

KALIE: Fine-Tuning Vision-Language Models for Open-World Manipulation without Robot Data

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Building generalist robotic systems involves effectively endowing
2 robots the capabilities to handle novel objects in an open-world setting. Inspired
3 by the advances of large pre-trained models, we propose Keypoint Affordance
4 Learning from Imagined Environments (KALIE), which adapts pre-trained Vi-
5 sion Language Models (VLMs) for robotic control in a scalable manner. Instead
6 of directly producing motor commands, KALIE controls the robot by predict-
7 ing point-based affordance representations based on natural language instructions
8 and visual observations of the scene. The VLM is trained on 2D images with
9 affordances labeled by humans, bypassing the need for training data collected on
10 robotic systems. Through an affordance-aware data synthesis pipeline, KALIE au-
11 tomatically creates massive high-quality training data based on limited example
12 data manually collected by humans. We demonstrate that KALIE can learn to ro-
13 bustly solve new manipulation tasks with unseen objects given only 50 example
14 data points. Compared to baselines using pre-trained VLMs, the our approach
15 consistently achieves superior performances.¹

16 **Keywords:** Vision-Language Model, Data Synthesis, Fine-Tuning, Manipulation

17 1 Introduction

18 The capability to handle an open set of objects, behaviors, and task specifications is essential to
19 the development of generalist robotic systems. Existing learning methods for robotic control can
20 require extensive amounts of data collected on embodied systems [1, 2, 3]. The diversity and quality
21 of the collected data determine the generalization capability that these methods can achieve, which
22 is subject to robotics expertise and manual labors that humans can provide. How can we endow
23 robots with generalizable skills for solving an open set of tasks in a scalable manner?

24 Large pre-trained models offer promising tools with their generalist visual understanding and com-
25 monsense reasoning abilities [4, 5, 6, 7]. Prior works have shown that pre-trained Large Language
26 Models (LLMs) and Vision Language Models (VLMs) can be directly applied to robotics control
27 through prompt engineering in a zero-shot manner [8, 9, 10, 11]. However, despite the capability of
28 generalizing to unseen tasks, such systems often suffer from instability and require significant hard-
29 coded domain knowledge to compensate the pre-trained models’ limited knowledge about robotic
30 control. While fine-tuning on robotic data can mitigate this issue and proves more sample-efficient
31 than training policies from scratch [12, 13], the largest available datasets [14, 15, 16] for robotic
32 control is still far from being comparable to the Internet-scale data used for pre-training the large
33 models with billions of parameters. It remains a grand challenge to effectively employ and adapt
34 pre-trained large models for robotic control.

35 In this paper, we aim to study an alternative solution to this challenge by training large models for
36 manipulation without data collected with on robots. Our key insight is to use visual affordances to

¹Project webpage: <https://sites.google.com/view/kalie-vlm>

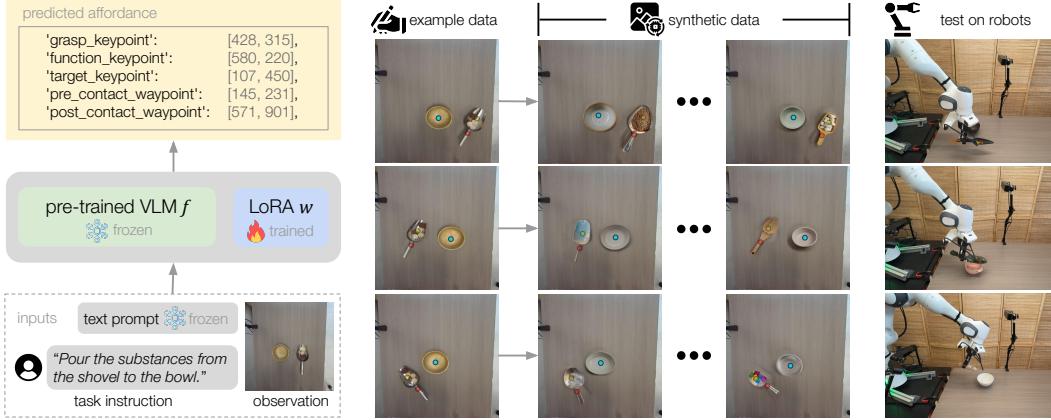


Figure 1: **Overview of KALIE.** By fine-tuning a pre-trained VLM, KALIE predicts the point-based affordance representation given the input task instruction and visual observation. Based on limited example data collected by humans, KALIE generates synthetic data with high diversity while preserving the task semantics and the keypoint annotations. The fine-tuned VLM can robustly generate motions for manipulation tasks with unseen objects and scene arrangements.

