

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

Hyperparameter-Tuned Object Grounding and Early-Stopped Motion Control Enhance 3D-LOTUS++ on GemBench

CVPR 2025 Gembench Challenge report

Team: MiRA Member: Poshenc Chen, Sin-Yi Chiu, Chuan-Yu Wu, Jia-Fong Yeh, Hung-Ting Su, Chih-Han Chen, Yan-Xiang Qiu Contact email: poshengchen@cmlab.csie.ntu.edu.tw

Abstract

001 This report presents our solution to the GemBench Challenge at CVPR 2025. Based on 3D-LOTUS++, 002 we modified the inference hyperparameters of the 003 VLMs and applied early stopping during the training 004 005 of the motion control policy. On the private test set, our method outperformed the best baseline by nearly 006 5%, showing significant gains in L2 (66.0) and L3 007 008 (49.4).

009 1. Introduction

010 In this challenge, GemBench offers a comprehensive suite of training and evaluation tasks designed to as-011 sess a model's generalization ability across diverse 012 013 manipulation scenarios. The training set comprises 16 distinct tasks and 31 task variants, encompassing 014 seven fundamental manipulation primitives. The test 015 set includes 44 tasks with a total of 92 variations, cat-016 egorized into four difficulty levels: novel placement, 017 novel rigid objects, novel articulated objects, and long-018 horizon tasks. A portion of the test set is kept private 019 and used for final evaluation on a hidden dataset. 020

021 Our submission is based on 3DLOTUS++. In 022 3DLOTUS++, we observed that the object grounding 023 inference occasionally filtered out the target object. 024 To address this issue, we adjusted the relevant hyper-025 parameters. Additionally, we identified an overfitting 026 problem during the training of the motion control pol-027 icy, which we mitigated by applying early stopping.

On the GemBench public benchmark, our method achieved the following performance scores across the four levels: L1: 68.32 ± 1.04 , L2: 63.57 ± 2.30 , L3: 41.95 \pm 1.18, and L4: 15.25 ± 1.17 . On the private dataset, our method attained L2: 66.0, L3: 49.4, and L4: 15.0.

034 1.1. Related Work

035 3DLOTUS++ [2] framework:

Task Planning We define six action primitives for036object manipulation. Human instructions are decomposed into a sequence of action primitives using the037LLaMA-3 8B [1] LLM. By providing a few in-context039examples, the LLM is guided to generate appropriate040action sequences.041

Object Grounding The OWLv2 [4] open-042 vocabulary detector generates high-confidence 043 bounding boxes and semantic embeddings from im-044 ages. These boxes are then segmented using SAM[3], 045 and combined with RGB-D data to produce 3D point 046 clouds. The point cloud is labeled into four categories: 047 goal object, goal target, robot, and obstacle. Among 048 them, the object and target are selected based on 049 OWLv2's semantic embeddings to match the textual 050 descriptions. 051

Motion Control The point cloud results from object grounding, along with the action primitives from task planning, are fed into a cross-attention encoder. This allows the model to make distinct predictions based on different targets and different actions.

2. Method

Object Grounding To refine the object grounding process during inference, we adjust several hyperparameters that control candidate box selection. The default parameters are defined as follows:

- threshold: Objectness score threshold. Only predictions with confidence scores higher than this value are retained.
- min_size_ratio: Minimum size ratio relative to the entire scene or input space. Objects smaller than this ratio are discarded to reduce noise.
- max_size_ratio: Maximum size ratio. Predictions exceeding this ratio are considered too large and removed.
- min_return_topk: Ensures that at least this number of object candidates are returned, regardless of filtering results.
- max_return_topk: Caps the number of returned object candidates to this value, preventing excessive
 074
 075

Table 1. Average performance (%) across GemBench Levels 2-4.

Method	L2	L3	L4
3D-LOTUS (baseline) 3D-LOTUS++ (baseline)	$\begin{array}{c} 49.9{\scriptstyle\pm2.2}\\ 64.5{\scriptstyle\pm0.9}\end{array}$	$\begin{array}{c} 38.1{\scriptstyle\pm1.1}\\ 41.5{\scriptstyle\pm1.8}\end{array}$	$\begin{array}{c} 0.3{\scriptstyle\pm0.3}\\ \textbf{17.4{\scriptstyle\pm0.4}}\end{array}$
MiRA (seed200 & 300) MiRA (ours)	66.0 ±0.1 63.6±2.3	$\begin{array}{c} 42.7{\scriptstyle\pm1.3} \\ 42.0{\scriptstyle\pm1.2} \end{array}$	16.3 ± 1.3 15.3 ± 1.2

 Table 2. Performance (%) across GemBench private dataset

Method	L2	L3	L4	avg.
3D-LOTUS (baseline) 3D-LOTUS++ (baseline)	13.0 58.0	52.2 41.1	0.0 17.2	21.7 38.8
MiRA (ours)	66.0	49.4	15.0	43.5

076 proposals.

