
Spatial RoboGrasp: Generalized Robotic Grasping Control Policy

Yiqi Huang*
ZhiCheng AI

Jiankai Sun
Stanford University
jksun@stanford.edu

Travis Davies*
ZhiCheng AI

Xiang Chen
Peking University

Jiahuan Yan
ZhiCheng AI

Luhui Hu
ZhiCheng AI
luhui@zhicheng-ai.com

Abstract

Achieving generalizable and precise robotic manipulation across diverse environments remains a critical challenge, largely due to limitations in spatial perception. While prior imitation-learning approaches have made progress, their reliance on raw RGB inputs and handcrafted features often leads to overfitting and poor 3D reasoning under varied lighting, occlusion, and object conditions. In this paper, we propose a unified framework that couples robust multimodal perception with reliable grasp prediction. Our architecture fuses domain-randomized augmentation, monocular depth estimation, and a depth-aware 6-DoF Grasp Prompt into a single spatial representation for downstream action planning. Conditioned on this encoding and a high-level task prompt, our diffusion-based policy yields precise action sequences, achieving up to 40% improvement in grasp success and 45% higher task success rates under environmental variation. These results demonstrate that spatially grounded perception, paired with diffusion-based imitation learning, offers a scalable and robust solution for general-purpose robotic grasping.

1 Introduction

Humans intuitively adapt their grasping strategies to novel objects and conditions. In contrast, enabling robots to replicate this innate ability, generalizing across diverse objects and dynamic environments, remains a crucial challenge. Recent advances in imitation learning, particularly diffusion-based policies Chi et al. [2023, 2024a], Fu et al. [2024], have shown promise in modeling the complexity and multimodality of manipulation tasks. However, these models often rely heavily on raw RGB inputs and handcrafted features, making them vulnerable to overfitting and weak spatial reasoning—especially when confronted with novel object geometries, changing viewpoints, lighting variations, or occlusions Selvaraju et al. [2019], Li et al. [2023]. Lacking explicit task grounding, such policies depend solely on data-driven patterns, which further limits their robustness and generalizability Selvaraju et al. [2019].

A major contributor to these limitations is the scarcity of perceptual diversity in existing datasets, which are often collected under controlled laboratory conditions Lin et al. [2024]. As a result, models trained in such environments tend to perform poorly when deployed in complex real-world settings. Meanwhile, autonomous driving research has demonstrated that spatial-temporal augmentations and multimodal perception strategies can significantly improve resilience to environmental variability Lan and Hao [2023], Raisuddin et al. [2023], Jin et al. [2024]. We propose to bring similar robustness principles into robotic manipulation through enhanced perception and spatial guidance.

*Authors contributed equally.

To this end, we introduce Spatial RoboGrasp, a unified framework designed to improve both generalization and precision by integrating domain-randomized image augmentation, monocular depth estimation, and structured 6-DoF Grasp Prompts. This multimodal perception module produces explicit spatial inputs that guide a diffusion-based policy toward accurate, contact-aware action generation—without the overhead of full point cloud processing. Conditioned on these enriched inputs and a high-level task prompt, our approach generates robust, goal-directed trajectories. We investigate the following research questions: (1) Can multimodal perception and spatial augmentations improve robustness under significant environmental variation? (2) Can explicit grasp affordances enhance few-shot generalization to novel objects? (3) How effectively can grasp-based spatial prompts guide policy behavior and improve task performance?

By bridging the gap between lab-controlled training and real-world deployment, this work builds on our previous efforts in grasp-guided imitation learning Huang et al. [2025] and robust multimodal visual perception Davies et al. [2024], integrating them into a unified spatial policy architecture for robotic manipulation in unstructured environments.

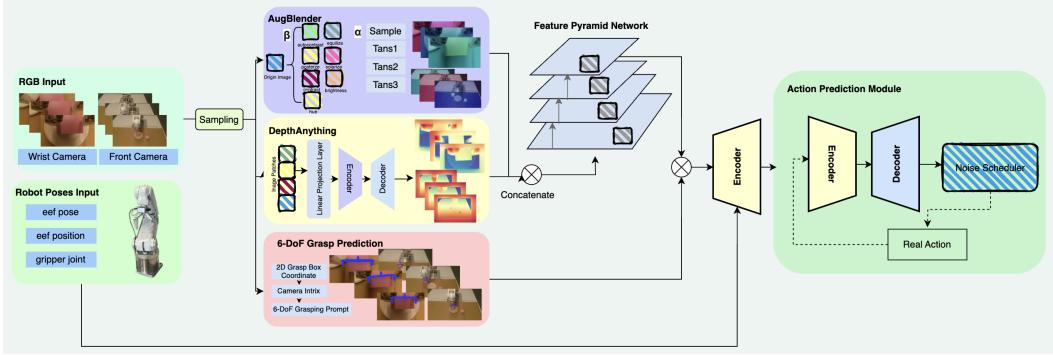


Figure 1: An overview Spatial RoboGrasp architecture, demonstrating the integration of image augmentations, depth estimation and a grasping prompt prediction modules. These observation conditions and robot state data to enhance generalizability and precision of grasping manipulation.