37 guide robotic control and leverage the broad knowledge incorporated in large pre-trained models
 38 to efficiently learn to predict the affordances. Built upon point-based affordance representations
 39 defined on 2D images from Manuelli et al. [17] and Fang et al. [10], we fine-tune a VLM on labeled
 40 affordance data. Humans can easily collect such affordance data by randomizing scenes for a target
 41 task (e.g., cleaning a table), taking images through the camera, and annotating the affordances on the
 42 image, which bypasses the need for collecting demonstration trajectories through teleoperation the
 43 robots or hard-coded policies. The main challenges are 1) how to repurpose the VLM pre-trained for
 44 visual-question answering for efficient affordance learning, and 2) how to efficiently utilize human
 45 supervisions to create training data that can cover diverse scenarios.

46 To this end, we propose Keypoint Affordance Learning from Imagined Environments (KALIE). As
 47 shown in Fig. 1, our approach fine-tunes a pre-trained VLM to the point-based affordance represen-
 48 tation given the input task instruction and visual observation. Starting with limited example data
 49 of manually arranged scenes and annotated affordance labels, KALIE automatically creates mas-
 50 sive and diverse synthetic images to scale up the training data. To stay faithful to the task semantics
 51 and keypoint annotations while diversifying the data distribution as much as possible, we propose an
 52 affordance-aware data synthesis pipeline using a pre-trained diffusion model [18, 19] with additional
 53 contexts. We evaluate KALIE on various manipulation tasks involving tool-use, articulated objects,
 54 and deformable objects. KALIE consistently solves target tasks with diverse unseen objects and
 55 initial arrangements and achieves superior performances compared to baselines using pre-trained
 56 VLMs [9, 10].

57 2 Keypoint Affordance Learning from Imagined Environments

58 We propose Keypoint Affordance Learning from Imagined Environments (KALIE) to adapt pre-
 59 trained Vision Language Models (VLMs) to acquire generalizable skills without robot experiences.
 60 In this section, we will first define the affordance prediction problem in the few-shot setting using
 61 point-based affordances labeled by human experts. Next, we will introduce a novel affordance-
 62 aware data synthesis recipe to diversify the training data, which automatically generates massive
 63 high-quality data based on the example data collected only for limited scenarios. Then, we will
 64 describe our VLM fine-tuning approach and discuss the key design options. Lastly, we summarize
 65 the overall system at the end of this section.

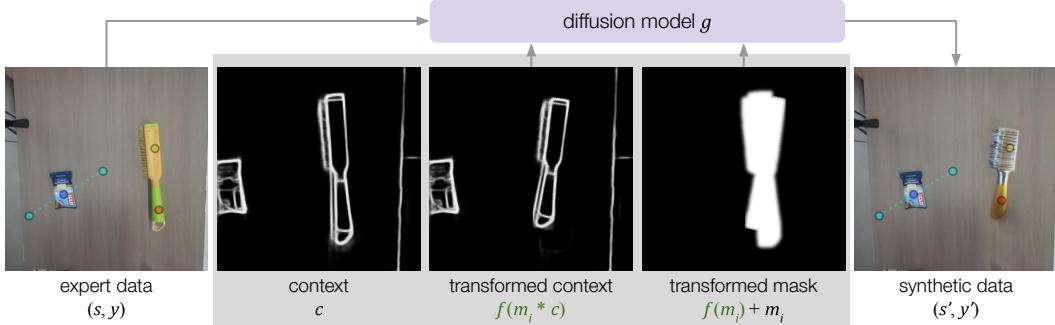


Figure 2: **Affordance-aware data synthesis.** KALIE employs the inpainting capability of a diffusion model to generate synthetic data. To diversify the scenes while staying faithful to the task semantics and the keypoint annotations, KALIE computes and transforms the context such as soft edges as additional inputs guiding the diffusion model.

66 2.1 Problem Statement

67 We consider the problem of open-world robotic manipulation involving unseen objects and initial
 68 arrangements of the scenes. As shown in Fig. 1, each task is single-stage and is specified by a
 69 free-form language description l , such as “*Using the brush to sweep the snack package.*” The robot
 70 observes an RGBD image from a third-person camera and performs a 6-DoF motion trajectory in
 71 open loop to complete the task.

72 To tackle this problem, we employ a VLM f pre-trained on Internet-scale data [6]. Following the
 73 practice of [10], we query the VLM to produce point-based affordance representations to guide a
 74 motion planner to generate motions. The VLM takes as inputs the prompt ρ , the task instruction l ,
 75 and the input image s and predict the affordance representation y as:

$$\hat{y} = f(\rho, l, s), \quad (1)$$

76 where y contains a set of keypoints, including the grasp keypoint, the function keypoint, the target
 77 keypoint, the pre-contact waypoint, and the post-contact waypoint, specified by 2D coordinates on
 78 the image (see Fig. 1). Additional properties such as the height and the orientation of the gripper will
 79 also be decided for each task. Based on the affordance representation, a low-level motion generator
 80 computes a motion trajectory to complete the task. We assume the desired motion to solve each
 81 task can be specified by the same subset of these points (e.g., the sweeping with a brush requires
 82 all five points, and drawer closing requires everything but the grasp point, as shown in Section 3.3),
 83 but the specific coordinates of these points depend on the objects and their poses with respect to the
 84 robot. More detail about definition of the affordance representation and the prompt design will be
 85 elaborated in the Appendix.

