

---

# Zero-Shot Robotic Manipulation with Pre-Trained Image-Editing Diffusion Models

---

Anonymous Author(s)

## Abstract

1 If generalist robots are to operate in truly unstructured environments, they need  
2 to be able to recognize and reason about novel objects and scenarios. Such ob-  
3 jects and scenarios might not be present in the robot’s own training data. We  
4 propose SuSIE, a method that leverages an image editing diffusion model to act  
5 as a high-level planner by proposing intermediate subgoals that a low-level con-  
6 troller attains. Specifically, we fine-tune InstructPix2Pix (Brooks et al., 2023)  
7 on robot data such that it outputs a hypothetical future observation given the  
8 robot’s current observation and a language command. We then use the same robot  
9 data to train a low-level goal-conditioned policy to reach a given image obser-  
10 vation. We find that when these components are combined, the resulting system  
11 exhibits robust generalization capabilities. The high-level planner utilizes its  
12 Internet-scale pre-training and visual understanding to guide the low-level goal-  
13 conditioned policy, achieving significantly better generalization than conventional  
14 language-conditioned policies. We demonstrate that this approach solves real  
15 robot control tasks involving novel objects, distractors, and even environments,  
16 both in the real world and in simulation. The project website can be found at  
17 <http://subgoal-image-editing.github.io>.

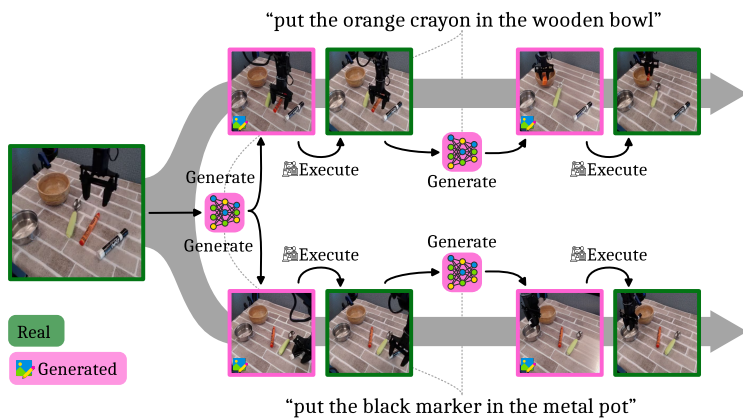


Figure 1: SuSIE leverages a pre-trained image editing model to generate future image subgoals based on a language commands. A low-level goal-reaching policy then executes the actions needed to reach each subgoal. Alternating this loop enables us to solve the task.

## 1 Introduction

19 A useful generalist robot must be able to — much like a person — recognize and reason about novel  
20 objects and scenarios it has never encountered before. For example, if a user instructs the robot to  
21 “hand me that jumbo orange crayon,” it ought to be able to do so even if it has never interacted

22 with a jumbo orange crayon before. In other words, the robot needs to possess not only the physical  
 23 capability to manipulate an object of that shape and size but also the semantic understanding to  
 24 reason about an object outside of its training distribution. As much as robotic manipulation datasets  
 25 have grown in recent years, it is unlikely that they will ever include every conceivable instance  
 26 of objects and settings, any more so than the life experiences of a person ever include physical  
 27 interactions with every type of object. While these datasets contain more than enough examples of  
 28 manipulating elongated cylindrical objects, they lack the broad semantic knowledge necessary to  
 29 ground the *particular* objects that robots will undoubtedly encounter during everyday operation.

30 In this paper, we develop an approach for leveraging a class of pre-trained image-editing models  
 31 (e.g., InstructPix2Pix (Brooks et al., 2023)) for improving motor control and policy execution. Our  
 32 key insight is to decompose the robotic control problem into two phases: first, synthesizing a “sub-  
 33 goal” that the robot must reach to complete the user-specified task, and then, attempting to reach  
 34 this subgoal via a goal-reaching robot controller. The first phase of this recipe incorporates semantic  
 35 information by fine-tuning an image-editing model on robot data such that, given the robot’s current  
 36 observation and a natural language command, the model generates a hypothetical *future* subgoal  
 37 that allows the robot to complete the command. We then employ a low-level goal-reaching policy to  
 38 reach this hypothetical future subgoal.

## 39 2 SuSIE: Subgoal Synthesis via Image Editing.

40 Our goal is to utilize semantic information from the Internet to improve language-guided robot  
 41 control in novel environments, scenes, and objects. How can we do this when models trained on  
 42 general-purpose Internet data do not provide guidance in selecting low-level actions? Our key insight  
 43 is that we can still utilize some sort of a pre-trained model for guiding low-level control if we could  
 44 decouple the robot control problem into two phases: (i) imagining subgoals that would need to  
 45 be attained to succeed at the task, and (ii) learning low-level control policies for reaching these  
 46 generated subgoals. Our method incorporates semantic information from Internet pre-training in  
 47 phase (i), by fine-tuning a text-guided image-editing model for subgoal generation. Phase (ii) is  
 48 accomplished via a goal-conditioned policy trained only on robot data. We describe each of these  
 49 phases in Section D in Appendix.

### 50 2.1 Phase (i): Synthesizing Subgoals From Image Editing Models

51 The primary component of our method is a generative model that, given a target task specified in  
 52 natural language, can guide the low-level controller towards a state that it must try to attain in order  
 53 to solve the task. One way to accomplish this is to train a generative model to produce an immediate  
 54 next way-point or subgoal frame. We can then incorporate semantic information from the Internet  
 55 into our algorithm by initializing this generative model with a suitable pre-trained initialization,  
 56 followed by fine-tuning it on multi-task, diverse robot data.

