
RoboGrasp: A Universal Grasping Policy for Robust Robotic Control

Yiqi Huang¹ Travis Davies¹ Jiahuan Yan¹ Xiang Chen² Yu Tian³ Luhui Hu¹

Abstract

Imitation learning and world models have shown significant promise in advancing generalizable robotic learning, with robotic grasping remaining a critical challenge for achieving precise manipulation. Existing methods often rely heavily on robot arm state data and RGB images, leading to overfitting to specific object shapes or positions. To address these limitations, we propose RoboGrasp, a universal grasping policy framework that integrates pretrained grasp detection models with robotic learning. By leveraging robust visual guidance from object detection and segmentation tasks, RoboGrasp significantly enhances grasp precision, stability, and generalizability, achieving up to 34% higher success rates in few-shot learning and grasping box prompt tasks. Built on diffusion-based methods, RoboGrasp is adaptable to various robotic learning paradigms, enabling precise and reliable manipulation across diverse and complex scenarios. This framework represents a scalable and versatile solution for tackling real-world challenges in robotic grasping.

1. Introduction

When a baby encounters an object for the first time, it can often grasp it instinctively. For robots, however, this task is far more complex. Policies trained for one object often fail to generalize to others. Recent advances in Behavior Cloning, particularly diffusion-based policies, have emerged as a promising solution, offering flexibility and expressiveness in handling complex, multi-modal action spaces (Pearce et al., 2023; Chi et al., 2023).

However, Behavior Cloning still face challenges in generalizing beyond their training environments, particularly in dynamic, cluttered settings with unseen or distractor ob-

jects. A key limitation lies in their reliance on raw sensor data for conditional input during training and inference (Chi et al., 2023; Ze et al., 2024). Without explicit task guidance, these policies depend on implicit patterns learned from data, limiting their robustness (Selvaraju et al., 2019).

To address this, we propose leveraging advancements in computer vision to enhance perception. Pre-trained vision models for tasks like object detection (Redmon et al., 2016), segmentation (Ravi et al., 2024), pose estimation (Huang et al., 2019), and depth estimation (Yang et al., 2024) can provide structured, task-relevant information. By integrating these models, robotic policies can focus on relevant objects and regions, even in cluttered environments, enabling scalable generalization to novel objects and tasks.

We introduce RoboGrasp, a universal grasping policy framework that integrates an auxiliary grasping-box detection model. This model identifies precise grasp regions, providing explicit spatial guidance for the robot arm. By conditioning the policy on these grasping boxes, RoboGrasp enhances generalizability and adaptability.

Our experiments explore two key questions: (1) Can RoboGrasp leverage grasping-based affordances for effective **few-shot learning** on new or unseen objects? (2) Can it use **grasping-based affordance prompts** as visual cues to define objectives and generate effective policies? These questions aim to evaluate the scalability of grasping-based affordances for robust, generalizable robot manipulation in real-world environments.

This work represents a step toward deploying robots in unstructured settings, reducing reliance on controlled lab data and improving adaptability for diverse, dynamic tasks.

2. Related Works

Recent advancements in robot policy planning have facilitated the democratization of Behavior Cloning (BC), extending its reach beyond specialized research labs (Zhao et al., 2023; Team et al., 2024; Chi et al., 2024). These approaches typically involve models that map sensor observations into trajectories of future robot poses. In this context, diffusion models have emerged as a powerful tool to address critical limitations of Behavior Cloning, such as covariate shift (Pomerleau, 1989), where robots fail to generalize be-

*Equal contribution ¹ZhiCheng AI, Hangzhou, China ²Peking University ³Harvard University. Correspondence to: Yiqi Huang <yiqi.huang.19@outlook.com>, Luhui Hu <luhuihu@gmail.com>.

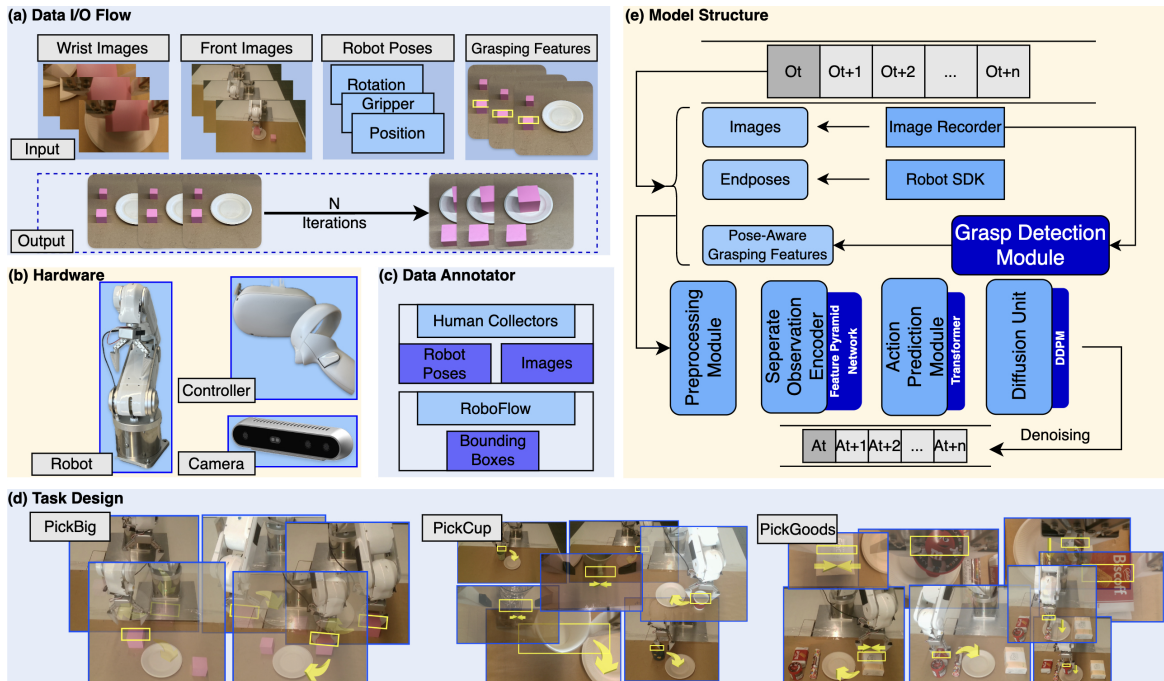


Figure 1. An overview RoboGrasp architecture, demonstrating the integration of grasping guidance, RGB images and robot state data to enhance generalizability and precision of grasping manipulation. (a) Data flow and datasets used for training and inference. (b) Hardware setup, including an industrial-grade robotic arm, RealSense cameras, and a Quest VR headset for data collection. (c) Annotation of demonstrations for grasping affordances. (d) Experimental task designs. (e) The RoboGrasp policy architecture.

