

Towards Generalizable Vision-Language Robotic Manipulation: A Benchmark and LLM-guided 3D Policy



PSI 🖈

Ricardo Garcia*, Shizhe Chen*, Cordelia Schmid

Inria Paris, École normale supérieure, PSL

Introduction



"put the frog toy in the top drawer"

Goal: Enhance the generalization capabilities of vision-language robotic manipulation policies.

Limitations of state-of-the-art methods:

- Train and test policies on the same task set
- Focus on a limited set of action skills

(pick-and-place)

Our contribution:

- 1. A comprehensive benchmark: covering 7 action skills and 4 generalization levels
- 2. A generalist **LLM-guided 3D policy**:
- + 3D-based robotic manipulation policy: more precise action prediction
- + Integration with LLMs and VLMs: improved generalization ability

GEMBench: GEneralizable Vision-Language Robotic Manipulation Benchmark

Real Robot Setups



Test set: 44 tasks (92 variations) / 4 levels of generalization.







"stack the (yellow/navy) cup on top of the (pink/yellow) cup"

"put the (strawberry/peach) in the box"







"put the pink mug on the middle part of the hanger"

"put the frog toy in the top drawer"

"open the top drawer'

Unseen Tasks



"place the yellow cup inside the "stack the (black/red) cup on top of the (orange/black) cup" red cup, then the cyan cup on top"







"put the (lemon/banana) in the box"

"put the grapes in the yellow plate, "put the tuna can in the box, then put the corn in the box" then put the banana in the pink plate"

The Proposed Method

3D-LOTUS policy:

- Efficient Point Transformer v3 backbone
- Improve precision via point-wise classification

3D-LOTUS++ framework:





1 LLaMA-7B performs task planning by splitting high-level goal instructions into a sequence of primitive actions.

(2) OwIViT v2 + SAM processes multiple views and the target object name for visual grounding.

③ 3D-LOTUS visuomotor policy predicts precise robot trajectory for a given target object and primitive action.

Experimental Results

Evaluation on RLBench-18Task

Achieved SoTA performance and faster training speed.

		Avg. SR \uparrow	Avg. Rank \downarrow	Train time \downarrow	Method	L1	L2	L3	L4	grounding	are the	main
C2F-ARM-BC [38] Hiveformer [17] PolarNet [2]		20.1 45.3 46.4	8.6 6.9 6.4	- - 8 0	Hiveformer [17]	$60.3_{\pm 1.5}$	$26.1_{\pm 1.4}$	$35.1_{\pm 1.7}$	$0.0_{\pm 0.0}$	robotic manipulation.		
Polarivet [2] 40.4 PerAct [18] 49.4 RVT [34] 62.9 Act3D [4] 65.0		6.2 4.4 4.3	128.0 8.0 40.0	3D diffuser actor [35] RVT-2 [37]	$77.7_{\pm 0.9}$ $91.9_{\pm 0.8}$ $89.1_{\pm 0.8}$	$37.1_{\pm 1.4} \\ 43.4_{\pm 2.8} \\ 51.0_{\pm 2.3}$	$38.5_{\pm 1.7}$ $37.0_{\pm 2.2}$ $36.0_{\pm 2.2}$	$\begin{array}{c} 0.1_{\pm 0.2} \\ 0.0_{\pm 0.0} \\ 0.0_{\pm 0.0} \end{array}$	Task Planning	Object Grounding	Avg.	
RVT2 [37] 3D diffuser actor [35] 3D-LOTUS		$81.4 \\ 81.3 \\ \textbf{83.1}_{\pm 0.8} \\ \textbf{83.1}_{\pm 0.8$	2.4 2.3 2.2	6.6 67.6 2.2 ³	3D-LOTUS 3D-LOTUS++	94.3 $_{\pm 1.4}$ 68.7 $_{\pm 0.6}$	$\begin{array}{c} 49.9_{\pm 2.2} \\ \textbf{64.5}_{\pm 0.9} \end{array}$	$\begin{array}{c} 38.1_{\pm 1.1} \\ \textbf{41.5}_{\pm 1.8} \end{array}$	$0.3_{\pm 0.3}$ 17.4 _{\pm 0.4}	GT GT LLM	GT VLM VLM	63.0 50.7 48.0
Real world results	Task Stack yellow cup in pink cup Stack navy cup in yellow cup Put strawberry in box Put peach in box Open drawer Put item in drawer Hang mug Avg.		PolarNet x cup 10/10 x cup 9/10 7/10 8/10 6/10 1/10 6/10	t3D-LOTUSTask9/10Stack red cup in yellow cup10/10Stack black cup in orange cup10/10Place the yellow cup inside the red cup,10/10then the cyan cup on top8/10Put lemon in box9/10Put banana in box3/10Put tuna can in box, then corn in box8/10Avg.		3D-LOTU 0/10 0/10 0/10 0/10 0/10 0/10 0/10 0/1	3D-LOTUS 3D-LOTUS++ 0/10 8/10 0/10 7/10 0/10 7/10 0/10 9/10 0/10 8/10 0/10 8/10 0/10 8/10 0/10 9/10 0/10 7.9/10		Project Webpage	CVPR 2025 Challenge & Workshop		
		Se	en Tasks		Unseen Tasks							