080

081

082

083

- use_nms: Boolean flag indicating whether to apply
 Non-Maximum Suppression (NMS) to reduce over lapping candidates.
 - nms_sigma: Gaussian parameter for soft-NMS or suppression radius for standard NMS.
 - nms_thresh: IoU threshold used during NMS to suppress overlapping boxes.

During evaluation, we observed that several target objects were mistakenly filtered out due to overly restrictive parameter settings in tasks objects (e.g., the cupboard in the Put In Cupboard task). So, we adjusted **max_return_topk** and **min_size_ratio** to mitigate this issue.

Motion Control During Motion Control training, 090 five loss terms are tracked: total, pos, rot, open, and 091 stop. Each converges at a different pace, resulting 092 093 in distinct optimal training lengths for the respective components. We observed that the key to achieving the 094 best performance is allowing the pos loss to converge 095 without overfitting. Therefore, our early-stopping cri-096 terion is defined solely by the behaviour of this pos 097 098 loss.

3. Experimental and Result

100 3.1. Experimental Setup

- 101 Object Grounding Our experiments use the set of hyperparameters:
- We keep proposals whose objectness score exceeds
 0.1 (threshold = 0.1).
- Box size filtering removes extremely small and large candidates: we discard any box whose projected area is < 0.15 % of the scene (min_size_ratio = 0.0015) or > 80 % (max_size_ratio = 0.8).
- To guarantee coverage, we always return at least one proposal but never more than twenty (min_return_topk = 1, max_return_topk = 20).
- Non-maximum suppression is enabled (use_nms = 113 True) with a Gaussian soft-NMS decay of $\sigma = 0.2$ and an IoU suppression threshold of 0.1 (nms_sigma

= 0.2, nms_thresh = 0.1).

115

127

128

129

130

131

132

133

134

135

136

The hyperparameters that differ from the 116 3D-LOTUS++ setup are max_return_topk and 117 min_size_ratio. We observed that the original settings 118 also caused some intended targets to be removed 119 (for example, the cupboard in the "Put In Cupboard" 120 task), which led to incorrect point-cloud labeling and 121 ultimately affected action prediction. To remedy this, 122 we increased max_return_topk from 10 to 20. This 123 change, however, introduced many small spurious 124 boxes, so we also raised min_size_ratio from 0 to 125 0.0015 126

Motion Control From the training loss curves, we noticed that the validation rotation loss began overfitting very early, model_step_30000 already gave the lowest value. The position loss converged at roughly 127000 steps and started overfitting by 140000 steps (the baseline checkpoint). After evaluating model_step_30000 and model_step_127000, we chose the better-performing model_step_127000 as our submission model.

3.2. result

We evaluated our method on the GemBench public test 137 set. As shown in Table 1, our approach demonstrates 138 a slight decrease in the overall average task success 139 rate across all seeds compared to the official baseline. 140 However, our method exhibits modest improvements 141 specifically on seed 200 & 300, suggesting potential 142 under certain initialization conditions. The detailed re-143 sults for each task are provided in the Supplementary 144 Material. 145

As shown in Table 2, our performance on the private test dataset improved by nearly 5% compared to the best-performing baseline. Specifically, when compared to the top baseline 3D-LOTUS++, our method achieved noticeable gains in both L2 and L3 tasks. However, performance on L4 still lags slightly behind. 152

CVPR

#

Based on 3D-LOTUS++, we tuned its hyperparame-154 ters and achieved better performance than the official 155 baseline on the private dataset. We identified short-156 comings in both the object-grounding and motion-157 control modules and introduced improvements. Al-158 159 though additional experiments are required to verify 160 that our adjustments are optimal, we believe these two components still offer considerable room for further 161

research toward more broadly generalized robotics.