57 What is a good pre-trained initialization for initializing this model? Our intuition is that since ac-  
 58 complishing a task is equivalent to “editing” the pixels of an image of the robot workspace under  
 59 controls prescribed by the language command, a favorable pre-trained initialization is provided by  
 60 a language-guided image-editing model. We instantiate our approach with Instruct pix2pix (Brooks  
 61 et al., 2023), though other image editing models could also be used. Formally, this model is given  
 62 by  $p_\theta(\mathbf{s}_{\text{edited}}|\mathbf{s}_{\text{orig}}, l)$ . Then, using the dataset  $\mathcal{D}$  of robot trajectories, we fine-tune  $p_\theta$  on tuples  
 63 containing a pair of images sampled from a trajectory and the corresponding language annotation:  
 64  $(\mathbf{s}_{\text{orig}} := \mathbf{s}_i^k, \mathbf{s}_{\text{edited}} := \mathbf{s}_j^k, l_k)$ , where  $\mathbf{s}_j$  is a state that appears after  $\mathbf{s}_i$  ( $j > i$ ). During fine-tuning,  
 65 we run gradient descent on the following objective, starting from  $\theta_0 := \theta_{\text{pre-trained}}$ :

$$\min_{\theta} - \mathbb{E}_{(\tau^k, l_k) \sim \mathcal{D}; \mathbf{s}_i^k \sim \tau^k; j \sim q(j|i)} [\log p_\theta(\mathbf{s}_j^k | \mathbf{s}_i^k, l_k)]. \quad (1)$$

66 We need to choose the distribution  $q$  over the time-step  $j$  given a state  $\mathbf{s}_i^k$  for fine-tuning the image-  
 67 editing model as in Equation 3. Since we model the next subgoal that the low-level controller should  
 68 attain, and since the depending upon the task, this subgoal could be arbitrarily close to the original  
 69 state  $\mathbf{s}_i$ , we require valid tuples  $(\mathbf{s}_i, \mathbf{s}_j, l)$  used for fine-tuning  $p_\theta$  in Equation 3 to have values of  $j$  in  
 70 a bounded interval around  $i$ , specifically we choose  $j \in [i, i + k]$ , where  $k$  is a fixed hyperparameter.  
 71

### 72 2.2 Phase (ii): Reaching Generated Sub-Goals with Goal-Conditioned Policies

73 In order to utilize the fine-tuned image-editing model to actually control the robot, we further need  
 74 to train a low-level controller to actually select suitable robot actions. In this section, we present

75 the design of our low-level controller, followed by a full description our test-time control procedure.  
 76 Since the image-editing model in SuSIE produces images of future subgoals conditioned on natural  
 77 language task descriptions, our low-level controller can simply be a language-agnostic goal-reaching  
 78 policy that aims to reach these generated subgoals.

79 **Training a goal-reaching policy.** Our goal-reaching policy is parameterized as  $\pi_\phi(\cdot|s_i, s_j)$ , where  
 80  $s_j$  is a future frame that the policy intends to reach, by acting at  $s_i$ . At test time, we only need the  
 81 low-level goal-conditioned policy to be proficient at reaching close-by states that lie within  $k$  steps  
 82 of a given state since the image editing model from phase (i) is also trained to produce subgoals  
 83 within  $k$  steps of any state. To train this policy, we run goal-conditioned behavioral cloning (GCBC)  
 84 on the robot data, utilized previously in phase (i). In addition, we can also leverage robot data  $\mathcal{D}'$   
 85 that does not contain language annotations. Formally, our training objective is given by:

$$\max_{\phi} \mathbb{E}_{\tau^i \sim \mathcal{D} \cup \mathcal{D}'; (s_i^k, a_i^k) \sim \tau^k; j \sim q(j|i)} [\log \pi_\phi(a_i^k | s_i^k, s_j^k)], \quad (2)$$

86 where  $q(j|i)$  is the distribution over future frames that we previously utilized in Equation 3.

87 **Test-time control with  $\pi_\phi$  and  $p_\theta$ .** Once both the goal-reaching policy  $\pi_\phi$  and the image editing  
 88 subgoal generation model  $p_\theta$  are trained, we utilize them together to solve new manipulation  
 89 tasks based on user-specified natural language commands. Given a new scene,  $s_0^{\text{test}}$ , and a lan-  
 90 guage task description  $l^{\text{test}}$ , SuSIE attempts to solve the task by iteratively generating subgoals  
 91 and commanding the low-level goal-reaching policy with these subgoals. At the start, we sample  
 92 the first subgoal  $\hat{s}_+^{\text{test}} \sim p_\theta(\cdot | s_0^{\text{test}}, l^{\text{test}})$ . Once the subgoal is generated, we then roll out the goal-  
 93 reaching policy  $\pi_\phi$ , conditioned on  $\hat{s}_+^{\text{test}}$ , for  $k$  time-steps, such that each action is chosen according  
 94 to  $a_j^{\text{test}} \sim \pi_\phi(\cdot | s_j^{\text{test}}, \hat{s}_+^{\text{test}})$ . After  $k$  time steps, given the current image  $s_k^{\text{test}}$ , we refresh the subgoal by  
 95 sampling from the image-editing model again and repeat the process. Note crucially that this recipe  
 96 does not require that the subgoal  $s^{\text{test}}$  be attained after  $k$  steps, as the generative model effectively  
 97 “replans” a new subgoal based on the current observation. Overall, at test time, we alternate between  
 98 obtaining a new subgoal from  $p_\theta$  and commanding the goal-reaching policy to attain this subgoal,  
 99 until a maximum number of allowed time steps. Pseudocode is provided in Algorithm 2.