2 Related Works

Recent advances in robot policy planning have broadened the scope of imitation learning (IL) beyond controlled lab environments Fu et al. [2024], Team et al. [2024], Chi et al. [2024b]. IL frameworks typically map sensory observations directly to action sequences. Among them, diffusion-based policies such as Diffusion Policy (DP) Chi et al. [2023] have shown strong performance in tackling covariate shift—where a robot’s training distribution diverges from deployment settings Pomerleau [1989], Zhou et al. [2022]. These models generate diverse and multimodal trajectories, improving policy robustness in complex and uncertain environments.

The scalability of such IL methods has been further fueled by large-scale expert demonstration datasets et al. [2024], enabling models like Robotics Diffusion Transformer (RDT) Liu et al. [2024b] and π_0 Black et al. [2024] to generalize to novel manipulation tasks with minimal supervision. CoinRobot Zhao et al. [2025] pushes this further by proposing a unified IL architecture that generalizes across robot embodiments and sensor configurations. However, despite this progress, their reliance on vast computational resources limits accessibility and practical deployment.

Beyond policy structure and data scale, real-world performance increasingly depends on robust spatial perception—especially in dynamic, visually diverse settings. Autonomous driving research highlights the importance of combining RGB data with spatial cues to maintain performance under adverse conditions Qian et al. [2021], Zheng et al. [2023], Dong et al. [2023]. In robotics, spatial understanding has traditionally been enhanced through depth sensors, LiDAR, and multi-view 3D reconstruction Huang et al. [2024], Shridhar et al. [2021], Li et al. [2024], ApolloAuto [2023], Liu [2023]. While effective, these approaches often incur high costs and demand precise calibration, limiting their feasibility for many applications.

Monocular depth estimation provides a promising alternative, enabling low-cost depth perception directly from RGB images Birkl et al. [2023], Zavadski et al. [2024], Yang et al. [2024]. Modern transformer-based models like Depth Anything V2 and DINOv2 Yang et al. [2024], Oquab et al. [2024] offer real-time performance and robustness, making them well-suited for robotic deployment. However, their integration into IL pipelines remains limited, especially in settings that demand resilience to environmental variability. SVP Davies et al. [2024] demonstrated that augmenting RGB data with monocular depth and domain-randomized corruptions (via AugBlender) improves model robustness across lighting conditions. Their findings show that structured visual perception can mitigate performance collapse under extreme exposure variation—a common limitation in RGB-only pipelines.

To address this gap, recent research emphasizes multimodal perception and data augmentation strategies. Domain-randomized augmentations such as AugMix Hendrycks et al. [2020], Lan and Hao [2023], Raisuddin et al. [2023], Jin et al. [2024] simulate realistic visual perturbations and, when combined with depth estimation, can help models learn more invariant spatial representations.

While robust perception is essential, manipulation performance ultimately hinges on accurate grasp planning. Affordance-based methods have emerged as a promising strategy by providing spatial priors that guide grasp selection Kleeberger et al. [2020]. Point-based affordances identify object locations but often lack sufficient detail for determining stable grasps Liu et al. [2024a], Tang et al. [2024], Huang et al. [2024]. In contrast, grasp-based affordances encode spatially grounded, actionable cues. Large-scale datasets like Grasp Anything Vuong et al. [2023] offer a path toward scalable training, but their integration with diffusion-based IL remains underexplored. Early efforts like GQCNN Mahler et al. [2017] have demonstrated the potential of combining affordance reasoning with robotic control. Grasp-based affordances, such as those used in RoboGrasp Huang et al. [2025], explicitly provide grasp poses or grasping boxes that enable consistent and precise manipulation. Their integration with diffusion policies improves few-shot and prompted grasp generalization, especially in tasks requiring spatial grounding.

These advances point to the need for a unified, lightweight framework that combines domain-randomized augmentation, monocular depth estimation, and spatially grounded grasp guidance. While each component has shown promise individually, their integration within a single pipeline for imitation learning remains underexplored. This gap motivates our proposed approach, which brings together these elements to enhance spatial reasoning and policy robustness in dynamic environments.

3 Methodology

Our approach addresses robust spatial perception and accurate grasping for robotic manipulation by integrating multimodal sensory information through a cohesive, structured pipeline. The architecture comprises four primary modules: AugFusion, Monocular Depth Estimation, Grasp Prediction, and a downstream Robotic Action Head as shown in Figure 1. RGB image inputs are augmented to simulate realistic environmental variability, enriched with depth information for precise spatial awareness, and further enhanced with grasp predictions for explicit spatial affordances. These components collaboratively produce a rich and robust observation embedding, which the diffusion-based robotic action head leverages to generate stable, accurate manipulation policies.

3.1 AugFusion

To improve perception robustness under real-world visual variability, we introduce AugFusion, a domain-randomized augmentation strategy that extends AugMix Hendrycks et al. [2020] by incorporating both in-distribution and out-of-distribution (OOD) RGB corruptions. Unlike AugMix, which focuses on mild, within-distribution transformations, AugFusion applies realistic corruptions such as lighting shifts, exposure changes, blur, and noise to simulate challenging conditions where RGB inputs degrade.

AugFusion uses a probabilistic mechanism controlled by parameter β to decide between mixing multiple augmentations or applying them sequentially. Mixing weights are drawn from a Dirichlet distribution (α), and a blending factor λ adjusts intensity (see Algorithm 1). This process encourages the model to shift reliance toward depth cues when RGB becomes unreliable.