yond their training data (Zhou et al., 2022). Diffusion-based policies, exemplified by Diffusion Policy (DP) (Chi et al., 2023), overcome these challenges by generating diverse and multi-modal action trajectories, significantly improving robustness in dynamic and unpredictable environments.

Recent large-scale robotic expert demonstration datasets (Collaboration et al., 2024) have fueled efforts to scale BC architectures. Works like Robotics Diffusion Transformer (RDT) (Liu et al., 2024b), Octo (Octo Model Team et al., 2024), and π_0 (Black et al., 2024) demonstrate that skills learned from diverse datasets can transfer to novel tasks, with some models achieving zero-shot generalization to grasping new objects. However, training large-scale models remains computationally expensive, limiting accessibility for resource-constrained settings.

Recent efforts have investigated point-based affordance representations (Liu et al., 2024a; Tang et al., 2024; Huang et al., 2024), where keypoints are used to identify task-relevant objects and guide the policy with structured information, often leveraging pre-trained vision models. While scalable, these approaches primarily convey object locations but lack actionable information on how to grasp or manipulate them effectively.

Grasping-based affordance representations offer a more comprehensive solution by encoding feasible grasping strate-

gies (Kleeberger et al., 2020), providing both spatial and actionable information. Datasets like Grasp Anything (Vuong et al., 2023) highlight the potential for scalable data collection in this domain. However, integrating grasping affordances with diffusion-based policies remains underexplored. Existing works such as GQCNN (Mahler et al., 2017) provide initial steps, but further research is needed to unlock the full potential of affordance-driven planning.

Our work bridges this gap by integrating grasping-based affordance representations with diffusion-based policies. By providing richer conditional inputs, we aim to improve the efficiency and generalization of robot planning models, particularly in resource-constrained settings.

3. RoboGrasp Policy

This section outlines the architecture of the RoboGrasp, an augmented variation of Diffusion Policy (DP) designed to incorporate grasp-specific information for improved robotic manipulation. Key enhancements include the integration of a Grasp Detection Module and modifications to the observation encoder. Hyperparameters, such as the number of historical timesteps (2) and predicted actions (16), remain consistent with the original DP framework.

The grasping box information includes, as shown in Fig-

ure 2, the x and y coordinate of the grasping box’s central point along with the height and width of the box. Normally the angle of rotation in relation to the camera’s orientation is also included, however since the robot arm used in this experiment cannot rotate, these parameter was considered redundant in experiments, and all objects were left in unrotated positions.

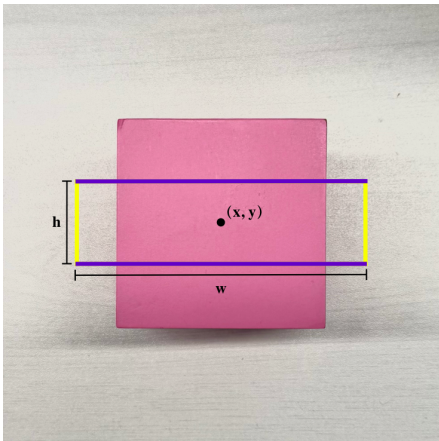


Figure 2. The anatomy of a grasping box. A region on an item indicating the region that can be grasped, along with the x , y coordinates of the box’s centroid and the box’s width and height.

3.1. Grasp Detection Module

The Grasp Detection Module leverages YOLOv11-m (Redmon et al., 2016) for its speed, simplicity, and generalizability. YOLOv11-m was fine-tuned on a custom-labeled dataset to predict the object class, 2D spatial coordinates of the grasping box’s center, and the box’s width and height. During policy training, labels generated by the Grasp Detection Module were directly utilized, while at inference, YOLOv11-m dynamically predicts grasping boxes for the observed data. To simplify grasp selection, the module outputs only the box with the highest confidence score for each run, as the task involves grasping a single object per experiment.

3.2. Observation Encoder

The observation encoder combines visual and low-dimensional data into a unified latent representation. A ResNet34-based feature pyramid encoder is employed for each camera view, processing multiview RGB data separately before concatenation. Low-dimensional inputs, such as the robot arm’s end pose and gripper sensor data, are incorporated following the original DP design. A novel augmentation introduces grasping box features—class label and spatial information—into the concatenated observation data.

This concatenated data is projected into a fixed-dimensional latent space, serving as a single token per timestep. To capture temporal dependencies, an untrained, lightweight transformer applies self-attention across the designated historical timesteps.

3.3. Diffusion Action Head

The action head utilizes a lightweight diffusion transformer, identical to that in DP, to predict actions over 16 timesteps. A DDIM scheduler (Nichol & Dhariwal, 2021) with a Cosine Beta noise schedule is used for denoising, ensuring efficient and smooth sampling.

Cross-attention mechanisms condition the noised actions on observation tokens, enabling the policy to integrate visual and spatial context effectively. Actions are linearly projected into the latent space for processing within the transformer and are subsequently reprojected into their original dimensions via dedicated linear layers.

