

ObjectVLA: End-to-End Open-World Object Manipulation Without Demonstration

Minjie Zhu^{12*} Yichen Zhu^{1*†} Jinming Li³ Zhongyi Zhou²
Junjie Wen² Xiaoyu Liu³ Chaomin Shen² Yaxin Peng³ Feifei Feng¹
¹Midea Group ²East China Normal University ³Shanghai University
*Equal contribution †Corresponding author

objectvla.github.io

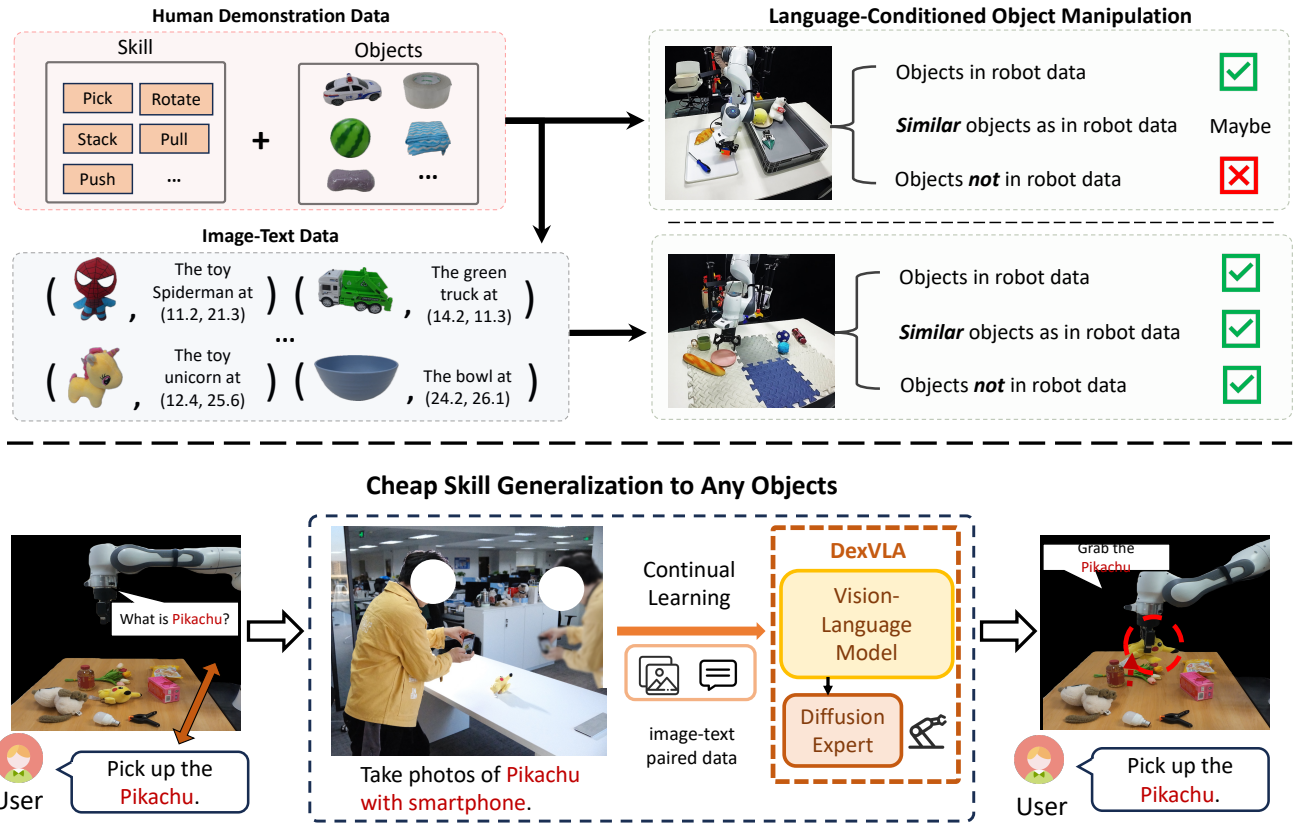


Figure 1. **A brief illustration of ObjectVLA.** Conventional imitation learning focuses on tasks that involve both skills and objects. While it performs well on seen objects and sometimes generalizes to similar ones (e.g., objects with changed colors), it typically fails with novel objects. By co-training with image-text data, our approach enables VLA models to generalize to any object present in the image-text dataset. **Additionally**, users can capture object images, automatically generate image-text data, and fine-tune a pre-trained VLA model with minimal resources to learn manipulation on novel objects.

Abstract

Imitation learning has proven to be highly effective in teaching robots dexterous manipulation skills. However, it typically relies on large amounts of human demonstration data, which limits its scalability and applicability in dynamic,

real-world environments. One key challenge in this context is object generalization—where a robot trained to perform a task with one object, such as “hand over the apple,” struggles to transfer its skills to a semantically similar but visually different object, such as “hand over the peach.” This gap in generalization to new objects beyond those in the

same category has yet to be adequately addressed in previous work on end-to-end visuomotor policy learning. In this paper, we present a simple yet effective approach for achieving object generalization through Vision-Language-Action (VLA) models, referred to as **ObjectVLA**. Our model enables robots to generalize learned skills to novel objects without requiring explicit human demonstrations for each new target object. By leveraging vision-language pair data, our method provides a lightweight and scalable way to inject knowledge about the target object, establishing an implicit link between the object and the desired action. We evaluate ObjectVLA on a real robotic platform, demonstrating its ability to generalize across 100 novel objects with a 64% success rate in selecting objects not seen during training. Furthermore, we propose a more accessible method for enhancing object generalization in VLA models—using a smartphone to capture a few images and fine-tune the pre-trained model. These results highlight the effectiveness of our approach in enabling object-level generalization and reducing the need for extensive human demonstrations, paving the way for more flexible and scalable robotic learning systems.