86 In this work, we consider tasks and objects which are challenging for VLMs to handle in zero-shot.
 87 In contrast to [10], we consider a few-shot learning setting, in which we fine-tune the pre-trained
 88 VLM to acquire and improve skills for unseen scenarios on limited and non-robot data. Notably,
 89 the point-based representations enable us to outsource motion generation to the low-level motion
 90 planner and focus on only the affordance prediction problem. Therefore, training only requires
 91 data of pairs of observed image s and the ground truth keypoints y collected by human experts or
 92 generated automatically as described below. We assume access to an example dataset \mathcal{D} containing a
 93 limited number of (s, y) pairs for each target task. In our experiments, we assume the most extreme
 94 case, in which the data is collected on a *single set of* objects for each task and the fine-tuned model
 95 is evaluated on unseen object sets.

96 2.2 Affordance-Aware Data Synthesis

97 To scale up the training of the VLM to enhance its generalization capabilities to unseen scenarios,
 98 we automatically synthesize a training dataset \mathcal{D}' to cover a wide range of environments. Each of

99 the new data points $(s', y') \in \mathcal{D}'$ is synthesized by modifying an existing data point $(s, y) \in \mathcal{D}$
 100 collected by human experts.
 101 Following the practice of [20], we compute the segmentation mask of the object in the scene using
 102 open-vocabulary segmentation [21] and then inpaint the masked region with a diffusion model [18].
 103 However, naively inpainting the masked region can lead to undesired results. Without a direct mech-
 104 anism to specify the geometric properties of the object, it would be hard to diversify the inpainted
 105 images in a way that can cover the desired distribution of testing scenarios. Moreover, there is
 106 usually a discrepancy between the appearance of the inpainted object and the original keypoint an-
 107 notation, introducing the need of manually re-labeling keypoints. To generate massive, high-quality
 108 data without additional manual labor, we need to ensure that the generated images stay faithful to
 109 the context of the target task and the annotated keypoints, while aggressively diversifying the envi-
 110 ronments as much as possible.
 111 To tackle this challenge, we design an affordance-aware data synthesis pipeline by leveraging addi-
 112 tional context images to guide the generation process, as shown in Fig. 2. Specifically, we employ
 113 the ControlNet [19] diffusion model g , which takes inputs as the input image s , the segmentation
 114 mask m , a context image c , and a language description of the object o and generates the new image
 115 as $s' \sim g(\cdot|s, m, c, o)$. We would like c to be a compact representation of the object’s geometric
 116 properties that provide clues about the affordances, while minimizing other detailed visual informa-
 117 tion to leave enough free reign for the diffusion model. By inpainting an image s' in accordance with
 118 c , we hope to obtain new objects of the same point-based affordances. In this work, we choose to
 119 use a soft edge map as c computed by an external image processing algorithm [22], which outlines
 120 the contours and parts of the object.
 121 To cover testing objects of unseen shapes, we introduce additional randomization operations to the
 122 geometry of the synthesized object. Directly modifying either the object’s appearance in the pixel
 123 space can easily affect other parts of the image and create artifacts. Instead, we propose to apply
 124 transformation $h(\cdot)$ on the compact context c before calling the diffusion model G to inpaint the
 125 image. The transformation function $h(\cdot)$ can include basic data augmentation operations such as
 126 random scaling, translation, and rotation, as well as additional operations such as elastic distortion.
 127 To modify the context c , we transform the region under the masked as $h(m * c)$. The region to be
 128 inpaint on the original image now becomes $m + h(m)$, which includes the transformed silhouette of
 129 the object $h(m)$ in addition to the white space left by removing the original object m . Accordingly,
 130 we also apply the same transformation to the annotated keypoints y with the context image to keep
 131 them consistent. Therefore, the sampling process with the transformation can be written as:

$$\begin{aligned} s' &\sim g(\cdot|s, h(m) + m, h(m * c), o) \\ y' &= h(y), \end{aligned} \tag{2}$$

132 where we slightly overload the notation by using $h(\cdot)$ to denote the transformation applied to both
 133 the images and the keypoint coordinates.