Integrating AugFusion-generated RGB data into our multimodal perception framework significantly broadens the model’s training distribution, enabling it to develop robust, invariant features essential for reliable real-world deployment under dynamic and unpredictable conditions.

Algorithm 1 AugFusion

Require: Image x , number of chains/augmentations k , parameter α , logic gate threshold β , mixing parameter λ

- 1: Randomly select $\xi \in [0, 1]$
- 2: Mixing weights: $w \leftarrow \text{Dirichlet}(\alpha)$
- 3: Augmentations: $A \leftarrow \{a_1, \dots, a_n\}$
- 4: $\lambda \leftarrow \begin{cases} 1, & \text{if } \xi < \beta \\ \lambda, & \text{otherwise} \end{cases}$
- 5: $x_t \leftarrow x$
- 6: **for** i in $\{1, \dots, k\}$ **do**
- 7: $x_{\text{aug}} \leftarrow x$
- 8: Randomly select $a \subseteq A$ such that $|a| = k$
- 9: Randomly select chain length $L \in \{1, \dots, k\}$
- 10: **if** $\xi > \beta$ **then**
- 11: **for** a_i in $a \subseteq \{a_1, \dots, a_L\}$ **do**
- 12: $x_{\text{aug}} \leftarrow a_i(x_{\text{aug}})$
- 13: **end for**
- 14: $x_t \leftarrow x_t + w_i \cdot x_{\text{aug}}$
- 15: **else**
- 16: $x_t \leftarrow a_i(x_t)$
- 17: **end if**
- 18: **end for**
- 19: $y \leftarrow \lambda \cdot x_t + (1 - \lambda) \cdot x$
- 20: **return** y

3.2 Monocular Depth Estimation Module

Integrating depth information significantly enhances robustness by enriching our model’s multimodal perception, allowing it to adaptively leverage complementary modalities during both training and inference. Depth data provides essential geometric context, which is especially beneficial for maintaining consistent performance under challenging perceptual variations, such as those caused by lighting changes. We specifically utilize monocular depth estimation due to its reliance solely on RGB images, enabling straightforward integration without the additional cost or complexity of dedicated hardware. Recent advances in monocular depth estimation models have demonstrated promising real-time capabilities combined with high accuracy and robustness, making them particularly suited to our robotic application Birk et al. [2023], Zavadski et al. [2024], Yang et al. [2024].

For our implementation, we adopt the Depth Anything V2 model, which leverages the transformer-based DINOv2 architecture explicitly optimized for robust monocular depth estimation Oquab et al. [2024]. Depth Anything V2 has shown impressive resilience to common image corruptions prevalent in robotic scenarios, aligning perfectly with our objective of achieving robust perception.

To enhance computational efficiency, we preprocess our entire training dataset by extracting depth maps from each RGB frame using the ViT-B-based variant of Depth Anything V2. This preprocessing step ensures tight alignment between RGB and depth modalities while significantly reducing memory consumption and accelerating the training pipeline. During inference, we employ the lighter ViT-S-based model variant, which maintains high-quality depth estimation at reduced computational costs, thereby achieving real-time inference suitable for deployment in resource-constrained robotic systems.

3.3 Grasp Prediction Module

To provide explicit grasp affordances for downstream manipulation, we employ a YOLO-based grasp detection module trained to predict 2D-oriented grasping boxes from RGB inputs. This module

significantly extends grasping capability by leveraging rich depth-aware embeddings derived from the monocular depth estimation model, enabling accurate grasp localization without the computational overhead of explicit point-cloud generation. Specifically, we fine-tuned a lightweight YOLOv11-m model on a custom-labeled dataset containing annotated graspable regions, where each annotation consists of a 5D grasp representation (as shown in Figure 2):

$$(x, y, w, h, \theta)$$

where (x, y) is the center of the grasp box in image coordinates, (w, h) denote its dimensions, and θ represents the in-plane rotation relative to the image frame.

While YOLO operates purely in 2D, we convert its predictions into full 6-degree-of-freedom (6-DoF) grasp poses through a lightweight, geometry-based postprocessing step. Given the estimated depth map from the monocular depth module and known camera intrinsics (f_x, f_y, c_x, c_y) , we project the 2D grasp center (x, y) into 3D space as follows:

$$x_{3D} = \frac{(x - c_x) \cdot z}{f_x}, \quad y_{3D} = \frac{(y - c_y) \cdot z}{f_y}, \quad z_{3D} = z$$

where f_x, f_y are focal lengths, (c_x, c_y) the camera’s principal point, and z the depth sampled from the predicted depth map at pixel (x, y) .

To determine the rotation component of the grasp pose, we use the predicted in-plane angle θ to construct a rotation matrix, assuming a top-down grasp direction (aligned along the camera optical axis Z):

$$R = [\mathbf{x}, \mathbf{y}, \mathbf{z}], \quad \text{where } \mathbf{x} = [\cos \theta, \sin \theta, 0], \quad \mathbf{z} = [0, 0, 1], \quad \mathbf{y} = \mathbf{z} \times \mathbf{x}$$

The resulting rotation matrix R can be converted to a quaternion representation (q_x, q_y, q_z, q_w) for efficient use in downstream robotic planning.