1. Introduction

Vision-language-action (VLA) models have emerged as a transformative paradigm for teaching robots dexterous skills, enabling them to replicate human behavior and master complex tasks [3–5, 27, 37]. However, a critical limitation persists: these models rely heavily on human demonstration data, which constrains their scalability and practicality in dynamic real-world environments [17, 33, 34]. For instance, a robot trained to execute “hand over the apple” often fails to generalize to analogous tasks like “hand over the peach,” despite conceptual similarity. This underscores the unresolved challenge of **object generalization** — adapting learned skills to novel, unseen objects — particularly when such objects lie **beyond the category of the teleoperated training data**. We name these objects as out-of-distribution (OOD) objects.

The core limitation stems from imitation learning’s tendency to learn fixed mappings from instruction and visual input to action. When encountering objects absent from teleoperation data, the model lacks mechanisms to associate the object’s name, visual features, and learned actions. To address this, we propose a framework that bridges visual-language semantics and robotic actions through localization-aware reasoning.

Our approach begins by curating a dataset of image-text pairs augmented with localization metadata (e.g., bounding boxes). This dataset is co-finetuned with teleoperated robot interaction data, while the robot data itself is enriched with localization-guided reasoning. By embedding localization

as a bridging representation, we create a unified pathway between visual-language inputs and robotic actions. This enables zero-shot object generalization: the model can recognize and manipulate novel objects—even those absent from robot training data—without task-specific retraining.

We designed rigorous real-robot experiments to validate the generalization capabilities of our framework, ObjectVLA. In these trials, six objects are positioned at distinct locations (left or right side of a table), with configurations spanning combinations of objects seen in robot interaction data or vision-language data. The robot is tasked with the instruction “move to the object”, achieving a 100% success rate for in-domain objects. To stress-test generalization, we evaluated 100 OOD objects, observing a 64% success rate. These experiments demonstrate that our method adapts to diverse novel object types when trained with vision-language priors.

The versatility of our approach is further demonstrated across diverse scenarios, including bin-picking and tasks requiring composite skills like pushing and rotating. Notably, our framework supports rapid adaptation to novel objects: by collecting smartphone-captured images and performing lightweight fine-tuning, the model generalizes to objects absent from the original dataset. These experiments underscore our method’s ability to reduce reliance on large-scale human demonstrations while achieving robust object generalization.

Our primary contribution is a unified pipeline for integrating vision-language datasets with robot interaction data, enabling end-to-end object generalization. Through systematic evaluation, we validate the framework’s performance on complex multi-stage tasks (e.g., bin-picking) and multi-skill manipulation (e.g., rotating, pushing), highlighting its universality. Despite some of the existing works, such as RT-2 [5] and ECoT [37] giving a glimpse of how co-finetuning can achieve simple object generalization, they neither elucidate the underlying mechanism of achieving such generalization nor address the boundary of their methodologies. In contrast, our approach — though simple and straightforward — demonstrates that training VLA models with a hybrid dataset of robot interaction data and image-text data significantly enhances generalization. This level of generalization goes significantly beyond previously demonstrated end-to-end approaches. Crucially, our framework enables practical deployment: even a small set of smartphone images and brief fine-tuning suffices to adapt the model to novel objects, significantly advancing real-world robotic flexibility.

2. Related Work

Vision-language-action models for robot control. Recent research has focused on developing generalist robot policies trained on increasingly expansive robot learning

datasets [9, 10, 16, 21, 26]. Vision-language-action models (VLAs) represent a promising approach for training such generalist policies [3, 8, 17, 25, 27, 35, 38, 40]. VLAs adapt vision-language models (VLMs) [1, 7, 14, 22–24, 32, 39, 41, 44], pre-trained on vast internet-scale image and text data, for robotic control. This approach offers several advantages: leveraging large vision-language model backbones, with billions of parameters, provides the necessary capacity for fitting extensive robot datasets. Furthermore, reusing weights pre-trained on internet-scale data enhances the ability of VLAs to interpret diverse language commands and generalize to novel objects and environments. However, current VLA models struggle to recognize open-world objects when these objects are absent from the robot interaction data [17, 34]. This is mainly due to VLMs essentially “overwrites” its previously acquired knowledge of open-world objects with robot-specific information.

Generalization in robot learning. In the realm of robot learning, generalization, particularly object generalization, remains a core challenge and active area of research. Many works leverage techniques such as domain randomization [13], meta-learning [15, 29], retrieval-augmented generation [43], extra modality [36, 45], and data augmentation to improve a robot’s ability to recognize and interact with novel objects unseen during training. For instance, domain randomization methods [13, 31] randomize visual and physical parameters during simulation training to force the agent to learn features invariant to these irrelevant details, leading to better real-world generalization. Furthermore, meta-learning approaches [11] aim to train models that can rapidly adapt to new objects with limited data, directly addressing the object generalization problem. Finally, data augmentation methods [18, 19], enhance the diversity of the training data, exposing the model to a wider range of object appearances and orientations, thereby promoting robustness and generalization to novel objects. There is also a field of work using large language models or vision-language models to do open-vocabulary manipulation [2, 20, 30, 42], combined with motion planning and robot learning methods. However, these approaches involve separate modules that are trained independently for different components. To the best of our knowledge, this work represents the first exploration of object generalization beyond specific categories within visuomotor policy learning.