134 2.3 Efficient Adaptation for Keypoint Affordance Prediction

135 We fine-tune the VLM f to predict point-based affordance representations. To adapt f , which is
 136 pre-trained for visual-question answering by predicting tokens, we need to make our design choices
 137 around two considerations. First, how to represent the point-based affordances so that we can effec-
 138 tively re-purpose or modify the VLM’s prediction head for affordance prediction. Second, how to
 139 perform sample-efficient fine-tuning to utilize the pre-trained VLM’s pre-trained capabilities while
 140 endowing it with the additional knowledge from the new dataset \mathcal{D} .
 141 We investigate two design options for the prediction head: 1) **Regression Head** [23] adds an addi-
 142 tional linear layer on top of the last hidden state of the VLM to directly predict the $x - y$ coordinates
 143 of the keypoint affordances. In input to the VLM is the image along with appropriate task instruc-
 144 tions and the last hidden state of the last token is used. 2) **Natural Language Affordance Prediction**
 145 fine-tunes the VLM to output a well-formatted natural language that includes text-based keypoint

146 affordances with an example shown in Fig. 1. Each keypoint affordance is represented by the $x - y$
 147 coordinates as integers normalized between a predefined range.
 148 These design options resort to the pre-trained weights and the new dataset in different manners.
 149 Empirically as shown in Fig. 3, we found those two design choices achieve similar performances.
 150 As Natural Language Affordance Prediction aligns more closely with other applications of VLMs
 151 such as Visual Question Answering, we choose to use this design option as our main method. During
 152 training, we convert the ground truth affordance label y into the corresponding format to compute the
 153 losses. For both design options, we use Low-Rank Adaptation (LoRA) [24] to fine-tune the VLM.
 154 We use a L2 regression loss for Regression Head and the cross entropy loss for Natural Language
 155 Affordance Prediction.

156 2.4 System Summary

157 The overall pipeline of KALIE is summarized in Algorithm 1. Starting with the original example
 158 dataset \mathcal{D} , KALIE automatically create \mathcal{D}' to fine-tune the VLM f . The data synthesis pipeline
 159 operates in an object-centric manner, iteratively process and inpaint each of the M objects in the
 160 scene as explained in Section 2.2, using each object’s segmentation mask m_i and description o_i .
 161 During this process, a VLM (which can be the same with f or a different model) is used to sample
 162 an alternative description of the object to guide the inpainting process of the diffusion model g . The
 163 generated data is combined with the example dataset to adapt f by optimizing the weights w as
 164 explained in Section 2.3. More implementation details will be explained in the Appendix.

Algorithm 1 Keypoint Affordance Learning from Imagined Environments (KALIE)

Inputs: Pre-trained VLM f , task instruction l , example dataset D , desired dataset size N .

```

1:  $D' \leftarrow \emptyset$ 
2: while  $|D'| < N$  do
3:   Sample  $(s, y)$  from  $D$ .
4:   Compute the context image  $c$  given  $s$ .
5:   Extract the object descriptions  $\{o_i\}_{i=1}^M$  for each object  $i$  given  $s$  and  $l$ .
6:    $s'_0 \leftarrow s$ 
7:   for  $i = 1, \dots, M$  do ( Section 2.2)
8:     Compute the segmentation mask  $m_i$  for the  $i$ -th object given  $s$  and  $o_i$ .
9:     Sample a new description  $o'_i$  given  $o_i$ .
10:    Transform the mask, the context as  $h(m_i)$  and  $h(m_i * c)$ .
11:    Transform the keypoints  $y_i$  on the  $i$ -th object as  $h(y_i)$ .
12:    Inpaint the image as  $s'_i \sim G(\cdot | s'_{i-1}, h(m_i) + m_i, h(m_i * c), o'_i)$ .
13:   end for
14:   Merge the transformed keypoints into  $y'$ .
15:   if  $s'_M$  pass the filter: then
16:      $D' \leftarrow D' \cup \{(s'_M, y')\}$ 
17:   end if
18: end while
19: Fine-tune  $f$  on  $D \cup D'$  by optimizing the weights  $w$  ( Section 2.3).
  
```

165 3 Experiments

166 We design our experiments to investigate the following questions: 1) Can KALIE synthesize diverse
 167 and high-quality data for affordance prediction? 2) Can KALIE fine-tune pre-trained VLMs to
 168 improve their performances on challenging manipulation tasks with unseen objects?

169 3.1 Experimental Setup

170 **Environments.** Our experiments are conducted in a real-world table-top manipulation environment
 171 with a 7-DoF robot arm and a parallel jaw gripper. A top-down RGBD camera is used to receive
 172 visual observation of the environment.

173 **Tasks.** We design various table-top manipulation tasks with a diverse set of daily objects

174 **Training and evaluation protocols.** We fine-tune CogVLM [6], an open-source VLM of 17 billion
175 parameters. To fine-tune the VLM, 50 input images are collected by humans through randomly
176 setting up the scenes based on the context of the task and annotating the affordance representations
177 through a graphical user interface. Though data synthesis with diffusion models, we generate an
178 additional set of 500 pairs of images and keypoints. For each target task, test the fine-tuned model
179 on 15 trials with three sets of unseen objects.