---

**Algorithm 1** SuSIE: Zero-Shot, Test-Time Execution

---

**Require:** subgoal model  $p_\theta(s_+ | s_t, l)$ , policy  $\pi_\phi(\cdot | s_t, s_+)$ , language command  $l^{\text{test}}$ , max episode  
 length  $T$ , goal sampling interval  $K$ , initial state  $s_0^{\text{test}}$

```

1:  $t \leftarrow 0$ 
2: while  $t \leq T$  do
3:   Sample  $s_+^{\text{test}} \sim p_\theta(s_+ | s_t^{\text{test}}, l^{\text{test}})$  ▷ Sample a new subgoal every  $K$  steps
4:   for  $j = 1$  to  $k$  do
5:     Sample  $a_t \sim \pi_\phi(\cdot | s_t^{\text{test}}, s_+^{\text{test}})$  ▷ Predict the action from current state and subgoal
6:     Execute  $a_t$ 
7:      $t \leftarrow t + 1$ 
8:   end for
9: end while

```

---

### 100 3 Can SuSIE Generate Plausible and Meaningful Subgoals?

101 We start by presenting qualitative examples of intermediate subgoals generated by the SuSIE image-  
 102 editing model in Figure 4. Even on previously unseen trajectories and language commands, the  
 103 model is able to produce visually high-quality and useful subgoals involving the gripper grasping  
 104 and moving objects. This is nontrivial since it requires the model to have not only the *semantic*  
 105 *knowledge* to detect which pixels in the image correspond to a given object, but also an understand-  
 106 ing of *dynamics* to predict how to move and rotate the gripper to grasp it.

### 107 4 Is the Synthesized Subgoal Useful for Completing New Commands?

108 **Real-world experimental setup and datasets.** We conduct our real-robot experiments on a Wid-  
 109 owX250 robot platform. Our robot dataset is BridgeData V2 (Walke et al., 2023). The dataset  
 110 contains over 60k trajectories, 24 environments, 13 skills, and hundreds of objects. Our evaluations  
 111 present three different scenarios 3, designed specifically to test the ability of various methods at dif-  
 112 ferent levels of open-world generalization: **Scene A:** this scene includes an environment and objects

	No. of Instructions Chained				
	1	2	3	4	5
HULC (Mees et al., 2022a)	0.43	0.14	0.04	0.01	0.00
MCIL (Lynch & Sermanet, 2020)	0.20	0.00	0.00	0.00	0.00
MdetrLC (Ge et al., 2023)	0.69	0.38	0.20	0.07	0.04
AugLC (Ge et al., 2023)	0.69	0.43	<b>0.22</b>	0.09	0.05
LCBC (Walke et al., 2023)	0.62	0.31	0.14	0.05	0.01
UniPi (Ours) (Du et al., 2023)	0.56	0.16	0.08	0.08	0.04
UniPi (HiP) (Ajay et al., 2023b)	0.08	0.04	0.00	0.00	0.00
SuSIE (Ours)	<b>0.75</b>	<b>0.46</b>	0.19	<b>0.11</b>	<b>0.07</b>

Table 1: **Comparison of SuSIE and other prior approaches on CALVIN.** SuSIE is able to chain together more instructions with a higher success rate than all of these prior methods.

113 that are well-represented in BridgeData V2; **Scene B**: this scene is situated in an environment with a  
114 seen tabletop but a novel background and distractors, where the robot must move a seen object (bell  
115 pepper) into a choice of seen container (orange pot) or unseen container (ceramic bowl); and **Scene**  
116 **C**: this scene includes a table texture unlike anything in BridgeData V2 and requires manipulating  
117 both seen and unseen objects.

118 **Real-world results.** We present performance of real-world evaluations in Table 2. Observe that in  
119 Scene A, SuSIE achieves the highest success rate on all three tasks, attaining an average success rate  
120 of 73% which improves over RT-2-X (Anonymous) by **69%**. In Scene B, SuSIE again outperforms  
121 other prior approaches on the two tasks, successfully grounding both the novel ceramic bowl and  
122 the previously seen orange pot. In the most challenging Scene C (unseen domain, unseen objects),  
123 SuSIE attains a success rate of 45%, outperforming MOO (Stone et al., 2023) by about **260%**.  
124 However, RT-2-X outperforms SuSIE in this scene. We believe that the superior performance of RT-  
125 2-X compared to SuSIE in Scene C is because it is a much larger 55B parameter model, initializes  
126 from a proprietary VLM, and is trained on much more data — including BridgeData V2, but also  
127 a vast quantity of additional tabletop manipulation. These differences in the amount of data and  
128 parameters put our method, which only utilizes BridgeData V2, at quite an unfair advantage against  
129 RT-2-X. Nevertheless, SuSIE is still able to recognize the novel objects and attain a high absolute  
130 success rate of 45%.

131 **Simulation results.** We also conducted simulation experiments in CALVIN benchmark (Mees et al.,  
132 2022b), where we compared to LCBC, UniPi (Du et al., 2023) and several other state-of-the-art  
133 methods, and showed that SuSIE outperforms all prior methods on this benchmark. Further details  
134 are described in Section E in Appendix.

## 135 5 Does SuSIE Improve Precision and Low-Level Skill Execution?

136 In this section, we visualize some evaluation rollouts from our experiments in Scene A to understand  
137 if SuSIE works merely because it enhances the generalization of the policy to semantic changes  
138 in the visual observation or if it actually does improve the precision of the low-level control by  
139 commanding meaningful subgoals. Observe in Figure 5 that the RT-2-X policy often produces  
140 actions that fail to precisely orient the gripper around the target object or close the gripper early. In  
141 contrast, policy executions obtained via SuSIE are more precise, and execute actions that attempt to  
142 match the gripper and object positions to the generated subgoal, allowing the policy to succeed at the  
143 task. To understand the contribution of the subgoal prediction towards improved precision, we also  
144 evaluate an oracle GCBC policy on a subset of tasks. This policy is trained on identical robot data  
145 as SuSIE; however, we at test time we command the policy with a real image of the completed task,  
146 which our method does not require. Observe  
147 that even then this GCBC oracle fails to accomplish the task due to issues with imprecise object  
148 localization and untimely gripper closing. Corroborated by numerical results in Table 3,  
149 these experiments validate our claim that utilizing subgoal prediction is crucial for enabling  
150 precise low-level skill execution and control.  
151  
152  
153

Task		LCBC/MOO	RT-2-X	Ours
Scene A	Eggplant on plate	0.4	0.3	<b>0.8</b>
	Carrot on plate	0.3	0.4	<b>0.7</b>
	Eggplant in pot	0.4	0.6	<b>0.7</b>
Average		0.37	0.43	<b>0.73</b>
Scene B	Bell pepper in pot	0.0	0.0	<b>0.2</b>
	Bell pepper in bowl	0.1	0.0	<b>0.4</b>
Average		0.05	0.00	<b>0.30</b>
Scene C	Toothpaste in bowl	0.0	<b>0.5</b>	<b>0.5</b>
	Crayon in bowl	0.0	<b>0.9</b>	0.6
	Spoon in bowl	0.3	<b>0.7</b>	0.4
	Bowl to top	0.2	<b>0.9</b>	0.3
	Average	0.13	<b>0.75</b>	0.45

Table 2: **Real-world performance.** SuSIE consistently achieves the best success rates in Scenes A and B, and is able to attain a high absolute success rate of 45% on the most challenging Scene C.