4. Experiments

Addressing the need for grasp-focused tasks is essential to overcome the limitations inherent in traditional robot learning experiments. Commonly, these experiments utilize similar object targets, allowing models to extensively learn from these specific objects and their associated task completion methods. However, real-world scenarios frequently present a diverse array of objects with varying sizes, types, and grasping requirements, challenging the model’s ability to generalize effectively. To bridge this gap, we designed three primary tasks, **PickBig**, **PickCup**, and **PickGoods** to evaluate the robot’s capability to perform accurate grasping across different object sizes, varied object types with distinct grasping strategies, few-shot learning abilities, and promptable grasping actions.

4.1. Task Description

PickBig: This task evaluates the robot’s ability to distinguish and grasp the larger of two nearly identical blocks, differing only in size. The blocks are placed in eight distinct positions within the workspace, introducing variability in both object dimensions and spatial arrangements. PickBig (see Figure 4) assesses the model’s capacity to adapt its grasping strategies to accommodate size differences and spatial diversity, ensuring accurate and stable grasps across various scenarios. A key challenge lies in defining the task’s goal, and the task aims to test whether providing grasping-based affordance regions helps clarify and achieve the objective more effectively. This focus on goal definition and adaptability makes PickBig a robust test of the model’s precision and responsiveness in goal-oriented grasping tasks

PickCup: This task focuses on the robot’s proficiency in

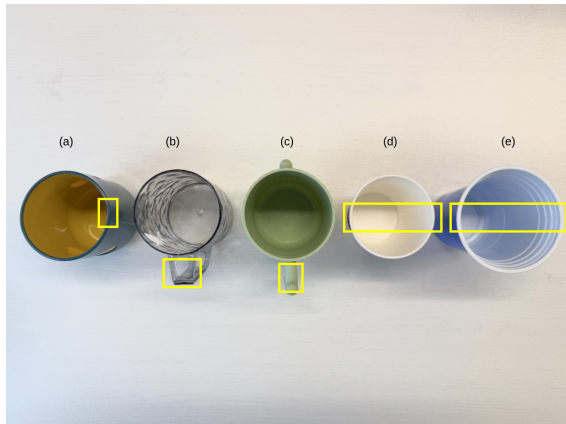


Figure 3. Grasping boxes for cups used in the experiments, shown in a bird’s-eye view. (a) illustrates a grasp by the wall of the cup, while (b) and (c) demonstrate grasps by the cup handles. (d) and (e) depict grasps over the cup’s diameter. (c) and (e) represent the cups used in the few-shot experiments.

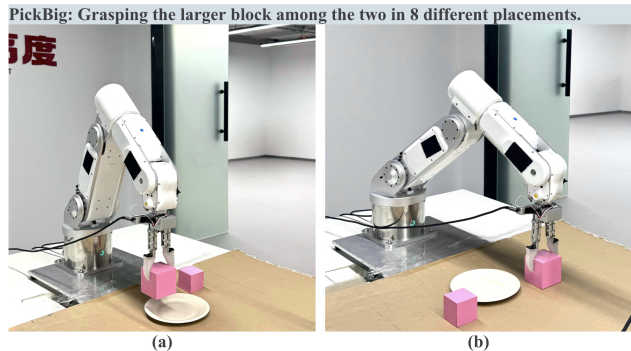


Figure 4. Placement Positions Generalizability experiment setup for PickBig. (a) and (b) show two of the eight placement positions. The objective of PickBig is to distinguish between two similarly shaped blocks and successfully grasp the larger one along its diameter.

grasping different types of cups using various grasping strategies, as shown in Figure 3. Three types of cups are employed, each presented with a distinct grasping pattern. Additionally, cups are placed in four distinct positions to introduce variability in grasping scenarios. To evaluate the model’s few-shot learning ability, additional instances are introduced with a limited number of demonstration trials (see Figure 5). This inclusion assesses the model’s capacity to generalize grasping strategies to new or less-represented objects with minimal training data, ensuring robustness and adaptability in handling diverse and unfamiliar cup types.

PickGoods: This task evaluates the robot’s ability to generalize grasping strategies across a wide variety of retail goods, simulating real-world retail environments (see Figure 6). A diverse set of retail items, varying in shape, size, and material, are selected to challenge the robot’s adaptabil-

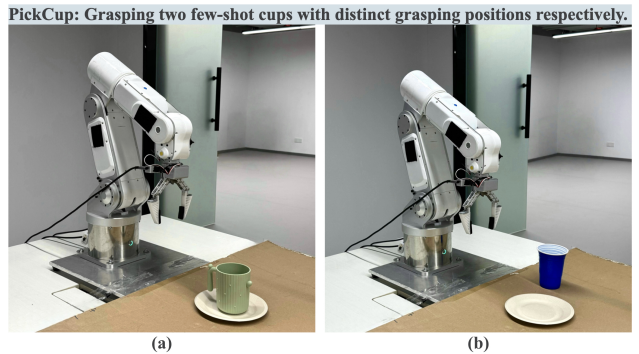


Figure 5. Few-shot experiment setup for PickCup task. The green mug in (a) represents the handle grasping few-shot task with only 5 demonstrations. The blue plastic cup in (b) represents the diameter grasping few-shot task with 10 demonstrations.

ity. Each item is grasped using a single, consistent grasping pattern, ensuring uniformity in the approach.

The PickGoods task specifically tests RoboGrasp’s ability to generate the correct grasping policy based on a provided *grasping-based affordance region*. This region serves as a *spatial prompt*, guiding the policy toward the desired goal. The prompt acts as a critical test of the policy’s responsiveness to predefined objectives, a feature that is notably absent in approaches like DP. Unlike DP, which relies solely on conditional sensor data without explicit goal specification, PickGoods incorporates a clear, goal-oriented prompt, enabling the policy to align its actions with the intended outcome.