180 **Baselines.** We compare KALIE with two baselines that employ VLMs for robotic manipulation,
181 **VoxPoser** [9] and **MOKA** [10]. In contrast to KALIE , which pre-trains a much smaller open-
182 source VLM, the two baselines use the OpenAI GPT-4V[4]. We use the collected example data as
183 the in-context examples in MOKA.

Methods	Table Sweeping	Drawer Closing	Towel Hanging	Trowel Pouring	USB Unplugging
VoxPoser [9]	3/15	8/15	1/15	0/15	0/15
MOKA [10]	9/15	9/15	5/15	7/15	2/15
KALIE (Ours)	13/15	15/15	12/15	11/15	9/15

Table 1: **Task success rates.** We evaluate KALIE on five manipulation tasks involving tool use, deformable objects, and articulated objects. KALIE robustly solve these tasks and consistently achieves superior performances compared to baselines.

184

3.2 Quantitative Results

185 **Comparative results.** We evaluate KALIE across 5 manipulation tasks involving tool objects, articulated
186 objects, and deformable objects. Compared to the baseline methods using pre-trained VLMs,
187 KALIE consistently achisves higher success rates. Baseline methods are often subject to critical
188 failures including perception errors, incorrect output formats, and unsuitable predicted affordances.
189 The VLM fine-tuned by KALIE demonstrate to robustly solve the tasks without these issues while
190 most failures are due to inaccurate point predictions.

191 **Ablative study.** In Fig. 3, we present some ablation experiment results examining the effectiveness
192 of Imagined Environments and our choice of Natural Language Affordance Prediction against Re-
193 gression Head. We report the test Mean Square Error (MSE) on a held-out set with unseen objects
194 and ground-truth annotations. The coordinates are normalized between 0 and 999 with respect to
195 the actual height and width of the image. The train set consists of 50 real images collected and 500
196 synthesized images, and the test set consists of 50 real images of novel objects excluded from the
197 train set. We observe that Imagined Environments significantly improve the test MSE and ours with
198 Natural Language Affordance Prediction works comparably with Regression Head, while staying
199 consistent with the general natural language interface of VLMs for vision-language tasks.

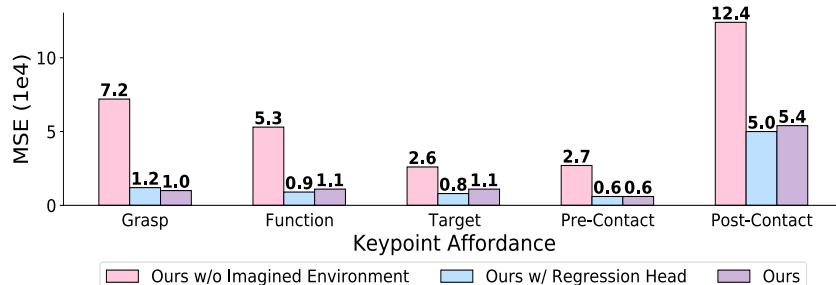


Figure 3: **Ablative study.** The ablation experiments are carried out on the sweeping task. Mean Square Error (MSE) for each keypoint affordance on the a test set of novel objects is reported.

200 **3.3 Qualitative Results**

201 In Fig. 4, we show examples of the synthetic images. Each row contains the data for a task. The
 202 first row is the example data and the remaining rows are the synthetic data. KALIE demonstrates to
 203 generate diverse samples while preserving the task context and the original affordance labels.

204 We show examples of task execution in Fig. 5. The VLM fine-tuned by KALIE successfully predict
 205 the affordance and solve each of the tasks.



Figure 4: **Data synthesis examples.** We show the example data and synthetic data generated by KALIE for each of the five tasks. Each row contains the data for a task. The first row is the example data and the remaining rows are the synthetic data.

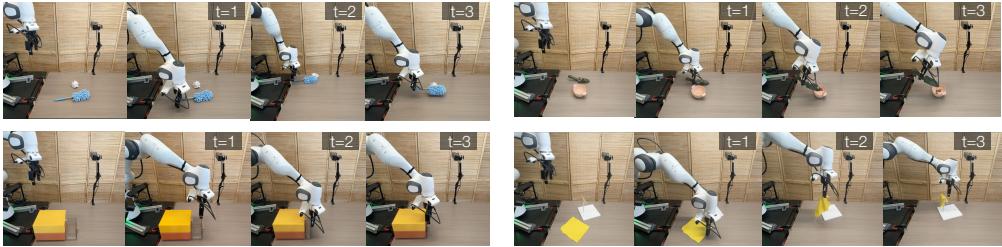


Figure 5: **Affordance prediction and task execution.** We show examples of the predicted affordance representations and the key frames of task execution. Across various manipulation tasks, KALIE demonstrates to robustly solve the task.

206 **4 Related Work**

207 **Adaptation of vision language models.** Due to their versatile nature, Vision Language Models
 208 (VLMs) can be effectively fine-tuned to accommodate a variety of downstream tasks [25, 26, 27, 28].