Table 3: **Comparison to GCBC with oracle goals.** Executing generated subgoals improves the performance of GCBC even when the latter is provided with a real goal image.

Task		GCBC	Ours
Scene A	Eggplant on plate	0.4	<b>0.8</b>
	Carrot on plate	0.4	<b>0.7</b>
	Eggplant in pot	0.5	<b>0.7</b>
CALVIN	8 tasks involving non-prehensile motion	0.16	<b>0.92</b>

154 **References**

- 155 Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B. Tenenbaum, Tommi S. Jaakkola, and Pulkit Agrawal.  
156 Is conditional generative modeling all you need for decision making? In *The Eleventh International  
157 Conference on Learning Representations, 2023a*. URL [https://openreview.net/  
158 forum?id=sPlfo2K9DFG](https://openreview.net/forum?id=sPlfo2K9DFG).
- 159 Anurag Ajay, Seungwook Han, Yilun Du, Shaung Li, Abhi Gupta, Tommi Jaakkola, Josh Tenen-  
160 baum, Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models  
161 for hierarchical planning. *arXiv preprint arXiv:2309.08587*, 2023b.
- 162 Anonymous. Rt-2-x. In *Under review*.
- 163 Mohammad Babaeizadeh, Mohammad Taghi Saffar, Suraj Nair, Sergey Levine, Chelsea Finn,  
164 and Dumitru Erhan. Fitvid: Overfitting in pixel-level video prediction. *arXiv preprint  
165 arXiv:2106.13195*, 2020.
- 166 Chethan Bhateja, Derek Guo, Dibya Ghosh, Anika Singh, Manan Tomar, Quan Ho Vuong, Yevgen  
167 Chebotar, Sergey Levine, and Aviral Kumar. Robotic offline rl from internet videos via value-  
168 function pre-training. 2023. URL [https://api.semanticscholar.org/CorpusID:  
169 262217278](https://api.semanticscholar.org/CorpusID:262217278).
- 170 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn,  
171 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics  
172 transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- 173 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choroman-  
174 ski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action  
175 models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023a.
- 176 Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho,  
177 Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. Do as i can, not as i say: Grounding  
178 language in robotic affordances. In *Conference on Robot Learning*, pp. 287–318. PMLR, 2023b.
- 179 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image  
180 editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
181 Recognition*, pp. 18392–18402, 2023.
- 182 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter  
183 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning  
184 via sequence modeling. *arXiv preprint arXiv:2106.01345*, 2021.
- 185 Zoey Chen, Sho Kiami, Abhishek Gupta, and Vikash Kumar. Genau: Retargeting behaviors to  
186 unseen situations via generative augmentation. *arXiv preprint arXiv:2302.06671*, 2023.
- 187 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shu-  
188 ran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint  
189 arXiv:2303.04137*, 2023.
- 190 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi  
191 Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language mod-  
192 els. *arXiv preprint arXiv:2210.11416*, 2022.
- 193 Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,  
194 Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multi-  
195 modal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- 196 Yilun Du, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Joshua B Tenenbaum, Dale Schu-  
197 urmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *arXiv  
198 preprint arXiv:2302.00111*, 2023.
- 199 Frederik Ebert, Chelsea Finn, Sudeep Dasari, Annie Xie, Alex Lee, and Sergey Levine. Visual fore-  
200 sight: Model-based deep reinforcement learning for vision-based robotic control. *arXiv preprint  
201 arXiv:1812.00568*, 2018.

- 202 Yuying Ge, Annabella Macaluso, Li Erran Li, Ping Luo, and Xiaolong Wang. Policy adaptation  
203 from foundation model feedback. In *Proceedings of the IEEE/CVF Conference on Computer  
204 Vision and Pattern Recognition*, pp. 19059–19069, 2023.
- 205 Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Gird-  
206 har, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in  
207 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
208 and Pattern Recognition*, pp. 18995–19012, 2022.
- 209 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning  
210 behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 211 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains  
212 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 213 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine.  
214 Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint  
215 arXiv:2304.10573*, 2023.
- 216 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint  
217 arXiv:2207.12598*, 2022.
- 218 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in  
219 Neural Information Processing Systems*, 33:6840–6851, 2020.
- 220 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P.  
221 Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High  
222 definition video generation with diffusion models, 2022a.
- 223 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J  
224 Fleet. Video diffusion models. *arXiv:2204.03458*, 2022b.
- 225 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as  
226 zero-shot planners: Extracting actionable knowledge for embodied agents. *arXiv preprint  
227 arXiv:2201.07207*, 2022a.
- 228 Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan  
229 Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through  
230 planning with language models. *arXiv preprint arXiv:2207.05608*, 2022b.
- 231 Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence  
232 modeling problem. In *Advances in Neural Information Processing Systems*, 2021.
- 233 Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for  
234 flexible behavior synthesis. In *International Conference on Machine Learning*, 2022.
- 235 Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Car-  
236 ion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the  
237 IEEE/CVF International Conference on Computer Vision*, pp. 1780–1790, 2021.
- 238 Ivan Kapelyukh, Vitalis Vosylius, and Edward Johns. Dall-e-bot: Introducing web-scale diffusion  
239 models to robotics. 2023.
- 240 Siddharth Karamcheti, Suraj Nair, Annie S Chen, Thomas Kollar, Chelsea Finn, Dorsa Sadigh,  
241 and Percy Liang. Language-driven representation learning for robotics. *arXiv preprint  
242 arXiv:2302.12766*, 2023.
- 243 Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International  
244 Conference on Learning Representations (ICLR)*, 2015.
- 245 Alex X Lee, Anusha Nagabandi, Pieter Abbeel, and Sergey Levine. Stochastic latent actor-critic:  
246 Deep reinforcement learning with a latent variable model. *Advances in Neural Information Pro-  
247 cessing Systems*, 33:741–752, 2020.