Evaluation on GemBench

3D-LOTUS++ performs better on more challenging generalization levels.

		Avg. SR ↑	Avg. Rank \downarrow	Train time ↓	Method	L1	L2	L3	L4	grounding	are the stor and the	main lizable
C2F-ARM-BC [38] Hiveformer [17] PolarNet [2]		20.1 45.3 46.4	8.6 6.9 6.4	- - 8.9	Hiveformer [17]	$60.3_{\pm 1.5}$	$26.1_{\pm 1.4}$	$35.1_{\pm 1.7}$	$0.0_{\pm 0.0}$	robotic ma	anipulation.	
PerAct [18] 49.4 RVT [34] 62.9 Act3D [4] 65.0		6.2 4.4 4.3	128.0 8.0 40.0	3D diffuser actor [35] RVT-2 [37]	$91.9_{\pm 0.8}$ $89.1_{\pm 0.8}$	$57.1_{\pm 1.4}$ $43.4_{\pm 2.8}$ $51.0_{\pm 2.3}$	$37.0_{\pm 2.2}$ $36.0_{\pm 2.2}$	$0.1 \pm 0.2 \\ 0.0 \pm 0.0 \\ 0.0 \pm 0.0$	Task Planning	Object Grounding	Avg.	
RVT2 [37] 3D diffuser actor [35]		81.4 81.3	2.4 2.3 2 2	6.6 67.6 2 2 ³	3D-LOTUS 3D-LOTUS++	94.3 _{±1.4} 68.7 _{±0.6}	$49.9_{\pm 2.2} \\ 64.5_{\pm 0.9}$	$38.1_{\pm 1.1} \\ \textbf{41.5}_{\pm 1.8}$	$0.3_{\pm 0.3}$ 17.4 $_{\pm 0.4}$	GT GT LLM	GT VLM VLM	63.0 50.7 48.0
Real world results	Task Stack yellow cup in pink cup Stack navy cup in yellow cup Put strawberry in box Put peach in box Open drawer Put item in drawer Hang mug		PolarNe k cup 10/10 w cup 9/10 7/10 8/10 6/10 6/10 6/10	t 3D-LOTUS 9/10 10/10 10/10 8/10 9/10 3/10 8/10 8/10	D-LOTUSTask9/10Stack red cup in yellow cup Stack black cup in orange cup10/10Place the yellow cup inside the red cup, then the cyan cup on top10/10Put lemon in box9/10Put lemon in box9/10Put banana in box9/10Put grapes in yellow plate, then banana in pink plate8.1/10Avg.		3D-LOTUS 3D-LOTUS++ 0/10 8/10 0/10 7/10 0/10 7/10 0/10 9/10 0/10 8/10 0/10 9/10 0/10 8/10 0/10 7/10 0/10 7/10 0/10 7/10 0/10 7/10 0/10 7/9/10		Project Webpage	CVPR 2025 Challenge & Workshop		
	Seen Tasks				Unsee	n Tasks						

Ablation on GemBench

The motion policy and object