3. Methodology

3.1. Notation and Motivation

Given a set of expert demonstrations that contain complex robot skill trajectories, we want to learn a visuomotor policy $\pi : \{\mathcal{O}_r, \mathcal{I}_r\} \mapsto \mathcal{A}$ that maps the visual observations $o_r \in \mathcal{O}_r$ and the language instruction $i_r \in \mathcal{I}_r$ to actions $a \in \mathcal{A}$. The action changes accordingly when the lan-

guage instruction and visual input change. The r denote the data in the human demonstration data. Typically, for each language instruction it contains robot skill such as “push” or “pick up” and the target object, which is denoted as $\{obj_r, skill_r\} \in i_r$. We then formally define the image-text data, where $\varphi : \{\mathcal{O}_v, \mathcal{I}_v\} \mapsto \mathcal{L}_v$, where we input the image $o_v \in \mathcal{O}_v$ and give a language instruction $i_v \in \mathcal{I}_v$, the model is output with the corresponding answer $l_r \in \mathcal{L}_v$. The notation v denotes image-text data.

In this work, we explore the generalization of objects, focusing on those that are not part of the robot interaction data but are present in image-text data.

3.2. Data Construction

image-text data construction. To explore the model’s ability to generalize to novel objects, we constructed a diverse image-text dataset. For the visual component, we collected 100 distinct objects that are not included in the robot interaction objects. Specifically, using three cameras mounted on the robot (see Figure 2), we captured 20 images per object, covering various poses and orientations to ensure diversity. For the textual component, we employ a fixed template, “Detecting the bounding box of object.”, as the question, and the corresponding bounding box as the answer. In total, our vision-language dataset comprises 2,000 image-text pairs.

Reasoning data construction. We utilize localization metadata to bridge the gap between image-text data and robot data, as previously mentioned. To establish this implicit link between image-text and action, we incorporate localization metadata into the robot data. This section details how we construct reasoning with localization for robot data.

For each task, we first identify target objects based on the language instructions. We then employ DinoX [28], a cutting-edge open-vocabulary object detector, to annotate the bounding boxes of these objects. DinoX can generate a bounding box given an object’s name. To ensure accuracy, we manually verify and correct any erroneous bounding boxes produced by DinoX. Since our workspace has two external camera views, which can result in different bounding boxes for the same object, we annotate only one (right camera in our experiments). Following Qwen2-VL [32], we use a fixed template, “<|object_ref_start|>{object}<|object_ref_end|><|box_start|>(x₁, y₁),(x₂, y₂)<|box_end|>.”, to represent the localization reasoning. This reasoning is generated before each action and injected into the policy model through a learnable module. For a detailed explanation of this injection module’s architecture, we refer readers to DiVLA [33], the base model used in our experiments. An example of constructed image-text data is at Figure 3.

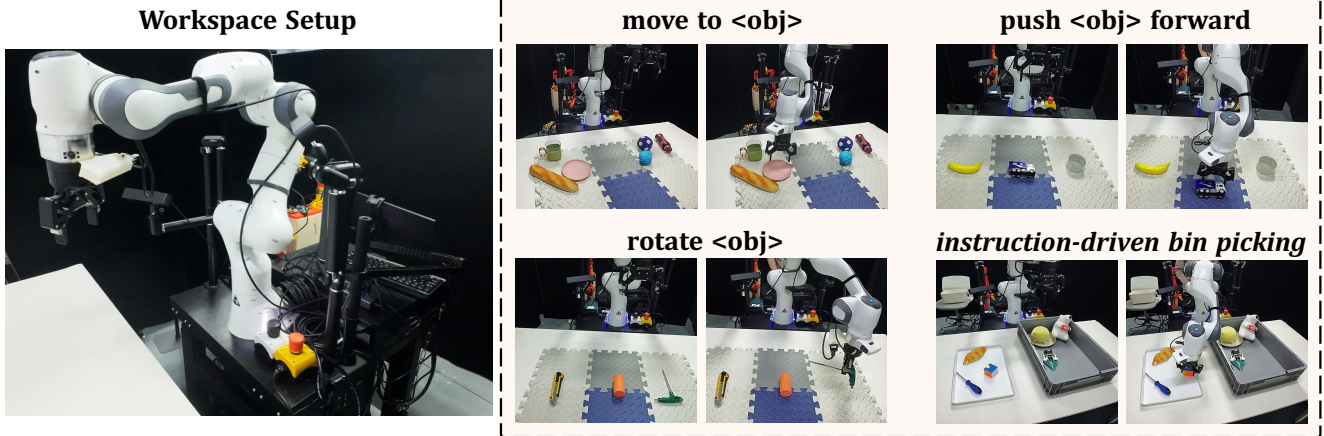


Figure 2. **Robot setup and examples for real-world manipulation tasks.** We evaluate ObjectVLA with 4 skills on a Franka robot arm equipped with two external Zed cameras and a Realsense 435i wrist camera.