Thus, the final predicted grasp pose is represented as:

$$g = (x_{3D}, y_{3D}, z_{3D}, q_x, q_y, q_z, q_w)$$

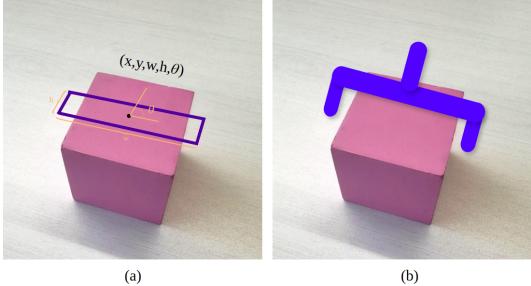


Figure 2: Grasp prompt visualization in the Pick-Big task. (a) Predicted oriented 2D grasp box from the RGB image. (b) Corresponding 6-DoF grasp prompt derived using camera intrinsics and a rotation matrix.

point cloud processing or retraining complex grasp regressors.

Furthermore, this setup remains fully compatible with existing grasp annotations and seamlessly integrates into our multimodal pipeline. The resulting 6-DoF grasp predictions provide actionable, contact-aware guidance to the diffusion-based policy, enhancing spatial reasoning, generalization, and manipulation precision in diverse, real-world environments.

3.4 Observation Encoder

The observation encoder fuses multiview RGB inputs, low-dimensional robot states, and grasp-specific features into a unified spatial-temporal representation for diffusion-based policy learning. Each fixed camera view is processed independently by a ResNet-34-based Feature Pyramid Network (FPN) Lin et al. [2017], He et al. [2015], enabling multi-scale feature extraction tailored to different viewpoints. Features from all views are pooled and concatenated into a comprehensive visual embedding.

This representation is augmented with AugFusion-processed RGBs, monocular depth maps, and 6-DoF grasp poses predicted by the grasp detection module. Additionally, robot end-effector state, gripper status, and task prompt embeddings are concatenated and linearly projected into a fixed-dimensional token per timestep.

A lightweight transformer then applies self-attention over tokens from the previous two timesteps, producing a spatially grounded, temporally aware encoding to condition the diffusion policy

3.5 Robotic Action Head

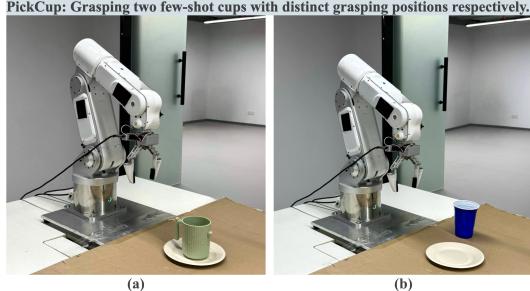


Figure 4: Few-shot PickCup setup. (a) Green mug with 5-shot handle grasping. (b) Blue plastic cup with 10-shot diameter grasping.

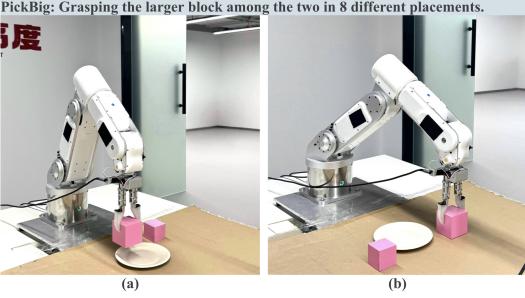


Figure 3: Placement generalization setup for the PickBig task. (a) and (b) illustrate two of the eight object configurations. The goal is to identify and grasp the larger of two similarly shaped blocks along its diameter.

Our action head adopts a diffusion-based policy Chi et al. [2024a], which models action generation as a denoising process from Gaussian noise to expert trajectories. Instead of direct regression, it refines sampled noise over 16 timesteps using a DDIM scheduler Nichol and Dhariwal [2021] with a cosine beta schedule, capturing multimodal and stochastic dynamics.

A transformer with cross-attention layers conditions the denoising on observation tokens. Actions are projected into a latent space, iteratively refined, and reprojected back via a linear head. Training uses Root Mean Square Error (RMSE) loss for its precision sensitivity, enabling the model to produce smooth and goal-aligned action sequences from multimodal inputs.

4 Experiments and Evaluation

Robotic manipulation in real-world settings demands precise grasping and strong generalization across diverse objects and conditions, yet prior imitation learning studies often rely on limited, controlled lab setups. To bridge this gap, we introduce a comprehensive experimental suite—PickBig, PickCup, and PickGoods—that evaluates generalization across object scales, categories, grasp strategies, and few-shot or prompt-driven scenarios. Our ablation studies analyze the impact of grasp affordances, monocular depth, and augmentation, while controlled exposure tests benchmark visual robustness under realistic lighting variation.

4.1 Task Description

PickBig: This task tests the robot’s ability to select and grasp the larger of two visually similar blocks placed in eight spatial configurations. The challenge lies in discerning subtle size differences and adapting the grasp accordingly (see Figure 3).

PickCup: This task evaluates generalization across diverse cup geometries and grasp types—handle, wall, and diameter grasps—across four positions. Few-shot trials with new cup types assess the policy’s ability to extend learned strategies with minimal data, guided by 6-DoF grasp representations (see Figure 4).

PickGoods: This task simulates open-world manipulation using varied consumer items. Each object is associated with a fixed grasp strategy, and a grasping prompt defines the target. This tests whether the policy can execute goal-directed actions when spatial intent is made explicit (see Figure 5).