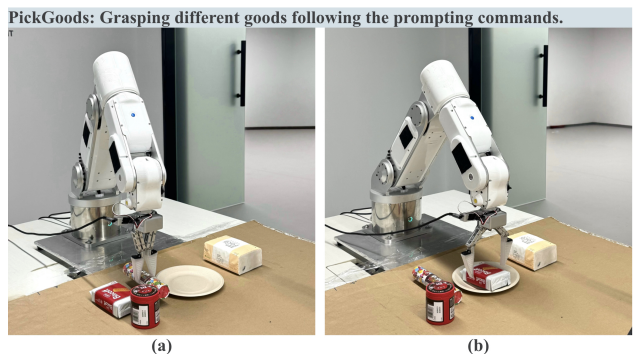


Figure 6. Promptable Grasping experiment setup for the PickGoods task. In (a), the grasping box for the chocolate bar is provided, while in (b), the grasping box for the biscuit is provided. The objective is to follow the grasping box prompts and successfully pick the specified item.

4.2. Data Processing

The datasets for each task are meticulously curated to reflect a wide range of real-world grasping scenarios:

- **PickBig**: Comprises 600 demonstration trials (8 placement positions \times 75 demo each) involving two blocks, positioned in eight distinct configurations.
- **PickCup**: Consists of 315 demonstration trials (3 cup types \times 4 placement positions \times 25 demo + 15 few-shot demo) involving three types of cups, mug, plastic cup, and paper cup, each subjected to three grasping patterns: handle grasping, cup wall grasping, and diameter grasping. Additionally, 15 extra demonstrations for mugs and paper cups is included to evaluate the model’s few-shot learning capabilities.
- **PickGoods**: Contains 400 demonstration trials where each of the retail good has 100 demonstrations.

The data preprocessing pipeline ensures high-quality and consistent grasping annotations for each task. A representative subset of approximately 500 frames is manually labeled with grasping boxes, covering all scenarios within the task. These annotations serve as the foundation for training the Grasp Detection Module. Fine-tuning the module on this dataset achieves a mean Average Precision (mAP) exceeding 98%.

Once trained, the Grasp Detection Module automatically generates grasping boxes for every frame in the collected video data, significantly enhancing annotation scalability while reducing human error. This uniform and reliable annotations support the training of our RoboGrasp policies. During inference, the Grasp Detection Module is seamlessly integrated to predict grasping boxes in real-time (see Figure 7). This enables the system to dynamically identify optimal grasping regions, ensuring precise and stable grasps across diverse object types and scenarios.

4.3. Experimental Design and Evaluations

Although we only use Diffusion Policy (DP) as a baseline, our solution can adapt to various robotic learning frameworks. To analyze how grasping-based affordances enhance model performance, we propose an ablation study across three primary tasks, **PickBig**, **PickCup**, and **PickGoods** to compare two model configurations:

1. **DP Model (Baseline)**: The standard diffusion-based visuomotor policy without any additional enhancements. It serves as a benchmark for measuring the impact of subsequent modifications.
2. **RoboGrasp Model**: Integrates grasping box annotations generated by the Grasp Detection Module to guide the diffusion process. This focuses the model on optimal grasping regions, improving accuracy and stability.

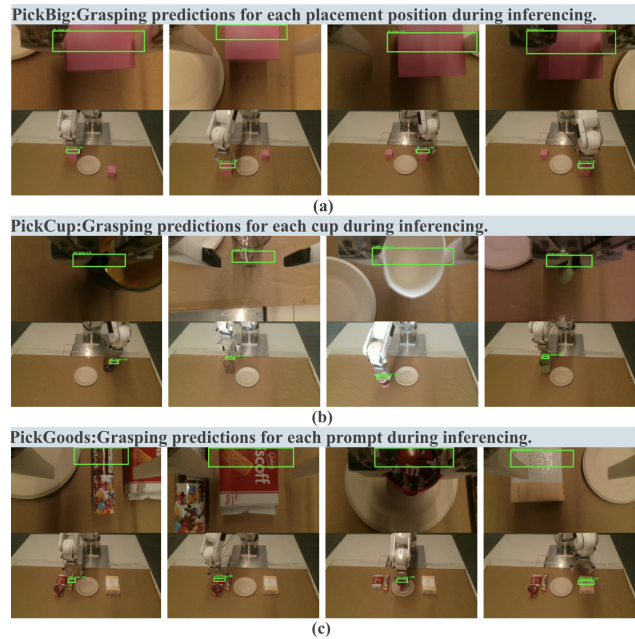


Figure 7. Real-time predictions from our pretrained grasp detection module across different tasks. (a) Demonstrates robust grasping predictions across various placement positions. (b) Highlights accurate detection for diverse grasping strategies. (c) Showcases effective prompting-based predictions.

4.3.1. TESTING STANDARDIZATION

To ensure a fair and reliable comparison between DP and RoboGrasp, strict standardization was maintained during inference testing. For each task, the placement positions and involved objects were identical across both models. This controlled setup eliminates variability in testing conditions, ensuring that any performance differences observed are attributable to the models themselves rather than inconsistencies in the experimental setup. Such ablation-based evaluations provide a robust framework for identifying and quantifying improvements in model performance.

4.3.2. EVALUATION METRICS

To quantitatively compare the two models, we employed the following metrics:

1. **Task Success Rate (TSR)**: The percentage of successful task completions across all tasks.
2. **Grasp Success Rate (GSR)**: Evaluates the effectiveness of grasping strategies by measuring action accuracy, consistency, and stability across diverse objects and scenarios. This metric is particularly useful for assessing the model’s ability to adapt and apply appropriate grasping strategies in varying contexts. GSR is defined as:

Table 1. Detailed Performance Comparison Across Tasks for Diffusion Policy and RoboGrasp. The table highlights task success rate and grasping success rate for each target object.