- 248 Donghwan Lee and Niao He. Stochastic primal-dual q-learning. *arXiv preprint arXiv:1810.08298*,  
249 2018.
- 250 Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and  
251 Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE*  
252 *International Conference on Robotics and Automation (ICRA)*, pp. 9493–9500. IEEE, 2023.
- 253 Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data.  
254 *arXiv preprint arXiv:2005.07648*, 2020.
- 255 Yecheng Jason Ma, William Liang, Vaidehi Som, Vikash Kumar, Amy Zhang, Osbert Bastani, and  
256 Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control. *arXiv*  
257 *preprint arXiv:2306.00958*, 2023.
- 258 Zhao Mandi, Homanga Bharadhwaj, Vincent Moens, Shuran Song, Aravind Rajeswaran, and Vikash  
259 Kumar. Cacti: A framework for scalable multi-task multi-scene visual imitation learning. *arXiv*  
260 *preprint arXiv:2212.05711*, 2022.
- 261 Oier Mees, Lukas Hermann, and Wolfram Burgard. What matters in language conditioned robotic  
262 imitation learning over unstructured data. *IEEE Robotics and Automation Letters*, 7(4):11205–  
263 11212, 2022a.
- 264 Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for  
265 language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics*  
266 *and Automation Letters (RA-L)*, 7(3):7327–7334, 2022b.
- 267 Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey  
268 Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Sim-  
269 ple open-vocabulary object detection. In *European Conference on Computer Vision*, pp. 728–755.  
270 Springer, 2022.
- 271 Vivek Myers, Andre He, Kuan Fang, Homer Walke, Philippe Hansen-Estruch, Ching-An Cheng,  
272 Mihai Jalobeanu, Andrey Kolobov, Anca Dragan, and Sergey Levine. Goal representations  
273 for instruction following: A semi-supervised language interface to control. *arXiv preprint*  
274 *arXiv:2307.00117*, 2023.
- 275 Suraj Nair and Chelsea Finn. Hierarchical foresight: Self-supervised learning of long-horizon tasks  
276 via visual subgoal generation. *arXiv preprint arXiv:1909.05829*, 2019.
- 277 Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhi Gupta. R3m: A universal  
278 visual representation for robot manipulation. *ArXiv*, abs/2203.12601, 2022.
- 279 Alexander Pashevich, Robin Strudel, Igor Kalevatykh, Ivan Laptev, and Cordelia Schmid. Learn-  
280 ing to augment synthetic images for sim2real policy transfer. In *2019 IEEE/RSJ International*  
281 *Conference on Intelligent Robots and Systems (IROS)*, pp. 2651–2657. IEEE, 2019.
- 282 Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. FiLM: Visual  
283 Reasoning with a General Conditioning Layer, December 2017.
- 284 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
285 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
286 models from natural language supervision. In *International conference on machine learning*, pp.  
287 8748–8763. PMLR, 2021.
- 288 Rafael Rafailov, Tianhe Yu, A. Rajeswaran, and Chelsea Finn. Offline reinforcement learning from  
289 images with latent space models. *Learning for Decision Making and Control (LADC)*, 2021.
- 290 Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with  
291 action-free pre-training from videos. In *International Conference on Machine Learning*, pp.  
292 19561–19579. PMLR, 2022.
- 293 Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic  
294 manipulation. In *Conference on Robot Learning*, pp. 894–906. PMLR, 2022.

- 295 Austin Stone, Ted Xiao, Yao Lu, Keerthana Gopalakrishnan, Kuang-Huei Lee, Quan Vuong, Paul  
296 Wohlhart, Brianna Zitkovich, Fei Xia, Chelsea Finn, et al. Open-world object manipulation using  
297 pre-trained vision-language models. *arXiv preprint arXiv:2303.00905*, 2023.
- 298 Ruben Villegas, Arkanath Pathak, Harini Kannan, Dumitru Erhan, Quoc V Le, and Honglak Lee.  
299 High fidelity video prediction with large stochastic recurrent neural networks. *Advances in Neural*  
300 *Information Processing Systems*, 32, 2019.
- 301 Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony Zhao,  
302 Philippe Hansen-Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang, Chelsea Finn, and  
303 Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot*  
304 *Learning (CoRL)*, 2023.
- 305 Philipp Wu, Alejandro Escontrela, Danijar Hafner, Pieter Abbeel, and Ken Goldberg. Daydreamer:  
306 World models for physical robot learning. In *Conference on Robot Learning*, pp. 2226–2240.  
307 PMLR, 2023.
- 308 Yinfei Yang, Daniel Cer, Amin Ahmad, Mandy Guo, Jax Law, Noah Constant, Gustavo Hernandez  
309 Abrego, Steve Yuan, Chris Tar, Yun-Hsuan Sung, et al. Multilingual Universal Sentence Encoder  
310 for Semantic Retrieval, July 2019.
- 311 Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in  
312 projected latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
313 *Pattern Recognition*, pp. 18456–18466, 2023a.
- 314 Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn.  
315 Combo: Conservative offline model-based policy optimization. *arXiv preprint arXiv:2102.08363*,  
316 2021.
- 317 Tianhe Yu, Ted Xiao, Austin Stone, Jonathan Tompson, Anthony Brohan, Su Wang, Jaspier Singh,  
318 Clayton Tan, Jodilyn Peralta, Brian Ichter, et al. Scaling robot learning with semantically imag-  
319 ined experience. *arXiv preprint arXiv:2302.11550*, 2023b.
- 320 Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual  
321 manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.