- PickBig: 600 trials across 8 placement configurations, distinguishing and grasping the larger of two similar blocks.
- PickCup: 315 trials covering three cup types with varying grasping strategies (handle, sidewall, diameter), plus 15 few-shot examples for underrepresented cups.
- PickGoods: 400 trials involving four retail items, each paired with a consistent grasp strategy to evaluate promptable grasping.

4.2 Data Processing

We curate task-specific datasets to support robust and generalizable grasp policy learning: Each demonstration includes synchronized RGB and monocular depth frames. A representative subset is annotated with 2D grasp boxes, used to fine-tune our grasp detector. These 2D grasp annotations are then converted into 6-DoF grasp poses using depth and fixed camera intrinsics—projecting 2D centers to 3D, estimating orientation from in-plane angles, and deriving gripper widths from box dimensions.

This automated grasp pipeline generates consistent 6-DoF annotations for the full dataset and enables real-time grasp prediction at inference. By eliminating the need for point cloud reconstruction, it supports scalable, accurate, and contact-aware grasp supervision integrated directly into our multimodal observation encoder.

4.3 Evaluation Metrics

To assess both task-level performance and environmental robustness, we employ a set of complementary metrics across all experiments. Specifically, we evaluate:

Task Success Rate (TSR): The percentage of successful task completions, measuring whether the robot accomplishes the overall objective.

Grasp Success Rate (GSR): The proportion of stable and accurate grasp executions, capturing action consistency and spatial precision. It is computed as:

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

To rigorously evaluate robustness under lighting variation, we conduct controlled tests across 10 discrete camera exposure levels (10 to 170 ms). For each exposure setting, every model is evaluated over 100-150 trials per task. Success is determined through consensus among 2–3 human evaluators, who assess whether the manipulation was executed correctly. Final scores are computed by averaging success rates across all exposure levels, providing a holistic view of the model’s adaptability and reliability in realistic, visually corrupted conditions.

4.4 Computation Resources

Our setup used a standard industrial robot arm with two RGB cameras: one front-facing and one wrist-mounted for close-up views. As most imitation learning methods are designed for simulation, we selected Diffusion Policy (DP) as our baseline due to its proven effectiveness in real-world settings. We retrained DP with extensive hyperparameter tuning and applied our improvements under similar configurations. Given the unreliability of loss curves in robotic learning, we trained each model for a fixed duration based on expert heuristics. Training typically took around 48 hours per model using a dual RTX 4090 GPU server.

5 Results and Discussion

Table 1 summarizes TSR and GSR across models and exposure levels. Our full model consistently outperforms all baselines across tasks and lighting conditions, validating the synergy of 6-DoF grasp affordances, monocular depth estimation, and AugFusion augmentation.

Each module contributes complementary benefits: depth alone boosts PickBig TSR by 34%, highlighting the value of geometric priors; AugFusion improves PickCup TSR by 24% under mid-range exposures; GraspPrompt enhances spatial precision in PickBig and PickGoods but struggles in clutter. When combined, these components yield the most robust performance, outperforming baselines by 15–30% across metrics.

5.1 Lighting Robustness Across Exposure Conditions

Our ablation study reveals the limitations of existing diffusion policy variants when operating under environmental variability. Across all tasks, baseline DP and its data-augmented variants (e.g., DP+Depth, DP+AugFusion) exhibit steep performance declines under extreme lighting—particularly at low (10–40 ms) and high (160–170 ms) exposures. In contrast, our method maintains high TSR and GSR across the full exposure range, peaking around natural lighting (80–120 ms) and degrading gracefully at extremes. This is a direct result of our structured perception pipeline and its robustness to visual corruptions.

For instance, in the PickBig task, our method achieves an average TSR of 82% and GSR of 81%, outperforming all four baselines by over 14% in both metrics. In PickCup, where lighting variation strongly affects visual cues for grasping handles or rims, our model sustains 82% TSR and 80% GSR, compared to DP’s 27% and 21%, respectively. This confirms our approach’s resilience to appearance changes through both visual and geometric embeddings.

5.2 Grasp Precision and Strategy Learning

The inclusion of GraspPrompts significantly improves performance across all grasping strategies. In PickCup, our model reliably distinguishes handle, wall, and diameter grasps with minimal demonstrations—even generalizing to unseen cups in few-shot settings. The prompt-driven PickGoods task further highlights our model’s capacity to attend to spatial objectives: while DP and other baselines perform near-zero due to ambiguity in cluttered scenes, our model reaches a 47% TSR and 44% GSR despite spatial complexity and distractors. This demonstrates the value of explicit spatial prompting for policy grounding.

5.3 Discussion on Generalization, Few-Shot, and Promptability

Generalization to Diverse Objects: Our results in PickBig show how spatial variation challenges policies relying on implicit low-dim state. Our model’s grasp predictions provide spatial priors that allow robust re-localization and differentiation, even for near-identical objects.

Few-shot Transfer: In PickCup, 5–10 demo generalization to new objects shows the model can abstract strategy from shape and appearance when supported by geometric priors.

Prompt-following: PickGoods highlights that spatial cues alone are insufficient without environmental diversity in training. Our model still succeeds by disambiguating goal directionality (e.g., chocolate vs. tissue), but future work should incorporate more clutter-aware prompting and pose diversity.