209 Specifically, to facilitate the prediction of spatially grounded outputs, previous research has
210 investigated methods such as sets of marks [29], scaffolding [4], and coordinate-based bounding
211 boxes [30, 31, 6]. Of these methods, coordinate-based references are particularly adaptable, ca-
212 pable of pinpointing any location within an image. Thus, we employ this approach for predicting
213 keypoints. Our contribution does not lie in proposing a new method for fine-tuning VLMs, but rather
214 in integrating and adapting VLM fine-tuning for robotic control.

215 **Robotic control with large models.** An increasing number of works have employed foundation
216 models in robotics through pre-training or fine-tuning on robot data manually collected through
217 teleoperation or scripted policies [12, 32, 13, 33]. Due to the lack of datasets that can cover the
218 vast complexity and diversity of robotics applications, most of the approaches focus on specific task
219 categories such as grasping and object rearrangements. The generalization capability of the trained
220 model is also constrained to the distribution of objects and environments covered by the manually
221 collected datasets. Alternatively, other works have attempted to combine and prompt pre-trained
222 models to solve unseen tasks in zero-shot manners [8, 9, 10, 11, 34]. However, the performance
223 of these works is usually subject to the capability of the pre-trained models, as well as non-trivial
224 expert knowledge and manual labor to design prompts and in-context examples. In contrast to these
225 approaches, the proposed method fine-tunes a VLM to robustly solve the target tasks. To avoid
226 the need of collecting extensive amount of robot data, our model train the VLM to predict the point-
227 based affordance representation from [10] and employ a diffusion model to automatically synthesize
228 massive, high-quality data.

229 **Data synthesis in robotics.** To alleviate the data bottleneck in robot learning, prior work has
230 investigated various approaches to augment and synthesize data. Data augmentation, especially
231 random image transformations, have been widely used to improve generalization to unseen visual
232 inputs [35, 36]. These random operations can effectively improve the models generalization capabili-
233 ties to unseen visual inputs during the test time, but cannot extend the model’s capabilities beyond
234 the coverage of the training distribution. Domain randomization has also been broadly used to train
235 robust models for robotic control [37, 38, 39, 40] when collecting simulated robot experiences. In
236 contrast, our method directly diversify the training data without the need for a physical simula-
237 tor. Recent works in computer vision and robotics have leveraged deep generative models, such as
238 diffusion models, to synthesize unseen environments by leveraging broad knowledge learned from
239 Internet-scale images [41, 42, 43, 44, 20, 45, 46, 47, 48]. While the prior work can easily generate
240 defected samples, we propose an affordance-aware data synthesis approach that enable the diffusion
241 model to generate diverse data with much higher quality and consistent affordance annotations.

242 5 Conclusion

243 In this work, we propose KALIE to fine-tune pre-trained VLMs to predict affordances for robotic
244 manipulation with open sets of objects and initial arrangements of scenes. Using the proposed
245 affordance-aware data synthesis pipeline, we generate massive high-quality data to scale up training
246 without extensive manual labors or domain expertise with robots. We fine-tune the VLM on the
247 generated data to predict the point-based affordance representations. Across various manipulation
248 tasks involving tool use, deformable objects, and articulated objects, we demonstrate KALIE can
249 robustly solve the task and consistently outperform baselines. We hope KALIE will inspire future
250 research towards adapting vision-language models for open-world robotic control.

251 **Limitations and future work.** The current affordance representation in KALIE is still limited to
252 single-arm table-top manipulation. To apply KALIE to more complicated scenarios, such as dy-
253 namic manipulation and whole-body control, we would need to extend the design of the affordance
254 representation. In addition, there is still a large discrepancy between the open-source VLM fine-
255 tuned by KALIE and the state-of-the-art VLMs without accessible fine-tuning APIs. In the future,
256 we hope to apply KALIE to more powerful VLMs to further improve the performance.