3.3. Training Strategy and Implementation Details

For our experiments, we adopt diffusion-based VLA, a widely used class of Vision-Language-Action (VLA) models exemplified by methods like π_0 [3] and TinyVLA [34]. We select diffusion-based VLA over auto-regressive alternatives due to its significantly faster inference speed, a critical advantage for real-time robotic applications (see FAST [27] for a detailed comparison). Specifically, we utilize DiVLA [33], a representative VLA architecture, co-train on a hybrid dataset comprising robot interaction data and the vision-language corpus. To balance task-specific adaptation and semantic generalization, we maintained a 10:1 data ratio (robot-to-image-text data) across all tasks. This ratio empirically proved sufficient for robust object generalization, aligning with prior findings on the benefits of co-training for VLA capabilities. Notably, increasing the proportion of robot data beyond this ratio led to a decline in in-domain task success rates. We hypothesize this stems from the limited capacity of the 2B-parameter DiVLA model compared to larger architectures like ECoT (7B) [37] and RT-2 (55B) [5], which can better absorb domain-specific data without overfitting.

4. Experiments

In this section, we examine the effectiveness of ObjectVLA for object generalization in embodied control. In section 4.1, we verify the effectiveness of our method in object generalization. In section 4.2 and 4.3, we illustrate how our model transfers skills to objects not present in robot interaction data but included in the vision-language corpus. In section 4.4, we show that even a small set of smartphone images and brief fine-tuning can effectively adapt the pre-trained model to novel objects.

Real robot setup. All experiments are conducted on a

Franka robot [12] equipped with a 7-degree-of-freedom arm and a gripper. We use two external ZED cameras and a wrist Realsense 435i camera to obtain real-world visual information. Our real-world robot setup is illustrated in Figure 2.

4.1. Validating Object Generalization

In this section, we conduct rigorous experiments to verify the object generalization capability of our method. We begin by describing the experimental setup and evaluation criteria. Next, we evaluate ObjectVLA on both in-distribution and out-of-distribution objects. Finally, we explore several interesting observations related to object generalization.

4.1.1. Experimental Setup

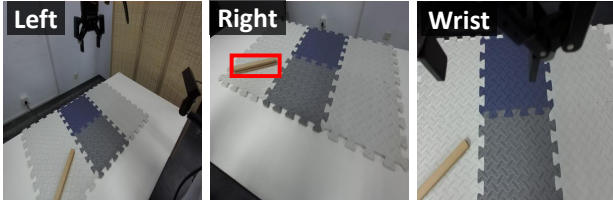
To verify the object generalization capability, we begin with a simple yet effective task, “Move to the object.”. In this task, we position objects on both sides of the robot, ensuring that each side has at least three objects on the table. The model is required to move toward the target object based on the given instruction. These objects are randomly chosen from a diverse set. For in-distribution (ID) evaluation, objects are only selected from the robot’s training data. And, for out-of-distribution evaluation, objects are randomly selected from either the robot’s training data or the vision-language data. A complete list of objects from both datasets is provided in the Appendix.

Evaluation criterion. We evaluate each object over 4 trials, with the target area’s side switching every two trials. We consider the model to have successfully recognized a novel object if and only if it moved toward the target object in all four trials. This criterion ensures that the model cannot achieve success simply by chance.

Experimental Results. Figure 5 presents the real-world experimental results for the “Move” task. Our ObjectVLA achieves a 100% success rate in ID evaluation. In the stress-

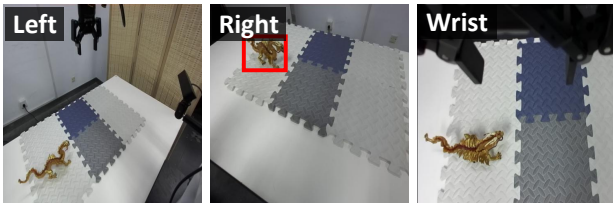
Example of Image-Text Paired Data

Photo taken by cameras from robot



Question Detecting the bounding box of *stick*.

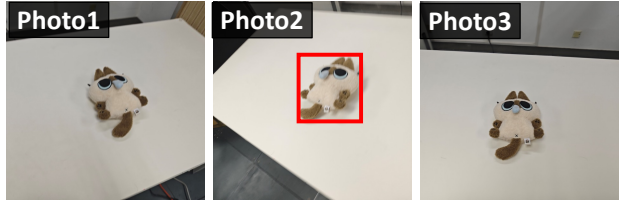
Answer <|object_ref_start|>stick<|object_ref_end|>
<|box_start|>(546,106),(557, 143)<|box_end|>.



Question Detecting the bounding box of *yellow dragon*.

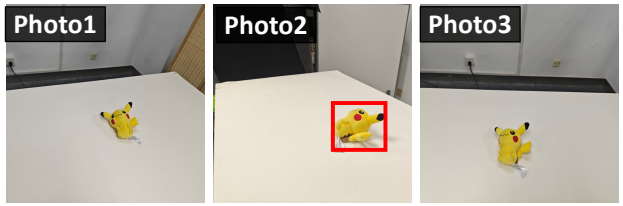
Answer <|object_ref_start|>yellow dragon<|object_ref_end|>
><|box_start|>(232,150),(389, 331)<|box_end|>.

Photo taken by smart-phone



Question Detecting the bounding box of *brown toy cat*.

Answer <|object_ref_start|>toy cat<|object_ref_end|>
<|box_start|>(421, 272),(657, 599)<|box_end|>.



Question Detecting the bounding box of *pikachu*.

Answer <|object_ref_start|>Pikachu<|object_ref_end|>
<|box_start|>(503, 523),(702, 694)<|box_end|>.

Figure 3. **Example of constructed image-text data.** *Left*: Photo taken by the robot’s camera. *Right*: Object captured with a smartphone.

test evaluation, our model successfully recognizes 64% of objects that are not present in the robot interaction data, confirming the effectiveness of co-training robot data with localization metadata.