Table 1: Success rates (%) for Task Success Rate (TSR) and Grasp Success Rate (GSR) across exposure levels. Each model spans two rows. AVG is the mean across all exposures. Bold font represents our better performance.

Task	Model	Metrics	10	20	40	60	80	100	120	140	160	170	AVG
PickBig	DP	TSR	9	19	27	53	59	87	68	39	33	28	42
		GSR	0	12	25	52	56	86	66	38	32	21	39
	+Depth	TSR	53	82	78	75	83	82	90	75	70	68	76
		GSR	53	80	77	73	83	80	87	73	67	65	74
	+AugFusion	TSR	51	65	72	75	80	89	51	61	60	58	67
		GSR	49	62	70	73	78	87	49	59	58	56	64
	+GraspPrompt	TSR	52	63	78	85	87	94	98	77	54	31	72
		GSR	51	60	75	84	83	94	96	72	51	29	70
	Ours	TSR	62	65	76	89	91	95	98	93	82	72	82
		GSR	59	65	72	88	91	95	98	91	80	67	81
PickCup	DP	TSR	0	0	0	53	75	82	40	17	3	0	27
		GSR	0	0	0	28	64	70	36	8	0	0	21
	+Depth	TSR	0	32	60	67	80	87	89	75	53	26	57
		GSR	0	27	41	55	64	76	80	71	48	22	48
	+AugFusion	TSR	0	32	50	57	64	76	81	69	53	26	51
		GSR	0	15	32	50	56	70	77	61	33	17	41
	+GraspPrompt	TSR	0	37	54	77	85	98	98	79	55	0	58
		GSR	0	35	52	73	82	98	98	70	49	0	56
	Ours	TSR	61	72	86	87	93	97	99	85	78	62	82
		GSR	60	69	85	85	93	96	99	84	73	60	80
PickGoods	DP	TSR	0	0	0	0	9	17	22	20	0	0	7
		GSR	0	0	0	0	8	8	22	10	0	0	5
	+Depth	TSR	0	0	0	16	25	34	41	28	0	0	14
		GSR	0	0	0	13	20	25	27	18	0	0	10
	+AugFusion	TSR	0	0	0	7	19	21	11	10	0	0	7
		GSR	0	0	0	3	8	14	5	2	0	0	3
	+GraspPrompt	TSR	0	0	0	19	31	49	52	36	12	0	20
		GSR	0	0	0	9	25	29	51	32	7	0	15
	Ours	TSR	26	33	48	52	57	63	65	51	44	29	47
		GSR	22	28	43	50	56	61	64	51	40	21	44

Our results demonstrate that combining structured 6-DoF grasp prompts, monocular depth, and robust visual augmentation enables a powerful spatial perception stack for generalizable and precise manipulation. The proposed approach consistently outperforms strong baselines in both task success and grasp execution, especially under environmental variation. This affirms our core claim: robust robotic grasping requires not just diverse inputs, but structured, spatially-grounded ones.

6 Conclusion

We present a unified framework that enhances both the spatial robustness and grasping generalization of robotic learning systems, enabling precise, contact-aware manipulation under diverse environmental conditions. By integrating AugFusion, monocular depth estimation, and 6-DoF grasp prediction into a cohesive perception module, our system significantly improves visuomotor policy performance without the need for additional 3D sensing hardware. Built atop a diffusion-based action model, this approach proves effective across a variety of tasks, object types, and lighting scenarios.

Through comprehensive experiments and ablation studies, we demonstrate that each module—depth, augmentation, and grasp prompt—contributes uniquely to policy robustness. Our full model achieves the highest task and grasp success rates across all exposure levels, showcasing strong few-shot generalization and reliable performance in unstructured, low-visibility environments. The results underscore the importance of spatially grounded perception in improving goal alignment, action precision, and environmental adaptability.

This work provides a scalable, plug-and-play solution for general-purpose robotic manipulation, bridging the gap between controlled lab setups and the complexities of real-world deployment.

7 Future Work

Building on our findings, future efforts should prioritize standardized benchmarks for imitation learning under real-world variability, similar to ImageNet-C or 4Seasons in computer vision. These would enable consistent evaluation of robustness and generalization across robotic systems. To streamline the grasping pipeline, an MLP grasp predictor could be trained end-to-end alongside depth and AugFusion features, replacing the standalone YOLO module. This approach may improve learning efficiency and offer tighter integration between visual cues and grasp prediction, while still supporting full 6-DoF outputs.

Improving temporal reasoning remains an open challenge. Our method processes short sequences but lacks long-horizon memory. Integrating memory-augmented architectures—such as Scene Memory Transformer or REMEMBR—could support policy continuity in dynamic or delayed-reward scenarios. Finally, extending our framework with language-conditioned prompting (e.g., via Grounding DINO or DINO-X) and applying grasp affordance prompts in world models or large-scale policy frameworks like ACT or Robotics Diffusion Transformer would further enhance task generality and semantic grounding. These directions aim to bridge spatial reasoning, multimodal inputs, and scalable policy learning in real-world environments.