TASK	OBJ. NUM	TARGET	DEMOS	POSITION. NUM	GRASPING STRATEGY	MODEL	TSR(%)	GSR (%)
PICKBIG	2	BIGGER BLOCK	600	8	DIAMETER	DP	67.5	66.25
						ROBOGRASP	97.5	96.25
PICKCUP	5	GREY MUG	100	4	HANDLE	DP	95	80
						ROBOGRASP	100	100
		BLUE PLASTIC CUP	100	4	WALL	DP	87.5	85
						ROBOGRASP	92.5	92.5
		RED PAPER CUP	100	4	DIAMETER	DP	95	70
				ROBOGRASP	100	100		
		BLUE PLASTIC CUP	10	2	DIAMETER	DP	60	55
					ROBOGRASP	100	95	
		GREEN MUG	5	1	HANDLE	DP	70	60
					ROBOGRASP	100	100	
PICKGOODS	4	MEIJI CHOCOLATE BAR	100	1	DIAMETER	DP	0	0
						ROBOGRASP	100	98
		LOTUS BISCUIT	100	1	DIAMETER	DP	0	0
						ROBOGRASP	0	0
		M&M	100	1	DIAMETER	DP	0	0
					ROBOGRASP	4	4	
		TISSUE	100	1	DIAMETER	DP	89.5	86
					ROBOGRASP	100	100	

$$GSR = \frac{\text{No. Successful Grasps}}{\text{No. Total Grasp Attempts}} \quad (1)$$

These metrics provide a comprehensive assessment of each model’s performance, capturing both high-level task success and fine-grained grasping proficiency.

5. Results and Discussion

Our experimental results clearly demonstrate RoboGrasp’s superiority over DP across all evaluated tasks. Table 1 provides a detailed comparison of TSR and GSR, offering insights into how RoboGrasp consistently outperforms DP across various objects, grasping strategies, and placement positions.

In the PickBig task, RoboGrasp achieves a task success rate of 97.5% and a grasp success rate of 96.25%, significantly outperforming DP’s 67.5% and 66.25%. The effectiveness of grasping box detections is evident as RoboGrasp excels in distinguishing the larger block across eight varied placement positions, adapting seamlessly to positional changes and size differences.

In the PickCup task, involving five distinct objects and diverse grasping strategies, RoboGrasp consistently achieves

near-perfect performance, with task success rates of 100% and grasp success rates ranging from 92.5% to 100%. This significantly surpasses DP, which struggles with inconsistent strategies, particularly for challenging cases like blue plastic cups and green mugs. RoboGrasp’s ability to transfer learned skills to few-shot objects highlights its adaptability.

For the PickGoods task, RoboGrasp’s prompt-based grasping successfully handles two out of four target objects, including chocolate bars and tissue packs. While DP struggles with object identification and inconsistent grasps, RoboGrasp utilizes grasping box prompts to focus on the target, achieving up to 100% task success and 98% grasp success rates for these objects, demonstrating adaptability despite the task’s challenges.

These results highlight the critical role of grasping box detections in improving task performance and generalization capabilities. RoboGrasp consistently delivers higher task success and grasping success rates across diverse objects, placements, and strategies, demonstrating its adaptability and robustness in various manipulation challenges.

5.1. Data Compensation

The size and diversity of the dataset play a critical role in the training and performance of robotic learning models

like DP and RoboGrasp. We identified several key factors influencing training and generalization:

State Space vs. Number of Demonstrations: In the PickBig task, the initial dataset comprised only 300 demonstrations. However, training results revealed a high variance in training loss and a large mean squared error (MSE) in action predictions for both DP and RoboGrasp. This inconsistency is likely due to the large state space introduced by eight distinct placement positions, which demands more data to adequately capture the variations. To address this, the dataset size was doubled to 600 demonstrations (75 per placement position), resolving the variance issue and improving training stability.

Similar vs. Distinct Objects: In contrast, the PickCups task demonstrated a different data requirement. For this task, only 25 demonstrations were collected for each cup at each placement position, and the models still achieved convergence during training. This outcome is likely because the cups in this task are distinct, making it easier for the model to differentiate between objects.

These results suggest that datasets need to be scaled up proportionally to the size of the state space and when the target objects are similar in appearance. The distinction between objects significantly impacts the model’s ability to generalize and perform effectively with fewer demonstrations.

5.2. Task Performance Analysis

This section systematically analyzes RoboGrasp’s capabilities and limitations across key dimensions critical to real-world deployment: (1) the interplay between state space complexity and dataset scale in policy optimization, (2) generalization capacity through few-shot skill transfer to novel objects, and (3) responsiveness to spatially grounded affordance cues for deriving context-aware policies. By dissecting performance variations across tasks, we identify how environmental constraints, object distribution sparsity, and affordance grounding collectively shape the system’s adaptability and failure modes.

5.2.1. ADDRESSING POSITIONING AND STATE SPACES

The PickBig task, with its eight distinct placement positions, requires a robust policy to handle a large state space. Increasing demonstrations from 300 to 600 resolved training instability, enabling better state coverage and a fairer comparison between DP and RoboGrasp.

DP struggles with overfitting to low-dimensional robot state representations, leading to repetitive, fixed movements and inconsistent performance. It often fails to adapt to positional changes or differentiate between similar-shaped blocks, limiting its effectiveness.

In contrast, RoboGrasp achieves a 33.75% higher TSR by leveraging grasping box detections as additional input (see Figure 8). These predictions explicitly guide RoboGrasp to target the correct block, enabling dynamic adaptation to positional changes and precise block differentiation. This results in significantly improved grasping accuracy and consistency.

The PickBig results highlight the importance of state space diversity in training and the critical role of grasping box detections in guiding robotic policies. RoboGrasp’s superior performance demonstrates its potential for generalizable and robust robotic manipulation

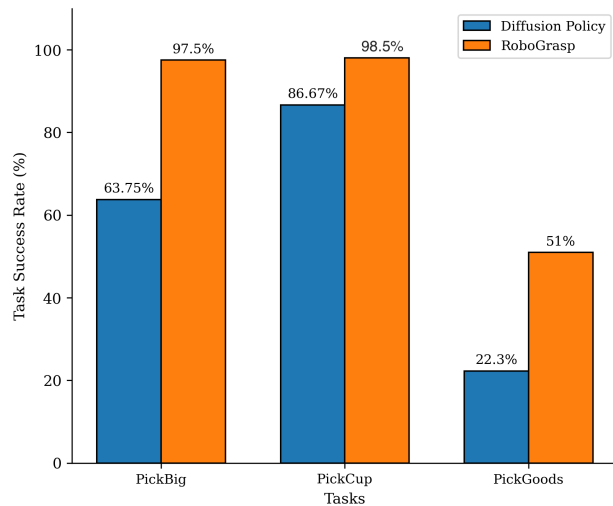


Figure 8. Comparison of average Task Success Rate for DP and RoboGrasp.