Ablation study. To further understand our method’s effectiveness, we conducted an ablation study. We found that object generalization relies heavily on two key factors: first, explicitly linking vision and language to action through bounding boxes. This provides a direct connection between the visual object, its linguistic description, and the required manipulation. Second, a reasoning process for the robot data should be designed that mirrors the structure of vision-language pair data. This allows the model to leverage the rich information encoded in pre-trained vision-language models.

To analyze the impact of these factors, we removed the reasoning module for robot data and eliminated bounding boxes for vision-language data. The VLA model is then co-finetuned with vision-language data and evaluated using the same criteria and test settings as our full method.

As illustrated in Figure 5, the model without bounding boxes achieves only a 19% success rate in OOD evaluation, representing a significant performance decline compared to our method, despite achieving a 100% success rate in the

ID test. This suggests that without explicit grounding and a structured reasoning process, the model struggles to differentiate objects in vision-language data, leading to confusion about object-instruction correspondence and appropriate action selection.

4.1.2. More Observations

Can VLA recognize unseen objects if only trained with teleoperated data? To further assess the importance of vision-language data, we evaluated a VLA model trained exclusively on robot data, without any vision-language co-finetuning. As shown in Figure 5, this model (DiVLA) achieved 8% accuracy, which is almost equivalent to random guessing. This stark outcome highlights the critical role of vision-language data in multimodal understanding.

While the VLA model’s backbone is pre-trained on internet-scale vision-language data, focusing solely on robot data during training leads to catastrophic forgetting. The model essentially “overwrites” its previously acquired knowledge of visual concepts with robot-specific information, hindering its ability to comprehend multimodal scenes. Consequently, even objects encountered during pre-training, such as Pikachu, remain unrecognizable to the VLA model without vision-language co-finetuning.

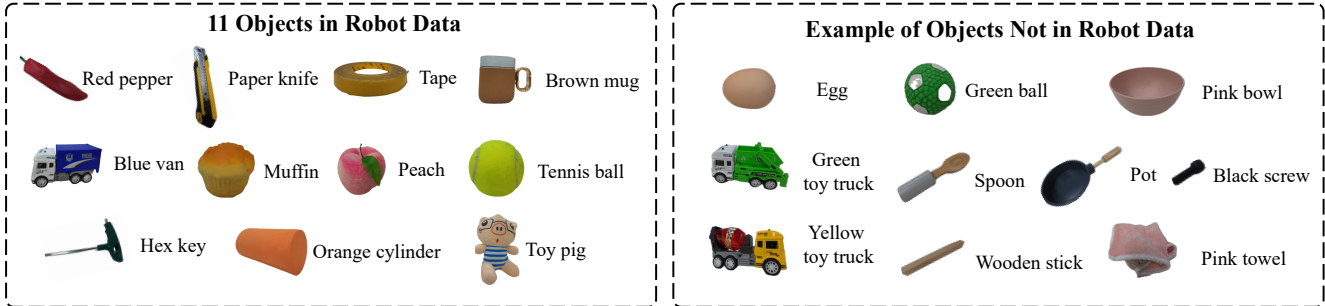


Figure 4. **Example Objects Used in Experiments.** *Left:* Objects present in the robot training data. *Right:* Examples of novel objects, not present in the robot data, but included in the image-text co-training dataset (see Appendix for a comprehensive list).

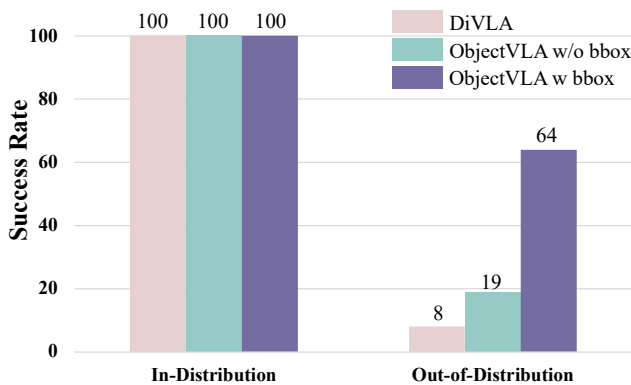


Figure 5. **Validation experiments on object generalization.** Our method achieved the best performance in both the in-distribution test setup and under visual changes. Each object is evaluated across 4 trials. We report the number of objects that were correctly identified in all four trials.

Table 1. **Experimental Results for rotate and push skills.** Our proposed ObjectVLA achieves high performance on both 5 in-distribution objects and 20 out-of-distribution objects. Each object is evaluated with three trials. We report the number of success trials.

Task	In-Distribution	Out-of-Distribution
Rotate	13/15	39/60
Push	12/15	52/60

4.2. Combining with More Skills

While the previous section employed a simple “move to” demonstration to validate the fundamental approach of our method, this section expands the evaluation to encompass more complex skills, specifically “push” and “rotate.” This broader assessment aims to demonstrate the generalizability of our method and its applicability beyond the “move to” task.

Experimental setup. In this experimental setup, we placed

three objects in front of the robot: one on the center, one on the right, and one on the left. The robot is instructed to either “rotate the object counterclockwise” or “push the object forward,” as illustrated in Figure 2. Following previous setup, we evaluate the model’s performance for both in-distribution (ID) and out-of-distribution (OOD) objects. Recognizing that some objects are inherently unsuitable for rotation or pushing actions (e.g., dishes), we conducted experiments on a curated set of 5 ID objects and 20 OOD objects. For each object in a skill, 40 demonstrations were collected, resulting in a total of 400 demonstrations.

