c\sqrt{\frac{\log(1/\delta)}{2N_{\text{total}}}}. \quad (6)$$

Crucially, since the Rademacher complexity for neural networks typically scales with  $\mathcal{O}(1/\sqrt{N})$ , and  $N_{\text{total}} \gg N$ , we have:

$$2L_\ell \mathfrak{R}_{N_{\text{total}}}(\mathcal{H}) + c\sqrt{\frac{\log(1/\delta)}{2N_{\text{total}}}} \ll 2L_\ell \mathfrak{R}_N(\mathcal{H}) + c\sqrt{\frac{\log(1/\delta)}{2N}}. \quad (7)$$

This theorem formally justifies that by increasing the diversity and quantity of training data through semantic augmentation, RoboAug reduces the estimation error gap, allowing the empirical risk to better approximate the true expected risk.

*3) Analysis of Region-Contrastive Learning:* While augmentation increases the sample size of the training dataset, the Region-Contrastive Learning (RCL) objective improves generalization by effectively constraining the hypothesis class  $\mathcal{H}$ . RCL enforces feature invariance with respect to task-irrelevant regions  $\mathbf{x}_{\text{scen}}$ .

Let  $\mathcal{H}_{\text{inv}} \subseteq \mathcal{H}$  denote the subset of policies that are invariant to variations in  $\mathbf{x}_{\text{scen}}$ , defined as  $\mathcal{H}_{\text{inv}} = \{\pi \in \mathcal{H} \mid \pi(\mathbf{x}_{\text{task}}, \mathbf{x}_{\text{scen}}) = \pi(\mathbf{x}_{\text{task}}, \mathbf{x}_{\text{scen}}^{\text{new}}), \forall \mathbf{x}_{\text{scen}} \in \mathcal{X}, \mathbf{x}_{\text{scen}}^{\text{new}} \in \mathcal{X}\}$ . The region-contrastive loss minimizes the distance between representations of the same task-relevant objects against different backgrounds, effectively regularizing the search space towards  $\mathcal{H}_{\text{inv}}$ .

*Corollary 6.2.1 (Complexity Reduction via RCL):* Since  $\mathcal{H}_{\text{inv}}$  is a proper subset of  $\mathcal{H}$ , its Rademacher complexity is strictly lower:

$$\mathfrak{R}_{N_{\text{total}}}(\mathcal{H}_{\text{inv}}) \leq \mathfrak{R}_{N_{\text{total}}}(\mathcal{H}). \quad (8)$$

Consequently, by optimizing the policy within this constrained invariant subspace, RoboAug further tightens the generalization bound:

$$\mathcal{R}(\pi_{\text{RCL}}) \leq \hat{\mathcal{R}}(\pi_{\text{RCL}}) + \underbrace{2L_\ell \mathfrak{R}_{N_{\text{total}}}(\mathcal{H}_{\text{inv}})}_{\text{Reduced Complexity}} + c \sqrt{\frac{\log(1/\delta)}{2N_{\text{total}}}}. \quad (9)$$

RoboAug achieves robust generalization by simultaneously reducing the error bound from two directions: expanding the denominator of the complexity term via *Semantic Augmentation* ( $N \rightarrow N_{\text{total}}$ ) and reducing the Rademacher complexity via *Region-Contrastive Learning* ( $\mathcal{H} \rightarrow \mathcal{H}_{\text{inv}}$ ).

#### H. Instantiations of Generalization Factors

In this section, we detail the specific configurations used in our three single-factor generalization experiments, covering background variations, distractor interference, and lighting conditions.

**Background Generalization.** We evaluated the policy on the UR-CornPot task using 170 distinct, unseen background textures. Based on visual complexity, we categorized these backgrounds into three types. Visualizations of these 170 background instances are provided in Figures 14 to 18.

- Regular Geometric: This category comprises 32 patterns characterized by basic shapes and lines, such as squares and rhombuses, arranged in a strictly ordered layout.
- Structured Design: This category consists of 35 patterns featuring more intricate motifs, including flowers, leaves, and animals. These patterns maintain a regular, tiled arrangement.
- Scattered Pattern: This category includes 103 complex patterns, such as toys and irregular polygons. Unlike the previous categories, these are distributed randomly

without a fixed grid, significantly increasing visual interference.

**Distractor Generalization.** As illustrated in Figure 19, we introduced visual clutter to test the model’s robustness against obstacles. We placed up to 10 distractor objects on the tabletop, effectively occupying the entire workspace. This setup creates a highly cluttered environment that poses a substantial challenge to the policy.

**Lighting Generalization.** As shown in Figure 20, we assessed the model’s performance under 20 different illumination settings. This set includes diverse static conditions as well as three scenarios involving dynamic lighting changes to evaluate adaptability.



Fig. 14: Visualization of the Regular Geometric background category.



Fig. 15: Visualization of the Structured Design background category.



Fig. 16: Visualization of the Scattered Pattern background category.



Fig. 17: Visualization of the Scattered Pattern background category.



Fig. 18: Visualization of the Scattered Pattern background category.

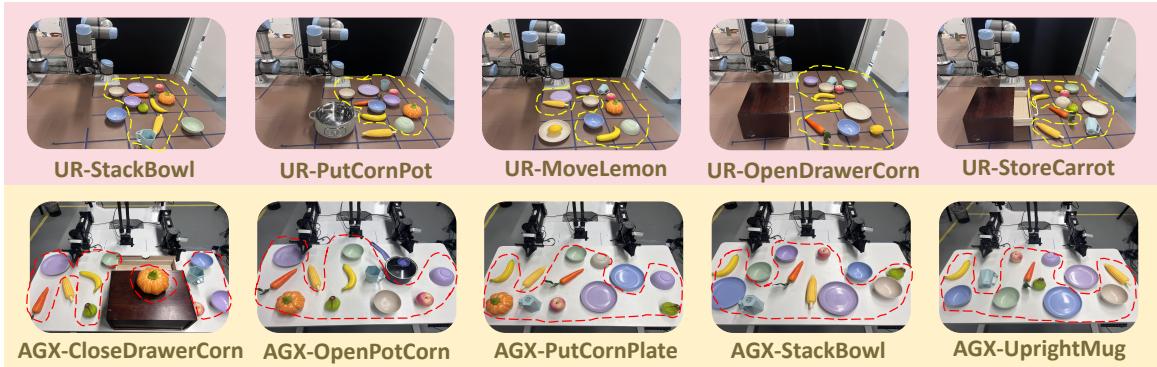


Fig. 19: Visualization of the ten distinct distractors used in the UR and AgileX tasks.



Fig. 20: Visualization of twenty distinct illumination conditions. The bottom row demonstrates dynamic lighting scenarios with multicolor changes at varying speeds.