5.2.2. FEW-SHOT LEARNING AND STRATEGY SELECTION

The PickCups task, with its four placement positions and distinct cup shapes, presents a smaller state space than PickBig, allowing DP to perform better by relying on a narrower diversity of policy action trajectories. However, RoboGrasp outperforms DP with a 23.33% higher GSR (see Figure 9). While DP can complete the pick-and-place task, it often employs inconsistent and mixed grasping strategies, such as using a wall grasp for one cup and a rim grasp for another, failing to generalize effectively. In contrast, RoboGrasp excels in few-shot learning, consistently selecting the correct grasping strategy for cups with only 5 or 10 demonstrations, even when their shapes differ from the primary training set. By leveraging grasping box predictions, RoboGrasp ensures precise and consistent manipulation, demonstrating its ability to generalize across geometric variations with minimal data and adapt to diverse tasks reliably

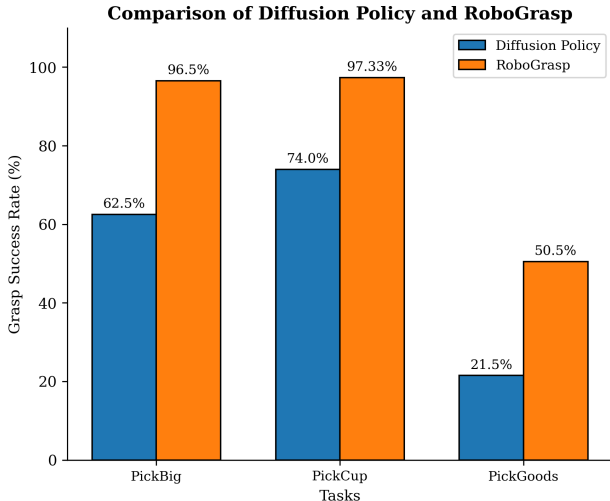


Figure 9. Comparison of average Grasp Success Rate for DP and RoboGrasp.

5.2.3. GRASPING-BASED AFFORDANCE PROMPT

The superior performance of RoboGrasp over DP in the PickBig task highlights the efficacy of grasping-based affordance prompts in guiding policy decisions for goal-oriented manipulation. In contrast, while RoboGrasp still significantly outperforms DP in the PickGoods task (28.7% improvement in TSR), its relative performance decline in this scenario stems from the inherent complexity of the environment: three closely adjacent candidate objects in a confined workspace created a multimodal action probability distribution. This complexity led the policy to prioritize proximal targets (e.g., grasping the first visible item) rather than strategically selecting optimal objects. Notably, the grasping-based affordances did still influence directional choices, such as prompting leftward motions for tissue retrieval or rightward motions for food items. The simpler PickBig environment—with only two candidate objects—resulted in a unimodal action distribution, enabling more deterministic and effective policy execution. However, the fixed spatial placement of objects in PickGoods raises concerns about potential over-reliance on robot pose priors rather than affordance-driven reasoning. To address these limitations, future work could enhance policy robustness by training on datasets with greater positional diversity, thereby reducing environmental bias and improving generalization to cluttered configurations.

6. Conclusion

RoboGrasp presents a novel approach to robotic grasping, addressing key challenges in precision and generalization by leveraging pretrained grasp detection models. By moving beyond the limitations of traditional reliance on robot

state data and RGB images, RoboGrasp enables robust and adaptable grasping across diverse scenarios, demonstrating significant improvements in grasp success rates.

Our experiments evaluated RoboGrasp’s ability to transfer grasping skills to new objects through few-shot learning and its capacity to utilize grasping boxes as visual prompts to define goals and generate effective robot policies. The results highlight RoboGrasp’s strong few-shot learning capabilities, achieving superior generalization to unseen items compared to DP. Furthermore, RoboGrasp exhibited a substantial performance boost in defining task objectives and generating policies with grasping box prompts, demonstrating its ability to outperform DP by a wide margin.

Built on diffusion-based methods, RoboGrasp is a flexible framework with the potential to scale across a variety of complex manipulation tasks. This work lays a strong foundation for the development of scalable, reliable, and generalizable robotic systems capable of addressing real-world challenges in dynamic and unstructured environments.

7. Future Directions

This work highlights several unexplored avenues with significant potential to advance robotic learning and grasping. Language prompting, inspired by methods such as Grounding DINO (Liu et al., 2024c) and DINO-X (Ren et al., 2024), remains underexplored in robotics. Integrating language commands to generate grasping boxes or guide manipulation tasks could greatly enhance flexibility and generalizability.

Similarly, employing grasp-guided approaches in other frameworks, such as ACT (Zhao et al., 2023), or large foundation models like Robotics Diffusion Transformer (RDT) (Liu et al., 2024b), presents an opportunity to scale robotic learning to broader tasks. Furthermore, incorporating grasping prompts into world models, akin to visual prompts in (Geng et al., 2024), could enhance real-world modeling and planning, boosting performance in dynamic and unstructured environments.

Additionally, evaluating these methods in simulation environments would provide a cost-effective way to assess robustness and scalability. Future work could also explore extending grasping capabilities by incorporating rotational parameters into grasping boxes and employing rotatable arms, enabling robots to handle more complex and precise manipulation tasks. These directions offer immense potential to improve precision, generalization, and adaptability in robotics.