Implementation details. We train one model for each skill to ensure that the model focuses more on understanding the objects rather than multi-task learning. We use the same image-text data of “move” task. Following established protocols from our prior work, this image-text dataset trained concurrently with the demonstration data for comprehensive evaluation. Each object was tested with 3 trials. In total, 150 trials were conducted. The training setting is provided in Appendix.

Results. As shown in Table 1, our method achieved high success rates on the robot interaction objects for both rotate and push skills. Analysis of the failed “rotate” trials revealed that the primary cause is the model’s inability to grasp the target object securely. When evaluating performance on out-of-distribution (OOD) objects, we observed a decrease in task completion rates compared to in-distribution objects, as expected. However, the model still successfully completed nearly two-thirds of the trials. Notably, in most failure cases, the model did not incorrectly identify the target object but rather failed to execute the skill completely. This was particularly evident in the “rotate” trials, where successful execution hinges on a secure grasp, a challenging requirement for unseen objects. Nevertheless, these experiments strongly support the claim that ObjectVLA can transfer learned skills, beyond basic pick and place, to novel objects within the framework we have developed. The results underscore the potential of ObjectVLA for generalized robotic manipulation, capable of adapting

Table 2. **Experimental results for bin picking.** Our proposed ObjectVLA achieves high performance on both 11 in-distribution objects and 50 out-of-distribution objects, with each object evaluated across 3 trials. We report the number of successful trials over total trials.

Method	In-Distribution	Out-of-Distribution
OpenVLA	14/33	17/150
ObjectVLA	21/33	87/150

to new objects and tasks beyond its initial training.

4.3. Instruction-Driven Bin Picking

To further evaluate ObjectVLA, we conducted experiments in a more practical scenario: end-to-end instruction-driven bin-picking. Unlike prior works (e.g., GR-2 [6] and DiVLA [33]) that execute bin-picking tasks without specific semantic instructions—typically limited to generic actions like transferring all objects from one container to another—we focus on a significantly more challenging setting [6, 33]. In our experiments, the robot is required to identify and retrieve a specific target object based on natural language instructions (e.g., "Pick the hexagonal bolt from the bin"). This scenario elevates the complexity of conventional bin-picking tasks by integrating cross-modal understanding (vision-to-language alignment) and fine-grained object discrimination that have multiple objects in the scene. Notably, the objects are randomly placed on the panel, which is a large area. Not only does the model need to figure out the object’s position, but also needs to be aware of its pose.

Implementation details. We collected new data within this environment. For robot interaction data, we collected 600 pick-and-place trajectories using the same "seen" objects as in previous experiments. For image-text data, we used half the number of objects from previous experiments, capturing 20 images of each. We compared our method against OpenVLA, a state-of-the-art VLA model, reporting success rates for both in-distribution and out-of-distribution objects. Evaluation consisted of three trials per object, totaling 183 trials per method. In each trial, at least two objects were randomly placed on the plate, and the model was instructed to pick and place a specific object according to the given instruction.

Results. Table 2 presents our experimental results. Bin picking, requiring object retrieval from random positions and poses, poses a significant challenge even for in-distribution objects. OpenVLA achieves a success rate of only 42.4% for in-distribution objects, significantly less than half. Surprisingly, it still completed roughly 10% of trials with out-of-distribution objects. This is likely due to some test objects sharing attributes with training objects (e.g., bread resembling a muffin, a green mug differ-

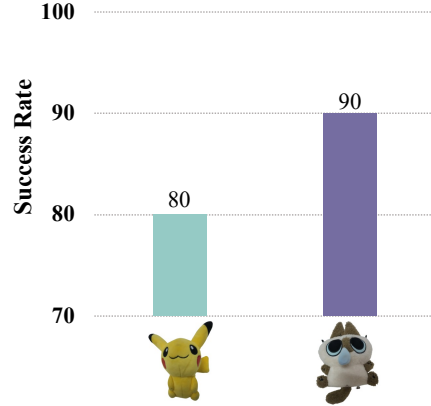


Figure 6. **Experimental results for smartphone captured objects and trained by continual learning.** We test two new objects. We took pictures of these two objects via smartphone and continually trained them on a pre-trained model. Each object was evaluated across 10 trials. We report the success rate for each object.

ing only in color from a brown training mug). In contrast, our method successfully completed 87 of 150 trials, including many completely novel objects, a 46.7% improvement over OpenVLA. This further emphasizes the necessity of co-training with both robot interaction and image-text data for effective object generalization.

4.4. Cheap Object Generalization via Smart-Phone Pictures and Continual Learning

The previous section demonstrated that our proposed co-training strategy enables the imitation learning method to generalize to any object by constructing corresponding image-text data for each object. This approach significantly enhances the model’s ability to handle a broader range of objects, without requiring extensive retraining. However, there are two key limitations that need to be addressed.

First, when a new object is introduced, the model must be trained from scratch to incorporate the new object into its understanding. This process can be highly cost-inefficient, as it involves retraining the model every time a novel object is added. In real-world applications, where objects are frequently introduced or changed, this limitation could significantly slow down deployment and increase operational costs. Second, our current image data is collected using cameras mounted on the robot, which ensures that there is no visual gap between the images captured by the robot’s cameras and the images input into the model. This setup works well in controlled environments but presents challenges in real-world scenarios. For instance, in order to capture the same images as the robot sees, you would need to replicate the exact camera positions and angles of the robot’s setup. This is not only cumbersome but also expen-

sive, as it requires building an identical system with matching camera views. Moreover, in environments where the robot is mobile or the scene is dynamic, maintaining consistent camera alignment becomes even more difficult.