257 **References**

- 258 [1] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen. Learning hand-eye coordination for
259 robotic grasping with deep learning and large-scale data collection. *The International Journal
260 of Robotics Research*, 37:421 – 436, 2016. URL <https://api.semanticscholar.org/CorpusID:13072941>.
- 262 [2] L. Pinto and A. K. Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and
263 700 robot hours. *2016 IEEE International Conference on Robotics and Automation (ICRA)*,
264 pages 3406–3413, 2015. URL <https://api.semanticscholar.org/CorpusID:3177253>.
- 265 [3] A. Brohan et al. Rt-1: Robotics transformer for real-world control at scale, 2023.
- 266 [4] OpenAI. Gpt-4 technical report, 2024.
- 267 [5] G. Team. Gemini: A family of highly capable multimodal models, 2024.
- 268 [6] W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song, J. Xu,
269 B. Xu, J. Li, Y. Dong, M. Ding, and J. Tang. Cogvlm: Visual expert for pretrained language
270 models, 2024.
- 271 [7] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny. Minigpt-4: Enhancing vision-language
272 understanding with advanced large language models, 2023.
- 273 [8] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as
274 policies: Language model programs for embodied control. In *IEEE International Conference
275 on Robotics and Automation*, pages 9493–9500. IEEE, 2023.
- 276 [9] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. Voxposer: Composable 3d value
277 maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023.
- 278 [10] K. Fang, F. Liu, P. Abbeel, and S. Levine. MOKA: Open-Vocabulary Robotic Manipulation
279 through Mark-Based Visual Prompting. In *Conference on Robot Learning (CoRL)*, 2024.
- 280 [11] D. Shah, B. Osinski, B. Ichter, and S. Levine. Lm-nav: Robotic navigation with large pre-
281 trained models of language, vision, and action. In *Conference on Robot Learning*, 2022. URL
282 <https://api.semanticscholar.org/CorpusID:250426345>.
- 283 [12] A. Brohan et al. Rt-2: Vision-language-action models transfer web knowledge to robotic
284 control, 2023.
- 285 [13] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna,
286 C. Xu, J. Luo, T. Kreiman, Y. Tan, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source
287 generalist robot policy. <https://octo-models.github.io>, 2023.
- 288 [14] E. Collaboration. Open x-embodiment: Robotic learning datasets and rt-x models, 2024.
- 289 [15] A. Khazatsky et al. Droid: A large-scale in-the-wild robot manipulation dataset, 2024.
- 290 [16] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch,
291 Q. Vuong, A. He, V. Myers, K. Fang, C. Finn, and S. Levine. Bridgedata v2: A dataset
292 for robot learning at scale, 2024.
- 293 [17] L. Manuelli, W. Gao, P. R. Florence, and R. Tedrake. kpam: Keypoint affordances for category-
294 level robotic manipulation. In *International Symposium of Robotics Research*, 2019. URL
295 <https://api.semanticscholar.org/CorpusID:80628296>.
- 296 [18] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image syn-
297 thesis with latent diffusion models, 2022.

- 298 [19] L. Zhang, A. Rao, and M. Agrawala. Adding conditional control to text-to-image diffusion
299 models, 2023.
- 300 [20] T. Yu, T. Xiao, A. Stone, J. Tompson, A. Brohan, S. Wang, J. Singh, C. Tan, M. Dee, J. Per-
301 alta, B. Ichter, K. Hausman, and F. Xia. Scaling robot learning with semantically imagined
302 experience. *ArXiv*, abs/2302.11550, 2023. URL <https://api.semanticscholar.org/CorpusID:257079001>.
- 304 [21] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
305 A. C. Berg, W.-Y. Lo, P. Dollár, and R. B. Girshick. Segment anything. 2023 IEEE/CVF
306 International Conference on Computer Vision (ICCV), pages 3992–4003, 2023. URL <https://api.semanticscholar.org/CorpusID:257952310>.
- 308 [22] X. Soria, Y. Li, M. Rouhani, and A. D. Sappa. Tiny and efficient model for the edge detection
309 generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision
(ICCV) Workshops*, pages 1364–1373, October 2023.
- 311 [23] W. Chen, O. Mees, A. Kumar, and S. Levine. Vision-language models provide promptable
312 representations for reinforcement learning, 2024.
- 313 [24] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, and L. Zhang.
314 Grounding dino: Marrying dino with grounded pre-training for open-set object detection, 2023.
- 315 [25] W. Hong, W. Wang, Q. Lv, J. Xu, W. Yu, J. Ji, Y. Wang, Z. Wang, Y. Zhang, J. Li, B. Xu,
316 Y. Dong, M. Ding, and J. Tang. Cogagent: A visual language model for gui agents, 2023.
- 317 [26] Y. Zhai, H. Bai, Z. Lin, J. Pan, S. Tong, Y. Zhou, A. Suhr, S. Xie, Y. LeCun, Y. Ma, and
318 S. Levine. Fine-tuning large vision-language models as decision-making agents via reinforce-
319 ment learning, 2024.
- 320 [27] B. Chen, C. Shu, E. Shareghi, N. Collier, K. Narasimhan, and S. Yao. Fireact: Toward language
321 agent fine-tuning, 2023.
- 322 [28] A. Zeng, M. Liu, R. Lu, B. Wang, X. Liu, Y. Dong, and J. Tang. Agenttuning: Enabling
323 generalized agent abilities for llms, 2023.
- 324 [29] J. Yang, H. Zhang, F. Li, X. Zou, C. Li, and J. Gao. Set-of-mark prompting unleashes extraor-
325 dinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*, 2023.
- 326 [30] P. Wang, S. Wang, J. Lin, S. Bai, X. Zhou, J. Zhou, X. Wang, and C. Zhou. One-peace:
327 Exploring one general representation model toward unlimited modalities, 2023.
- 328 [31] K. Chen, Z. Zhang, W. Zeng, R. Zhang, F. Zhu, and R. Zhao. Shikra: Unleashing multimodal
329 llm’s referential dialogue magic, 2023.
- 330 [32] D. Driess, F. Xia, M. S. M. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tomp-
331 son, Q. H. Vuong, T. Yu, W. Huang, Y. Chebotar, P. Sermanet, D. Duckworth, S. Levine,
332 V. Vanhoucke, K. Hausman, M. Toussaint, K. Greff, A. Zeng, I. Mordatch, and P. R. Florence.
333 Palm-e: An embodied multimodal language model. In *International Conference on Machine
Learning*, 2023. URL <https://api.semanticscholar.org/CorpusID:257364842>.
- 335 [33] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipu-
336 lation. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.
- 337 [34] Y. Hu, F. Lin, T. Zhang, L. Yi, and Y. Gao. Look before you leap: Unveiling the power of
338 gpt-4v in robotic vision-language planning. *arXiv preprint arXiv:2311.17842*, 2023.
- 339 [35] I. Kostrikov, D. Yarats, and R. Fergus. Image augmentation is all you need: Regularizing
340 deep reinforcement learning from pixels. *ArXiv*, abs/2004.13649, 2020. URL <https://api.semanticscholar.org/CorpusID:216562627>.