References

- Black, K., Brown, N., Driess, D., Esmail, A., Equi, M., Finn, C., Fusai, N., Groom, L., Hausman, K., Ichter, B., Jakubczak, S., Jones, T., Ke, L., Levine, S., Li-Bell, A., Mothukuri, M., Nair, S., Pertsch, K., Shi, L. X., Tanner, J., Vuong, Q., Walling, A., Wang, H., and Zhilinsky, U. π_0 : A vision-language-action flow model for general robot control, 2024. URL <https://arxiv.org/abs/2410.24164>.
- Chi, C., Feng, S., Du, Y., Xu, Z., Cousineau, E., Burchfiel, B., and Song, S. Diffusion policy: Visuomotor policy learning via action diffusion. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- Chi, C., Xu, Z., Pan, C., Cousineau, E., Burchfiel, B., Feng, S., Tedrake, R., and Song, S. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots, 2024. URL <https://arxiv.org/abs/2402.10329>.
- Collaboration, E., O’Neill, A., Rehman, A., Gupta, A., Madhukuri, A., Gupta, A., Padalkar, A., Lee, A., Pooley, A., Gupta, A., Mandelkar, A., Jain, A., Tung, A., Bewley, A., Herzog, A., Irpan, A., Khazatsky, A., Rai, A., Gupta, A., Wang, A., Kolobov, A., Singh, A., Garg, A., Kembhavi, A., Xie, A., Brohan, A., Raffin, A., Sharma, A., Yavary, A., Jain, A., Balakrishna, A., Wahid, A., Burgess-Limerick, B., Kim, B., Schölkopf, B., Wulfe, B., Ichter, B., Lu, C., Xu, C., Le, C., Finn, C., Wang, C., Xu, C., Chi, C., Huang, C., Chan, C., Agia, C., Pan, C., Fu, C., Devin, C., Xu, D., Morton, D., Driess, D., Chen, D., Pathak, D., Shah, D., Büchler, D., Jayaraman, D., Kalashnikov, D., Sadigh, D., Johns, E., Foster, E., Liu, F., Ceola, F., Xia, F., Zhao, F., Frujeri, F. V., Stulp, F., Zhou, G., Sukhatme, G. S., Salhotra, G., Yan, G., Feng, G., Schiavi, G., Berseth, G., Kahn, G., Yang, G., Wang, G., Su, H., Fang, H.-S., Shi, H., Bao, H., Amor, H. B., Christensen, H. I., Furuta, H., Bharadhwaj, H., Walke, H., Fang, H., Ha, H., Mordatch, I., Radosavovic, I., Leal, I., Liang, J., Abou-Chakra, J., Kim, J., Drake, J., Peters, J., Schneider, J., Hsu, J., Vakil, J., Bohg, J., Bingham, J., Wu, J., Gao, J., Hu, J., Wu, J., Wu, J., Sun, J., Luo, J., Gu, J., Tan, J., Oh, J., Wu, J., Lu, J., Yang, J., Malik, J., Silvério, J., Hejna, J., Booher, J., Tompson, J., Yang, J., Salvador, J., Lim, J. J., Han, J., Wang, K., Rao, K., Pertsch, K., Hausman, K., Go, K., Gopalakrishnan, K., Goldberg, K., Byrne, K., Oslund, K., Kawaharazuka, K., Black, K., Lin, K., Zhang, K., Ehsani, K., Lekkala, K., Ellis, K., Rana, K., Srinivasan, K., Fang, K., Singh, K. P., Zeng, K.-H., Hatch, K., Hsu, K., Itti, L., Chen, L. Y., Pinto, L., Fei-Fei, L., Tan, L., Fan, L. J., Ott, L., Lee, L., Weihs, L., Chen, M., Lepert, M., Memmel, M., Tomizuka, M., Itkina, M., Castro, M. G., Spero, M., Du, M., Ahn, M., Yip, M. C., Zhang, M., Ding, M., Heo, M., Srirama, M. K., Sharma, M., Kim, M. J., Kanazawa, N., Hansen, N., Heess, N., Joshi, N. J., Suenderhauf, N., Liu, N., Palo, N. D., Shafiqullah, N. M. M., Mees, O., Kroemer, O., Bastani, O., Sanketi, P. R., Miller, P. T., Yin, P., Wohlhart, P., Xu, P., Fagan, P. D., Mitrano, P., Sermanet, P., Abbeel, P., Sundaresan, P., Chen, Q., Vuong, Q., Rafailov, R., Tian, R., Doshi, R., Mart’in-Mart’in, R., Baijal, R., Scalise, R., Hendrix, R., Lin, R., Qian, R., Zhang, R., Mendonca, R., Shah, R., Hoque, R., Julian, R., Bustamante, S., Kirmani, S., Levine, S., Lin, S., Moore, S., Bahl, S., Dass, S., Sonawani, S., Tulsiani, S., Song, S., Xu, S., Haldar, S., Karamcheti, S., Adebola, S., Guist, S., Nasiriany, S., Schaal, S., Welker, S., Tian, S., Ramamoorthy, S., Dasari, S., Belkhal, S., Park, S., Nair, S., Mirchandani, S., Osa, T., Gupta, T., Harada, T., Matsushima, T., Xiao, T., Kollar, T., Yu, T., Ding, T., Davchev, T., Zhao, T. Z., Armstrong, T., Darrell, T., Chung, T., Jain, V., Kumar, V., Vanhoucke, V., Zhan, W., Zhou, W., Burgard, W., Chen, X., Chen, X., Wang, X., Zhu, X., Geng, X., Liu, X., Liangwei, X., Li, X., Pang, Y., Lu, Y., Ma, Y. J., Kim, Y., Chebotar, Y., Zhou, Y., Zhu, Y., Wu, Y., Xu, Y., Wang, Y., Bisk, Y., Dou, Y., Cho, Y., Lee, Y., Cui, Y., Cao, Y., Wu, Y.-H., Tang, Y., Zhu, Y., Zhang, Y., Jiang, Y., Li, Y., Li, Y., Iwasawa, Y., Matsuo, Y., Ma, Z., Xu, Z., Cui, Z. J., Zhang, Z., Fu, Z., and Lin, Z. Open x-embodiment: Robotic learning datasets and rt-x models, 2024. URL <https://arxiv.org/abs/2310.08864>.
- Geng, D., Herrmann, C., Hur, J., Cole, F., Zhang, S., Pfaff, T., Lopez-Guevara, T., Doersch, C., Aytar, Y., Rubinstein, M., Sun, C., Wang, O., Owens, A., and Sun, D. Motion prompting: Controlling video generation with motion trajectories, 2024. URL <https://arxiv.org/abs/2412.02700>.
- Huang, J., Zhu, Z., and Huang, G. Multi-stage hrnet: Multiple stage high-resolution network for human pose estimation, 2019. URL <https://arxiv.org/abs/1910.05901>.
- Huang, W., Wang, C., Li, Y., Zhang, R., and Fei-Fei, L. Rekep: Spatio-temporal reasoning of relational keypoint constraints for robotic manipulation. *arXiv preprint arXiv:2409.01652*, 2024.
- Kleeberger, K., Bormann, R., Kraus, W., and Huber, M. F. A survey on Learning-Based robotic grasping. *Current Robotics Reports*, 1(4):239–249, December 2020.
- Liu, F., Fang, K., Abbeel, P., and Levine, S. Moka: Open-world robotic manipulation through mark-based visual prompting, 2024a. URL <https://arxiv.org/abs/2403.03174>.
- Liu, S., Wu, L., Li, B., Tan, H., Chen, H., Wang, Z., Xu, K., Su, H., and Zhu, J. Rdt-1b: a diffusion foundation