Therefore, in this section, we test a simpler and more cost-effective approach: using a smart-phone camera to collect images from various perspectives. The model is then continuously trained on the pre-trained weights using this more accessible data collection method. As shown in Figure 6, we test with two objects: Pikachu and a brown toy cat. For each object, we capture 21 images and follow the same data construction pipeline discussed earlier. We train these objects in a bin-picking environment. Our results demonstrate that the model is able to recognize and successfully grasp the objects with a high success rate, 80% for Pikachu and 90% success rate for the toy cat. More importantly, we only need to continue training the model for 1 epoch. Because the collected data size is small, the training process can be extremely fast and can be finished up to ten minutes. This validates the effectiveness of our approach, showing that simple smartphone image collection combined with continuous learning enables open-world object manipulation in an end-to-end model. This experiment demonstrates that our method is flexible and cost-effective, making it a plug-and-play solution for existing VLA models, enabling them to generalize to virtually any object.

5. Conclusion

In this work, we present ObjectVLA, a Vision-Language-Action framework that addresses object generalization in robotic manipulation. By integrating vision-language datasets with robot interaction data, our method establishes a unified pipeline that bridges semantic understanding and physical action execution. This enables zero-shot generalization to over 100 novel objects with a 64% success rate, even when objects differ in category, appearance, or fine-grained attributes (e.g., color, shape). Our framework demonstrates that lightweight co-training with image-text priors and localization-aware reasoning can unlock robust cross-modal alignment. Key to our success is the ability to adapt rapidly to real-world scenarios: using just a few smartphone-captured images and quick continual fine-tuning, robots generalize to unseen objects without costly human demonstrations. We validate our approach across diverse tasks—including bin-picking, rotating and pushing—showcasing its versatility and practicality. Our results highlight a path toward scalable robotic learning systems that reduce dependence on large-scale teleoperation data while maintaining high performance.

6. Limitation

There are still a number of limitations in this work. Specifically, the image-text data was collected by the authors using either a robot-mounted camera or a smartphone. While we have not yet explored the feasibility of leveraging internet-sourced image-text data, it presents an intriguing avenue for future research. Specifically, investigating the necessary degree of visual similarity between internet images and target objects for effective skill transfer would be valuable. Our primary focus here is to introduce a novel pipeline that enables deep learning models to transfer skills to new objects without explicit demonstrations. Determining the limits of this transferability, particularly concerning the permissible visual gap between training and target objects, remains an open question for future investigation. Currently, our method struggles to generalize to novel backgrounds and lighting conditions. We believe the visual gap between our collected image-text data and the robot’s operational environment contributes to this challenge. Bridging this gap to improve generalization is a key focus for future development.

References

- [1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024. 3
- [2] 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. 3
- [3] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky. π_0 : A vision-language-action flow model for general robot control, 2024. 2, 3, 4
- [4] 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.
- [5] 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. 2, 4
- [6] 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*, 2024. 7

- [7] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198, 2024. 3
- [8] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023. 3
- [9] Sudeep Dasari, Oier Mees, Sebastian Zhao, Mohan Kumar Srirama, and Sergey Levine. The ingredients for robotic diffusion transformers. *arXiv preprint arXiv:2410.10088*, 2024. 3
- [10] Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, Chenxi Wang, Junbo Wang, Haoyi Zhu, and Cewu Lu. Rh20t: A comprehensive robotic dataset for learning diverse skills in one-shot. *arXiv preprint arXiv:2307.00595*, 2023. 3
- [11] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017. 3
- [12] Sami Haddadin. The franka emika robot: A standard platform in robotics research. *IEEE Robotics & Automation Magazine*, 2024. 4
- [13] Stephen James, Paul Wohlhart, Mrinal Kalakrishnan, Dmitry Kalashnikov, Alex Irpan, Julian Ibarz, Sergey Levine, Raia Hadsell, and Konstantinos Bousmalis. Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12627–12637, 2019. 3
- [14] Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models. *arXiv preprint arXiv:2402.07865*, 2024. 3
- [15] Rituraj Kaushik, Timothée Anne, and Jean-Baptiste Mouret. Fast online adaptation in robotics through meta-learning embeddings of simulated priors. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5269–5276. IEEE, 2020. 3
- [16] 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. 3
- [17] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 2, 3
- [18] Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020. 3
- [19] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. *Advances in neural information processing systems*, 33:19884–19895, 2020. 3
- [20] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023. 3
- [21] Fanqi Lin, Yingdong Hu, Pingyue Sheng, Chuan Wen, Jiacheng You, and Yang Gao. Data scaling laws in imitation learning for robotic manipulation, 2024. 3
- [22] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 3
- [23] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [24] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 3
- [25] Dantong Niu, Yuvan Sharma, Giscard Biamby, Jerome Quenum, Yutong Bai, Baifeng Shi, Trevor Darrell, and Roei Herzig. Llarva: Vision-action instruction tuning enhances robot learning. *arXiv preprint arXiv:2406.11815*, 2024. 3
- [26] 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. 3
- [27] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025. 2, 3, 4
- [28] Tianhe Ren, Yihao Chen, Qing Jiang, Zhaoyang Zeng, Yuda Xiong, Wenlong Liu, Zhengyu Ma, Junyi Shen, Yuan Gao, Xiaoke Jiang, et al. Dino-x: A unified vision model for open-world object detection and understanding. *arXiv preprint arXiv:2411.14347*, 2024. 3
- [29] Gerrit Schoettler, Ashvin Nair, Juan Aparicio Ojea, Sergey Levine, and Eugen Solowjow. Meta-reinforcement learning for robotic industrial insertion tasks. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9728–9735, 2020. 3
- [30] Austin Stone, Ted Xiao, Yao Lu, Keerthana Gopalakrishnan, Kuang-Huei Lee, Quan Vuong, Paul Wohlhart, Sean Kirmani, Brianna Zitkovich, Fei Xia, et al. Open-world object manipulation using pre-trained vision-language models. *arXiv preprint arXiv:2303.00905*, 2023. 3
- [31] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 23–30, 2017. 3