## 322 A Introduction

323 A useful generalist robot must be able to — much like a person — recognize and reason about novel  
324 objects and scenarios it has never encountered before. For example, if a user instructs the robot to  
325 “hand me that jumbo orange crayon,” it ought to be able to do so even if it has never interacted  
326 with a jumbo orange crayon before. In other words, the robot needs to possess not only the physical  
327 capability to manipulate an object of that shape and size but also the semantic understanding to  
328 reason about an object outside of its training distribution. As much as robotic manipulation datasets  
329 have grown in recent years, it is unlikely that they will ever include every conceivable instance  
330 of objects and settings, any more so than the life experiences of a person ever include physical  
331 interactions with every type of object. While these datasets contain more than enough examples of  
332 manipulating elongated cylindrical objects, they lack the broad semantic knowledge necessary to  
333 ground the *particular* objects that robots will undoubtedly encounter during everyday operation.

334 How can we imbue this semantic knowledge into language-guided robotic control? One approach  
335 to do this would be to utilize pre-trained models trained on vision and language to initialize differ-  
336 ent components in the robotic learning pipeline. Recent efforts attempt to do this, for example, by  
337 initializing robotic policies with pre-trained vision-language encoders (Brohan et al., 2023a) or uti-  
338 lizing pre-trained models for generating semantic scene augmentation (Chen et al., 2023; Yu et al.,  
339 2023b). While these methods bring semantic knowledge into robot learning, it remains unclear  
340 if these approaches realize the full potential of Internet pre-training in improving low-level motor  
341 control and policy execution, or whether they simply improve visual generalization of the policy.

342 In this paper, we develop an approach for leveraging a class of pre-trained image-editing models  
343 (e.g., InstructPix2Pix (Brooks et al., 2023)) for improving motor control and policy execution. Our  
344 key insight is to decompose the robotic control problem into two phases: first, synthesizing a “sub-  
345 goal” that the robot must reach to complete the user-specified task, and then, attempting to reach  
346 this subgoal via a goal-reaching robot controller. The first phase of this recipe incorporates semantic  
347 information by fine-tuning an image-editing model on robot data such that, given the robot’s current  
348 observation and a natural language command, the model generates a hypothetical *future* subgoal  
349 that allows the robot to complete the command. We then employ a low-level goal-reaching policy  
350 to reach this hypothetical future subgoal. Crucially, our image-editing model does not need to un-  
351 derstand *how* to achieve this future subgoal, and on the other hand, the policy only needs to infer  
352 visuo-motor relationships to determine the correct actuation and does not require an understanding  
353 of the semantics. Furthermore, such subgoals can significantly simplify the task by inferring likely  
354 poses for the arm or intermediate sub-steps, such as grasping an object when the command requires  
355 repositioning it to a new location (see Figure 2). In fact, we observe in our experiments that while  
356 existing approaches often fail due to imprecise understanding of obstacles or object orientations,  
357 following the generated subgoals enables our method to perform well in such scenarios.

358 The main contribution of our work is **SUBgoal Synthesis via Image Editing (SuSIE)**, a simple and  
359 scalable method for incorporating semantic information in pre-trained models to improve robotic  
360 control. The pre-trained image editing model is used with minimal modification, requiring only  
361 fine-tuning on robot data. The low-level goal-conditioned policy is trained with standard supervised  
362 learning, and faces the comparatively easier problem of reaching nearby image subgoals; this typi-  
363 cally only requires attending to a single object or the arm position, ignoring most parts of the scene.  
364 Together, we find that this approach solves real robot control tasks involving novel objects, novel  
365 distractors, and even novel scenes, all of which are not observed at all in the robot training data.

## 366 B Related Work

367 **Incorporating semantic information from vision-language pre-trained models.** Prior works that  
368 incorporate semantic information from vision-language pre-trained models into robot learning can  
369 be classified into two categories. The first category aims to improve visual scene understanding in  
370 robot policies with semantic information from VLMs. For instance, GenAug (Chen et al., 2023),  
371 ROSIE (Yu et al., 2023b), DALL-E-Bot (Kapelyukh et al., 2023), and CACTI (Mandi et al., 2022)  
372 use text-to-image generative models to produce semantic augmentations of a given scene with novel  
373 objects and arrangements and train the robot policy on the augmented data to enable it to perform  
374 well in a similar scene. MOO Stone et al. (2023) utilizes a pre-trained object detector to extract  
375 bounding boxes that guide the robot policy towards the object of interest. Other works directly  
376 train language and image-conditioned policies (Brohan et al., 2022, 2023a; Shridhar et al., 2022), by

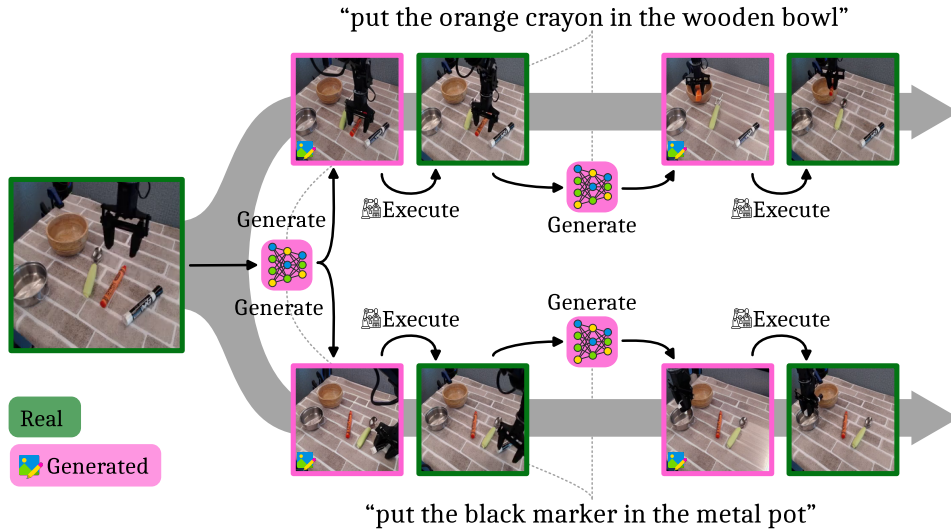


Figure 2: **SuSIE leverages a pre-trained image editing model to generate future image subgoals based on a language commands.** A low-level goal-reaching policy then executes the actions needed to reach each subgoal. Alternating this loop enables us to solve the task.