- model for bimanual manipulation, 2024b. URL <https://arxiv.org/abs/2410.07864>.
- Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Jiang, Q., Li, C., Yang, J., Su, H., Zhu, J., and Zhang, L. Grounding dino: Marrying dino with grounded pre-training for open-set object detection, 2024c. URL <https://arxiv.org/abs/2303.05499>.
- Mahler, J., Liang, J., Niyaz, S., Laskey, M., Doan, R., Liu, X., Ojea, J. A., and Goldberg, K. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics, 2017. URL <https://arxiv.org/abs/1703.09312>.
- Nichol, A. and Dhariwal, P. Improved denoising diffusion probabilistic models, 2021. URL <https://arxiv.org/abs/2102.09672>.
- Octo Model Team, Ghosh, D., Walke, H., Pertsch, K., Black, K., Mees, O., Dasari, S., Hejna, J., Xu, C., Luo, J., Kreiman, T., Tan, Y., Chen, L. Y., Sanketi, P., Vuong, Q., Xiao, T., Sadigh, D., Finn, C., and Levine, S. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.
- Pearce, T., Rashid, T., Kanervisto, A., Bignell, D., Sun, M., Georgescu, R., Macua, S. V., Tan, S. Z., Momennejad, I., Hofmann, K., and Devlin, S. Imitating human behaviour with diffusion models, 2023. URL <https://arxiv.org/abs/2301.10677>.
- Pomerleau, D. Alvin: An autonomous land vehicle in a neural network. In Touretzky, D. (ed.), *Proceedings of (NeurIPS) Neural Information Processing Systems*, pp. 305 – 313. Morgan Kaufmann, December 1989.
- Ravi, N., Gabeur, V., Hu, Y.-T., Hu, R., Ryali, C., Ma, T., Khedr, H., Rädle, R., Rolland, C., Gustafson, L., Mintun, E., Pan, J., Alwala, K. V., Carion, N., Wu, C.-Y., Girshick, R., Dollár, P., and Feichtenhofer, C. Sam 2: Segment anything in images and videos, 2024. URL <https://arxiv.org/abs/2408.00714>.
- Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. You only look once: Unified, real-time object detection, 2016. URL <https://arxiv.org/abs/1506.02640>.
- Ren, T., Chen, Y., Jiang, Q., Zeng, Z., Xiong, Y., Liu, W., Ma, Z., Shen, J., Gao, Y., Jiang, X., Chen, X., Song, Z., Zhang, Y., Huang, H., Gao, H., Liu, S., Zhang, H., Li, F., Yu, K., and Zhang, L. Dino-x: A unified vision model for open-world object detection and understanding, 2024. URL <https://arxiv.org/abs/2411.14347>.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2):336–359, October 2019. ISSN 1573-1405. doi: 10.1007/s11263-019-01228-7. URL <http://dx.doi.org/10.1007/s11263-019-01228-7>.
- Tang, G., Rajkumar, S., Zhou, Y., Walke, H. R., Levine, S., and Fang, K. Kalie: Fine-tuning vision-language models for open-world manipulation without robot data, 2024. URL <https://arxiv.org/abs/2409.14066>.
- Team, A. ., Aldaco, J., Armstrong, T., Baruch, R., Bingham, J., Chan, S., Draper, K., Dwibedi, D., Finn, C., Florence, P., Goodrich, S., Gramlich, W., Hage, T., Herzog, A., Hoech, J., Nguyen, T., Storz, I., Tabanpour, B., Takayama, L., Tompson, J., Wahid, A., Wahrburg, T., Xu, S., Yaroshenko, S., Zakka, K., and Zhao, T. Z. Aloha 2: An enhanced low-cost hardware for bimanual teleoperation, 2024. URL <https://arxiv.org/abs/2405.02292>.
- Vuong, A. D., Vu, M. N., Le, H., Huang, B., Huynh, B., Vo, T., Kugi, A., and Nguyen, A. Grasp-anything: Large-scale grasp dataset from foundation models, 2023. URL <https://arxiv.org/abs/2309.09818>.
- Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., and Zhao, H. Depth anything: Unleashing the power of large-scale unlabeled data, 2024. URL <https://arxiv.org/abs/2401.10891>.
- Ze, Y., Zhang, G., Zhang, K., Hu, C., Wang, M., and Xu, H. 3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- Zhao, T. Z., Kumar, V., Levine, S., and Finn, C. Learning fine-grained bimanual manipulation with low-cost hardware, 2023. URL <https://arxiv.org/abs/2304.13705>.
- Zhou, K., Liu, Z., Qiao, Y., Xiang, T., and Loy, C. C. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–20, 2022. ISSN 1939-3539. doi: 10.1109/tpami.2022.3195549. URL <http://dx.doi.org/10.1109/TPAMI.2022.3195549>.