- [32] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 3, 10
- [33] Junjie Wen, Minjie Zhu, Yichen Zhu, Zhibin Tang, Jinming Li, Zhongyi Zhou, Chengmeng Li, Xiaoyu Liu, Yaxin Peng, Chaomin Shen, et al. Diffusion-vla: Scaling robot foundation models via unified diffusion and autoregression. *arXiv preprint arXiv:2412.03293*, 2024. 2, 3, 4, 7, 10
- [34] Junjie Wen, Yichen Zhu, Jinming Li, Minjie Zhu, Kun Wu, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin Shen, Yaxin Peng, et al. Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation. *arXiv preprint arXiv:2409.12514*, 2024. 2, 3, 4
- [35] Junjie Wen, Yichen Zhu, Jinming Li, Zhibin Tang, Chaomin Shen, and Feifei Feng. Dexvla: Vision-language model with plug-in diffusion expert for general robot control. *arXiv preprint arXiv:2502.05855*, 2025. 3
- [36] Weirui Ye, Fangchen Liu, Zheng Ding, Yang Gao, Oleh Rybkin, and Pieter Abbeel. Video2policy: Scaling up manipulation tasks in simulation through internet videos. *arXiv preprint arXiv:2502.09886*, 2025. 3
- [37] Michał Zawalski, William Chen, Karl Pertsch, Oier Mees, Chelsea Finn, and Sergey Levine. Robotic control via embodied chain-of-thought reasoning. *arXiv preprint arXiv:2407.08693*, 2024. 2, 4
- [38] Michał Zawalski, William Chen, Karl Pertsch, Oier Mees, Chelsea Finn, and Sergey Levine. Robotic control via embodied chain-of-thought reasoning. *arXiv preprint arXiv:2407.08693*, 2024. 3
- [39] Han Zhao, Min Zhang, Wei Zhao, Pengxiang Ding, Siteng Huang, and Donglin Wang. Cobra: Extending mamba to multi-modal large language model for efficient inference. *arXiv preprint arXiv:2403.14520*, 2024. 3
- [40] Zhongyi Zhou, Yichen Zhu, Minjie Zhu, Junjie Wen, Ning Liu, Zhiyuan Xu, Weibin Meng, Ran Cheng, Yaxin Peng, Chaomin Shen, et al. Chatvla: Unified multimodal understanding and robot control with vision-language-action model. *arXiv preprint arXiv:2502.14420*, 2025. 3
- [41] Minjie Zhu, Yichen Zhu, Xin Liu, Ning Liu, Zhiyuan Xu, Chaomin Shen, Yaxin Peng, Zhicai Ou, Feifei Feng, and Jian Tang. Mipha: A comprehensive overhaul of multi-modal assistant with small language models. *arXiv preprint arXiv:2403.06199*, 2024. 3
- [42] Yifeng Zhu, Arisrei Lim, Peter Stone, and Yuke Zhu. Vision-based manipulation from single human video with open-world object graphs. *arXiv preprint arXiv:2405.20321*, 2024. 3
- [43] Yichen Zhu, Zhicai Ou, Xiaofeng Mou, and Jian Tang. Retrieval-augmented embodied agents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17985–17995, 2024. 3
- [44] Yichen Zhu, Minjie Zhu, Ning Liu, Zhiyuan Xu, and Yaxin Peng. Llava-phi: Efficient multi-modal assistant with small language model. In *Proceedings of the 1st International Workshop on Efficient Multimedia Computing under Limited*, pages 18–22, 2024. 3
- [45] Yichen Zhu, Zhicai Ou, Feifei Feng, and Jian Tang. Any2policy: Learning visuomotor policy with any-modality. *Advances in Neural Information Processing Systems*, 37: 133518–133540, 2025. 3

7. Appendix

7.1. Evaluation Metrics

For real robot task, we record the percentage of trials where the robot successfully completes the assigned task. This is a fundamental metric for any robot experiment. Multiple trials are conducted for the evaluation.

7.2. Implementation Details

All experiments were conducted on eight NVIDIA A800 GPUs using the Adam optimizer with a constant learning rate of $2e-5$ and a global batch size of 128. Training proceeded for 50,000 steps, with the final checkpoint selected based on validation performance. Unless otherwise stated, the ratio of robot data to visual-text data was 10:1. Empirically, we observed that increasing the proportion of robot data significantly degraded manipulation performance. Our base model is DiVLA[33] with a Qwen2-VL-2B backbone [32]. As our focus is developing a co-training method for novel object generalization, we retained the original model architecture.

7.3. Example of Objects Used in Experiments

We provide a comprehensive list of out-of-distribution objects and names that we used for training and evaluation, which are shown in Figure 7 and Figure 8.



Figure 7. The out-of-distribution objects used in our experiments (part 1).



Figure 8. The out-of-distribution objects used in our experiments (part 2).