References

- ApolloAuto. Apollo: An open autonomous driving platform. <https://github.com/ApolloAuto/apollo>, 2023. Accessed: October 18, 2023.
- Reiner Birlk, Diana Wofk, and Matthias Müller. Midas v3.1 – a model zoo for robust monocular relative depth estimation, 2023. URL <https://arxiv.org/abs/2307.14460>.
- 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. URL <https://arxiv.org/abs/2410.24164>.
- Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- 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, 2024a. URL <https://arxiv.org/abs/2303.04137>.
- Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots, 2024b. URL <https://arxiv.org/abs/2402.10329>.
- Travis Davies, Jiahuan Yan, Xiang Chen, Yu Tian, Yuetong Zhuang, Yiqi Huang, and Luhui Hu. Spatially visual perception for end-to-end robotic learning, 2024. URL <https://arxiv.org/abs/2411.17458>.
- Yinpeng Dong, Caixin Kang, Jinlai Zhang, Zijian Zhu, Yikai Wang, Xiao Yang, Hang Su, Xingxing Wei, and Jun Zhu. Benchmarking robustness of 3d object detection to common corruptions in autonomous driving, 2023. URL <https://arxiv.org/abs/2303.11040>.
- Embodiment Collaboration et al. Open x-embodiment: Robotic learning datasets and rt-x models, 2024. URL <https://arxiv.org/abs/2310.08864>.
- Zipeng Fu, Tony Z. Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation, 2024. URL <https://arxiv.org/abs/2401.02117>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. URL <https://arxiv.org/abs/1512.03385>.
- Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty, 2020. URL <https://arxiv.org/abs/1912.02781>.

- Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang, and Li Fei-Fei. Rekep: Spatio-temporal reasoning of relational keypoint constraints for robotic manipulation, 2024. URL <https://arxiv.org/abs/2409.01652>.
- Yiqi Huang, Travis Davies, Jiahuan Yan, Xiang Chen, Yu Tian, and Luhui Hu. Robograsp: A universal grasping policy for robust robotic control, 2025. URL <https://arxiv.org/abs/2502.03072>.
- Yeying Jin, Beibei Lin, Wending Yan, Yuan Yuan, Wei Ye, and Robby T. Tan. Enhancing visibility in nighttime haze images using guided apsf and gradient adaptive convolution, 2024. URL <https://arxiv.org/abs/2308.01738>.
- Kilian Kleeberger, Richard Bormann, Werner Kraus, and Marco F Huber. A survey on Learning-Based robotic grasping. *Current Robotics Reports*, 1(4):239–249, December 2020.
- Gongjin Lan and Qi Hao. End-to-end planning of autonomous driving in industry and academia: 2022-2023, 2023. URL <https://arxiv.org/abs/2401.08658>.
- Jinhan Li, Yifeng Zhu, Yuqi Xie, Zhenyu Jiang, Mingyo Seo, Georgios Pavlakos, and Yuke Zhu. Okami: Teaching humanoid robots manipulation skills through single video imitation. In *8th Annual Conference on Robot Learning (CoRL)*, 2024.
- Siyuan Li, Xun Wang, Rongchang Zuo, Kewu Sun, Lingfei Cui, Jishiyu Ding, Peng Liu, and Zhe Ma. Robust visual imitation learning with inverse dynamics representations, 2023. URL <https://arxiv.org/abs/2310.14274>.
- Fanqi Lin, Yingdong Hu, Pingyue Sheng, Chuan Wen, Jiacheng You, and Yang Gao. Data scaling laws in imitation learning for robotic manipulation. *arXiv preprint arXiv:2410.18647*, 2024.
- Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection, 2017. URL <https://arxiv.org/abs/1612.03144>.
- Fangchen Liu, Kuan Fang, Pieter Abbeel, and Sergey Levine. Moka: Open-world robotic manipulation through mark-based visual prompting, 2024a. URL <https://arxiv.org/abs/2403.03174>.
- Langechuan Patrick Liu. The practice of mass production autonomous driving. Presented at the CVPR 2023 E2EAD Workshop, 2023. Available at <https://opendrivelab.com/e2ead/cvpr23>.
- 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, 2024b. URL <https://arxiv.org/abs/2410.07864>.
- Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. 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>.
- Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models, 2021. URL <https://arxiv.org/abs/2102.09672>.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. URL <https://arxiv.org/abs/2304.07193>.
- Dean Pomerleau. Alvinn: An autonomous land vehicle in a neural network. In D.S. Touretzky, editor, *Proceedings of (NeurIPS) Neural Information Processing Systems*, pages 305 – 313. Morgan Kaufmann, December 1989.
- Kun Qian, Shilin Zhu, Xinyu Zhang, and Li Erran Li. Robust multimodal vehicle detection in foggy weather using complementary lidar and radar signals. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 444–453, 2021. doi: 10.1109/CVPR46437.2021.00051.
- Abu Mohammed Raisuddin, Tiago Cortinhal, Jesper Holmlad, and Eren Erdal Aksoy. 3d-outdet: A fast and memory efficient outlier detector for 3d lidar point clouds in adverse weather. October 2023. doi: 10.36227/techrxiv.24297166.v1. URL <http://dx.doi.org/10.36227/techrxiv.24297166.v1>.
- 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. *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>.

Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation, 2021. URL <https://arxiv.org/abs/2109.12098>.

Grace Tang, Swetha Rajkumar, Yifei Zhou, Homer Rich Walke, Sergey Levine, and Kuan Fang. Kalie: Fine-tuning vision-language models for open-world manipulation without robot data, 2024. URL <https://arxiv.org/abs/2409.14066>.