377 utilizing frozen or fine-tuned off-the-shelf VLMs (Driess et al., 2023; Radford et al., 2021) on robot  
 378 data to produce action sequences (Brohan et al., 2023a).

379 While these approaches do utilize pre-trained models, we find in our experiments, that pre-training  
 380 using VLMs (e.g., Brohan et al. (2023a)) does not necessarily enhance low-level motor control,  
 381 in the sense that learned policies often localize the object or move the gripper imprecisely (see  
 382 Figure 5). On the other hand, our approach is able to incorporate benefits of pre-training in syn-  
 383 thesizing subgoals that carefully steer the motion of the low goal-conditioned policy, improving its  
 384 precision. Also, while our approach can be directly applied in *unstructured* open-world settings, ap-  
 385 plying GenAug (Chen et al., 2023), MOO (Stone et al., 2023), and ROSIE (Yu et al., 2023b) requires  
 386 additional information about the scene, such as clean object bounding boxes or 3D object meshes.  
 387 This significantly restricts their applicability to scenarios where this additional information is not  
 388 available: for example, GenAug is not applicable in our real-world experiments since 3D object  
 389 meshes for new target objects are not available. Distinct from our approach for utilizing generative  
 390 models, other works design representation learning objectives for vision-language pre-training for  
 391 control (Nair et al., 2022; Ma et al., 2023; Karamcheti et al., 2023; Bhateja et al., 2023), but these  
 392 methods still need to utilize limited amounts of data from the target task to learn a policy.

393 The second category of approaches also incorporates semantic information from pre-trained models  
 394 for planning. Most approaches in this category use pre-trained models to imagine visual (Du et al.,  
 395 2023; Ajay et al., 2023b) or textual plans (Brohan et al., 2023b; Huang et al., 2022a,b; Liang et al.,  
 396 2023), which then inform a low-level robot control policy. Low-level policies conditioned on text  
 397 suffer from a grounding problem, which our approach circumvents entirely since the low-level con-  
 398 trol policy only observes image-based plans. Perhaps the most related are UniPi (Du et al., 2023)  
 399 and HiP (Ajay et al., 2023b), which train video models to generate a sequence of frames achieving  
 400 the target task, and then extract robot actions from an inverse dynamics model. Our approach does  
 401 *not* attempt to generate full videos (i.e., *all* frames in a rollout), but only the next waypoint that a  
 402 low-level policy must achieve to solve the commanded task. While this difference might appear  
 403 small, it has major implications: modeling an entire video puts a very high burden on the generative  
 404 model, requiring the frames to obey strict physical consistency. Unfortunately, we find that current  
 405 video models often produce temporally inconsistent frames (“hallucinations”), which only confuse  
 406 the low-level controller, inhibiting it from completing the task. Indeed, the control evaluations in  
 407 such prior works often focus on simpler simulated environments. Our method provides more free-  
 408 dom to the low-level controller to handle the physical aspects of the task over a longer time interval  
 409 while providing higher-level guidance at a level that is suitable to the diffusion model’s ability to  
 410 preserve physical plausibility. In our experiments, we find that our method significantly improves  
 411 over a reimplementations of UniPi (Du et al., 2023).

412 **Classical model-based RL and planning with no pre-training.** The idea behind our approach  
 413 is also related to several methods in the deep RL literature that do not use pre-trained models and  
 414 generally do not study language-guided control. For instance, (Hafner et al., 2019; Lee et al., 2020;  
 415 Yu et al., 2021; Wu et al., 2023; Rafailov et al., 2021; Hafner et al., 2023) train action-conditioned  
 416 dynamics models and run RL in the model. While our approach also models multi-step dynamics,  
 417 our model is not conditioned on an action input. Removing the dependency on an action input  
 418 enables us to de-couple the fine-tuning of the (large) image-editing model from the policy entirely,  
 419 improving simplicity and time efficiency. APV (Seo et al., 2022) trains an action agnostic dynamics  
 420 model from videos but fine-tunes it in a loop with the policy with actions, and hence, does not enjoy  
 421 the above benefits. Finally, these model-based RL methods do not exhibit zero-shot generalization  
 422 abilities to new tasks, which is an important capability that our method enjoys. Our approach is  
 423 also related to several video prediction methods (Ebert et al., 2018; Lee & He, 2018; Babaeizadeh  
 424 et al., 2020; Villegas et al., 2019) but utilizes a better neural network architecture (i.e., diffusion  
 425 models instead of LSTMs and CNNs). Most related is to our method is hierarchical visual foresight  
 426 (HVF) (Nair & Finn, 2019): while HVF utilizes MPC to find an action, our approach simply utilizes  
 427 a goal-reaching policy thereby eliminating the cost of running MPC with large dynamics models.

428 Our approach is also related to several prior works that utilize generative models for planning in a  
 429 single-task setting, with no pre-training. Trajectory transformer (TT) (Janner et al., 2021), decision  
 430 transformer (DT) (Chen et al., 2021), and their extensions condition the policy on the target return  
 431 or goal. While diffusion-based variants of these methods (Janner et al., 2022; Ajay et al., 2023a) use  
 432 diffusion models to model long-term rollout distributions over states, actions, and rewards, they still  
 433 require training data from the target task to learn a policy, unlike our zero-shot planning approach.