163 References

- 164 [1] AI@Meta. Llama 3 model card. 2024. 1
- [2] Ricardo Garcia, Shizhe Chen, and Cordelia Schmid. Towards generalizable vision-language robotic manipulation: A benchmark and llm-guided 3d policy. In *Interna- tional Conference on Robotics and Automation (ICRA)*,
 2025. 1
- [3] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 1
- [4] Matthias Minderer, Alexey Gritsenko, and Neil
 Houlsby. Scaling open-vocabulary object detection. Advances in Neural Information Processing Systems, 36:
 72983–73007, 2023. 1

Ħ

182

Hyperparameter-Tuned Object Grounding and Early-Stopped Motion Control Enhance 3D-LOTUS++ on GemBench

Supplementary Material

Table 3. Performance on GemBench Level 1.

Method	Avg.	Close Fridge+0	Close Jar+15	Close Jar+16	CloseLaptop Lid+0	Close Microwave+0	LightBulb In+17	LightBulb In+19	Open Box+0	Open Door+0	Open Drawer+0
3D-LOTUS (baseline)	94.3 ±3.5	96±3.7	100 ± 0.0	100 ± 0.0	98±2.5	98±4.0	84±7.4	85±9.5	99 ±2.0	77±2.5	83±18.7
3D-LOTUS++ (baseline)	$68.7{\pm}0.6$	$95{\pm0.0}$	100 ± 0.0	$99{\scriptstyle\pm2.0}$	28±2.5	87±5.1	55±10.5	45±8.9	55±8.9	79±9.7	68±11.2
MiRA (ours)	68.32±6.01	95±3.16	98±2.45	99±2	84±2	84±7.35	42±5.1	48±12.08	66±10.68	79±6.63	52±11.22
Method	Open Drawer+2	Pick& Lift+0	Pick& Lift+2	Pick& Lift+7	PickUp Cup+8	PickUp Cup+9	PickUp Cup+11	Push Button+0	Push Button+3	Push Button+4	PutIn Cupboard+0
3D-LOTUS (baseline)	93±6.0	99±2.0	100 ±0.0	100 ±0.0	93±4.0	94±3.7	94±4.9	99±2.0	100 ± 0.0	100 ± 0.0	89±5.8
3D-LOTUS++ (baseline)	75±4.5	97±6.0	94±3.7	$93{\pm}5.1$	88±6.6	88±6.6	91±4.9	100 ± 0.0	100 ± 0.0	100 ± 0.0	1 ± 2.0
MiRA (ours)	55±13.04	98±2.45	93±6.78	90±4.47	91±3.74	84±7.35	88±6	100 ± 0	100 ± 0	100 ± 0	9±5.83
Method	PutIn Cupboard+3	PutMoney InSafe+0	PutMoney InSafe+1	Reach& Drag+14	Reach& Drag+18	Slide Block+0	Slide Block+1	Stack Blocks+30	Stack Blocks+36	Stack Blocks+39	
3D-LOTUS (baseline)	72±11.2	94±3.7	94±3.7	100 ± 0.0	100 ± 0.0	100 ± 0.0	100 ± 0.0	91±6.6	90±4.5	90±5.8	
3D-LOTUS++ (baseline)	2±2.5	22±6.8	16±4.9	94±3.7	62±8.7	100 ± 0.0	65±5.5	86±5.8	20±4.5	28±13.6	
MiRA (ours)	15±10.49	22±9.8	4±3.74	56±8.6	27±2.45	100 ±0	73±8.12	65±8.37	59±11.58	42±10.77	

Table 4. Performance on GemBench Level 2.

Method	Avg.	Push Button+13	Push Button+15	Push Button+17	Pick& Lift+14	Pick& Lift+16	Pick& Lift+18	PickUp Cup+10	PickUp Cup+12	PickUp Cup+13
3D-LOTUS (baseline)	$49.9{\scriptstyle\pm2.2}$	99±2.0	100 ± 0.0	100 ± 0.0	3±2.5	18±8.7	33±9.3	89±3.7	78±8.7	57±7.5
3D-LOTUS++ (baseline)	64.5 ± 0.9	99±2.0	100 ± 0.0	99±2.0	94±3.7	96±3.7	95±3.2	79±4.9	89±9.7	$84{\pm}10.2$
MiRA (ours)	63.57±6.33	100 ± 0	100 ± 0	100 ± 0	$89{\pm}10.68$	91±3.74	86±8.6	78±9.27	86±10.68	82±10.3
Method	Stack Blocks+24	Stack Blocks+27	Stack Blocks+33	Slide Block+2	Slide Block+3	Close Jar+3	Close Jar+4	LightBulb In+1	LightBulb In+2	Lamp On+0
3D-LOTUS (baseline)	13±8.1	40±9.5	69±5.8	1±2.2	1±2.2	71±5.8	90±4.5	24±4.9	81±6.6	$0{\pm}0.0$
3D-LOTUS++ (baseline)	22±9.3	83±7.5	63±7.3	27±9.8	5±3.2	$98{\pm}2.5$	98±2.5	56±9.7	43±7.5	2 ± 2.0
MiRA (ours)	46±5.83	74±7.35	68±6.78	45±9.49	16±6.63	97±4	86±7.35	46±11.58	63±6.78	1 ± 2
Method	Reach& Drag+5	Reach& Drag+7	PutCube InSafe+0	Pick&Lift Cylinder+0	Pick&Lift Star+0	Pick&Lift Moon+0	Pick&Lift Toy+0	PutIn Cupboard+7	PutIn Cupboard+8	
3D-LOTUS (baseline)	95±4.5	18 ± 10.8	25±5.5	69±6.6	93±6.0	80±4.2	66±3.7	$0{\pm}0.0$	$0{\pm}0.0$	
3D-LOTUS++ (baseline)	94±2.0	64±12.4	37±5.1	91±2.0	94±3.7	$29{\pm}6.6$	71±2.0	$0{\pm}0.0$	$0{\pm}0.0$	
MiRA (ours)	42±17.2	45±0	28±11.66	88±9.27	99±2	43±5.1	75±7.07	6±3.74	$0{\pm}0$	