- 342 [36] M. Laskin, K. Lee, A. Stooke, L. Pinto, P. Abbeel, and A. Srinivas. Reinforcement learning
343 with augmented data. *Advances in neural information processing systems*, 33:19884–19895,
344 2020.
- 345 [37] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization
346 for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ*
347 *international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.
- 348 [38] B. Mehta, M. Diaz, F. Golemo, C. J. Pal, and L. Paull. Active domain randomization. In
349 *Conference on Robot Learning*, pages 1162–1176. PMLR, 2020.
- 350 [39] I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino,
351 M. Plappert, G. Powell, R. Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint*
352 *arXiv:1910.07113*, 2019.
- 353 [40] H. Qi, B. Yi, S. Suresh, M. Lambeta, Y. Ma, R. Calandra, and J. Malik. General in-hand object
354 rotation with vision and touch. In *Conference on Robot Learning*, pages 2549–2564. PMLR,
355 2023.
- 356 [41] A. Kar, A. Prakash, M.-Y. Liu, E. Cameracci, J. Yuan, M. Rusiniak, D. Acuna, A. Torralba, and
357 S. Fidler. Meta-sim: Learning to generate synthetic datasets. In *Proceedings of the IEEE/CVF*
358 *International Conference on Computer Vision*, pages 4551–4560, 2019.
- 359 [42] Y. Zhang, H. Ling, J. Gao, K. Yin, J.-F. Lafleche, A. Barriuso, A. Torralba, and S. Fidler.
360 Datasetgan: Efficient labeled data factory with minimal human effort. In *Proceedings of*
361 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10145–10155,
362 2021.
- 363 [43] Y. Wang, Z. Xian, F. Chen, T.-H. Wang, Y. Wang, K. Fragkiadaki, Z. Erickson, D. Held, and
364 C. Gan. Robogen: Towards unleashing infinite data for automated robot learning via generative
365 simulation. *arXiv preprint arXiv:2311.01455*, 2023.
- 366 [44] M. Yang, Y. Du, K. Ghasemipour, J. Tompson, D. Schuurmans, and P. Abbeel. Learning
367 interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 2023.
- 368 [45] J. Devarajan, A. Kar, and S. Fidler. Meta-sim2: Unsupervised learning of scene structure
369 for synthetic data generation. In *Computer Vision–ECCV 2020: 16th European Conference,*
370 *Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII 16*, pages 715–733. Springer, 2020.
- 371 [46] E. Wood, T. Baltrušaitis, C. Hewitt, S. Dziadzio, T. J. Cashman, and J. Shotton. Fake it till you
372 make it: face analysis in the wild using synthetic data alone. In *Proceedings of the IEEE/CVF*
373 *international conference on computer vision*, pages 3681–3691, 2021.
- 374 [47] S. W. Kim, B. Brown, K. Yin, K. Kreis, K. Schwarz, D. Li, R. Rombach, A. Torralba, and
375 S. Fidler. Neuralfield-lmd: Scene generation with hierarchical latent diffusion models. In
376 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
377 8496–8506, 2023.
- 378 [48] D. Li, H. Ling, S. W. Kim, K. Kreis, S. Fidler, and A. Torralba. Bigdatasetgan: Synthesizing
379 imangenet with pixel-wise annotations. In *Proceedings of the IEEE/CVF Conference on*
380 *Computer Vision and Pattern Recognition*, pages 21330–21340, 2022.