ALOHA 2 Team, Jorge Aldaco, Travis Armstrong, Robert Baruch, Jeff Bingham, Sankt Chan, Kenneth Draper, Debidatta Dwibedi, Chelsea Finn, Pete Florence, Spencer Goodrich, Wayne Gramlich, Torr Hage, Alexander Herzog, Jonathan Hoech, Thinh Nguyen, Ian Storz, Baruch Tabanpour, Leila Takayama, Jonathan Tompson, Ayzaan Wahid, Ted Wahrburg, Sichun Xu, Sergey Yaroshenko, Kevin Zakka, and Tony Z. Zhao. Aloha 2: An enhanced low-cost hardware for bimanual teleoperation, 2024. URL <https://arxiv.org/abs/2405.02292>.

An Dinh Vuong, Minh Nhat Vu, Hieu Le, Baoru Huang, Binh Huynh, Thieu Vo, Andreas Kugi, and Anh Nguyen. Grasp-anything: Large-scale grasp dataset from foundation models, 2023. URL <https://arxiv.org/abs/2309.09818>.

Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2, 2024. URL <https://arxiv.org/abs/2406.09414>.

Denis Zavadski, Damjan Kalšan, and Carsten Rother. Primedepth: Efficient monocular depth estimation with a stable diffusion preimage, 2024. URL <https://arxiv.org/abs/2409.09144>.

Yu Zhao, Huxian Liu, Xiang Chen, Jiankai Sun, Jiahuan Yan, and Luhui Hu. Coinrobot: Generalized end-to-end robotic learning for physical intelligence, 2025. URL <https://arxiv.org/abs/2503.05316>.

Tianyue Zheng, Ang Li, Zhe Chen, Hongbo Wang, and Jun Luo. Autofed: Heterogeneity-aware federated multimodal learning for robust autonomous driving, 2023. URL <https://arxiv.org/abs/2302.08646>.

Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, page 1–20, 2022. ISSN 1939-3539. doi: 10.1109/tpami.2022.3195549. URL <http://dx.doi.org/10.1109/TPAMI.2022.3195549>.

A Technical Appendices and Supplementary Material

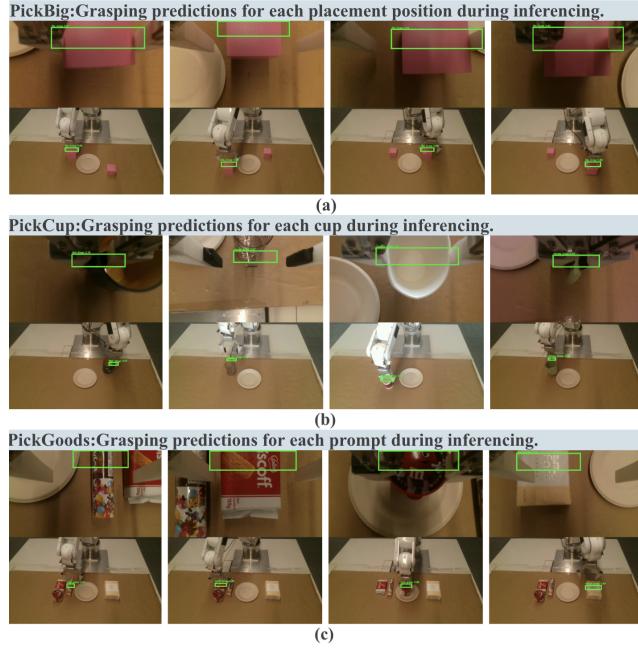


Figure 6: Real-time orientated-2D grasp box predictions across tasks using Spatial RoboGrasp. It will be later converted to 6-DoF grasp prompts: (a) Robust grasp detection across varied placement configurations in PickBig. (b) Accurate strategy-specific predictions for diverse cup geometries in PickCup. (c) Prompt-guided grasp localization for goal-driven manipulation in PickGoods.

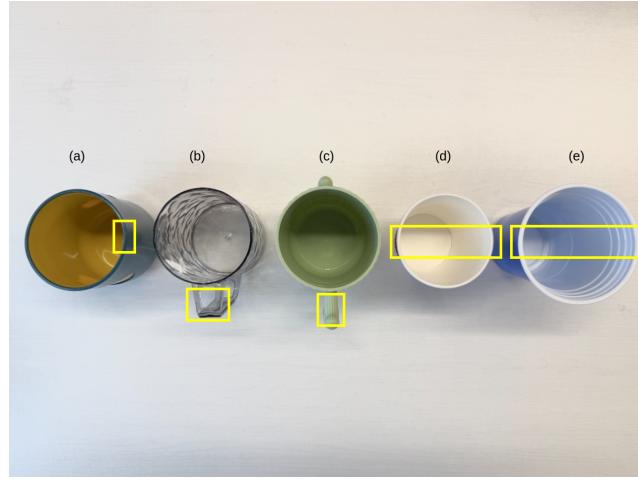


Figure 7: Illustration of grasping strategies for the PickCup task (top-down view). (a) shows a sidewall grasp; (b) and (c) illustrate handle grasps; (d) and (e) depict diameter grasps. (c) and (e) correspond to the cups used in the few-shot evaluation.