Ħ

#

Method	Avg.	Close Door+0	Close Box+0	Close Fridge2+0	CloseLaptop Lid2+0	Close Microwave2+0	Open Door2+0	Open Box2+0
3D-LOTUS (baseline)	38.1±1.1	0±0.0	58±8.1	36±9.7	54±10.7	85±7.1	42±6.8	11±6.6
3D-LOTUS++ (baseline)	41.5±1.8	1±2.0	$29{\pm}8.6$	93±2.5	50±9.5	99±2.0	52 ± 10.3	16 ± 8.0
MiRA (ours)	41.95±6.07	1±2	9±3.74	92±4	40±10	$99{\pm}2$	57±8.12	23 ± 12.08
Method	Open Drawer2+0	Open Drawer3+0	OpenDrawer Long+0	OpenDrawer Long+1	OpenDrawer Long+2	OpenDrawer Long+3	Toilet SeatUp+0	Open Fridge+0
3D-LOTUS (baseline)	90±3.2	22±8.1	56±13.9	33±11.2	17±8.1	75±6.3	0±0.0	4±5.8
3D-LOTUS++ (baseline)	70±5.5	41±4.9	72±4.0	52 ± 10.8	23±8.1	78±5.1	8 ± 5.1	$0{\pm}0.0$
MiRA (ours)	71±4.9	19±5.83	81±6.63	61±4.9	37±8.12	69±5.83	12±9.8	$0{\pm}0$
Method	OpenLaptop Lid+0	Open Microwave+0	PutMoney InSafe+2	Open Drawer+1	Close Drawer+0	Close Grill+0		
3D-LOTUS (baseline)	100±0.0	0±0.0	0±0.0	0±0.0	87±8.1	29±6.6		
3D-LOTUS++ (baseline)	86±6.6	$0{\pm}0.0$	13±8.1	$0{\pm}0.0$	69±5.8	19±13.9		
MiRA (ours)	75±10	$0{\pm}0$	22±7.48	$0{\pm}0$	72±7.48	41±14.63		

Table 5. Performance on GemBench Level 3.

Table 6. Performance on GemBench Level 4.

Method	Avg.	Push Buttons4+1	Push Buttons4+2	Push Buttons4+3	TakeShoes OutOfBox+0	PutItems InDrawer+0	PutItems InDrawer+2
3D-LOTUS (baseline)	0.3±0.3	3±4.0	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$
3D-LOTUS++ (baseline)	17.4 ± 0.4	76±7.4	$49{\pm}8.6$	37±8.1	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$
MiRA (ours)	15.25±3.56	67±8.12	21±6.63	34±5.83	$0{\pm}0$	$0{\pm}0$	$0{\pm}0$
Method	PutItems InDrawer+4	Tower4+1	Tower4+3	Stack Cups+0	Stack Cups+3	PutAllGroceries InCupboard+0	
3D-LOTUS (baseline)	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$	
3D-LOTUS++ (baseline)	$0{\pm}0.0$	$17{\pm}10.8$	$30{\pm}13.4$	$0{\pm}0.0$	$0{\pm}0.0$	$0{\pm}0.0$	
MiRA (ours)	$0{\pm}0$	23±6.78	38±15.36	$0{\pm}0$	$0{\pm}0$	$0{\pm}0$	