## 434 C Preliminaries and Problem Statement

435 We consider the problem setting of language-conditioned robotic control. Specifically, we want  
 436 a robot to accomplish the task described by a novel language command. We study this problem  
 437 in the context of learning from a dataset  $\mathcal{D}$  of language-labeled robot trajectories, and optionally,  
 438 an additional dataset,  $\mathcal{D}'$  of robot data, which is not annotated any task labels (e.g., play data).  
 439 Formally,  $\mathcal{D} = \{(\tau^1, l_1), (\tau^2, l_2), \dots, (\tau^N, l_N)\}$ , where each rollout  $\tau^i$  consists of a sequence of  
 440 scenes (or states)  $s_k^i \in \mathcal{S}$  and actions  $a_k^i \in \mathcal{A}$  that were executed while collecting this data, i.e.,  
 441  $\tau^i = (s_0^i, a_0^i, \dots, s_k^i, a_k^i, \dots)$ , following the standard assumptions of a Markov decision process.  $l_i$   
 442 is a natural language command describing the task accomplished in the trajectory.  $\mathcal{D}'$  is organized  
 443 similarly to  $\mathcal{D}$ , but does not contain any language annotations  $l_i$ . At test time, given a new scene  $s_0^{\text{test}}$   
 444 and a new natural language description  $l^{\text{test}}$  of a task, we evaluate a method in terms of its success  
 445 rate at accomplishing this task starting from this scene,  $s_0^{\text{test}}$ .

## 446 D SuSIE: Subgoal Synthesis via Image Editing

447 Our goal is to utilize semantic information from the Internet to improve language-guided robot  
 448 control in novel environments, scenes, and objects. How can we do this when models trained on  
 449 general-purpose Internet data do not provide guidance in selecting low-level actions? Our key insight  
 450 is that we can still utilize some sort of a pre-trained model for guiding low-level control if we could  
 451 decouple the robot control problem into two phases: **(i)** imagining subgoals that would need to  
 452 be attained to succeed at the task, and **(ii)** learning low-level control policies for reaching these  
 453 generated subgoals. Our method incorporates semantic information from Internet pre-training in  
 454 phase (i), by fine-tuning a text-guided image-editing model for subgoal generation. Phase (ii) is  
 455 accomplished via a goal-conditioned policy trained only on robot data. We describe each of these  
 456 phases below and then summarize the resulting robot controller.

### 457 D.1 Phase (i): Synthesizing Subgoals From Image Editing Models

458 The primary component of our method is a generative model that, given a target task specified in  
 459 natural language, can guide the low-level controller towards a state that it must try to attain in order  
 460 to solve the task. One way to accomplish this is to train a generative model to produce an immediate  
 461 next way-point or subgoal frame. We can then incorporate semantic information from the Internet  
 462 into our algorithm by initializing this generative model with a suitable pre-trained initialization,  
 463 followed by fine-tuning it on multi-task, diverse robot data.

464 What is a good pre-trained initialization for initializing this model? Our intuition is that since ac-  
 465 complishing a task is equivalent to “editing” the pixels of an image of the robot workspace under  
 466 controls prescribed by the language command, a favorable pre-trained initialization is provided by  
 467 a language-guided image-editing model. We instantiate our approach with Instruct pix2pix (Brooks  
 468 et al., 2023), though other image editing models could also be used. Formally, this model is given  
 469 by  $p_\theta(\mathbf{s}_{\text{edited}}|\mathbf{s}_{\text{orig}}, l)$ . Then, using the dataset  $\mathcal{D}$  of robot trajectories, we fine-tune  $p_\theta$  on tuples  
 470 containing a pair of images sampled from a trajectory and the corresponding language annotation:  
 471  $(\mathbf{s}_{\text{orig}} := \mathbf{s}_i^k, \mathbf{s}_{\text{edited}} := \mathbf{s}_j^k, l_k)$ , where  $\mathbf{s}_j$  is a state that appears after  $\mathbf{s}_i$  ( $j > i$ ). During fine-tuning,  
 472 we run gradient descent on the following objective, starting from  $\theta_0 := \theta_{\text{pre-trained}}$ :

$$\min_{\theta} - \mathbb{E}_{(\tau^k, l_k) \sim \mathcal{D}; \mathbf{s}_i^k \sim \tau^k; j \sim q(j|i)} [\log p_\theta(\mathbf{s}_j^k | \mathbf{s}_i^k, l_k)] . \quad (3)$$

473 We need to choose the distribution  $q$  over the time-step  $j$  given a state  $\mathbf{s}_i^k$  for fine-tuning the image-  
 474 editing model as in Equation 3. Since we model the next subgoal that the low-level controller should  
 475 attain, and since the depending upon the task, this subgoal could be arbitrarily close to the original  
 476 state  $\mathbf{s}_i$ , we require valid tuples  $(\mathbf{s}_i, \mathbf{s}_j, l)$  used for fine-tuning  $p_\theta$  in Equation 3 to have values of  $j$  in  
 477 a bounded interval around  $i$ , specifically we choose  $j \in [i, i+k]$ , where  $k$  is a fixed hyperparameter.

## 478 D.2 Phase (ii): Reaching Generated Sub-Goals with Goal-Conditioned Policies

479 In order to utilize the fine-tuned image-editing model to actually control the robot, we further need  
 480 to train a low-level controller to actually select suitable robot actions. In this section, we present  
 481 the design of our low-level controller, followed by a full description our test-time control procedure.  
 482 Since the image-editing model in SuSIE produces images of future subgoals conditioned on natural  
 483 language task descriptions, our low-level controller can simply be a language-agnostic goal-reaching  
 484 policy that aims to reach these generated subgoals.

485 **Training a goal-reaching policy.** Our goal-reaching policy is parameterized as  $\pi_\phi(\cdot | \mathbf{s}_i, \mathbf{s}_j)$ , where  
 486  $\mathbf{s}_j$  is a future frame that the policy intends to reach, by acting at  $\mathbf{s}_i$ . At test time, we only need the  
 487 low-level goal-conditioned policy to be proficient at reaching close-by states that lie within  $k$  steps  
 488 of a given state since the image editing model from phase (i) is also trained to produce subgoals  
 489 within  $k$  steps of any state. To train this policy, we run goal-conditioned behavioral cloning (GCBC)  
 490 on the robot data, utilized previously in phase (i). In addition, we can also leverage robot data  $\mathcal{D}'$   
 491 that does not contain language annotations. Formally, our training objective is given by:

$$\max_{\phi} \mathbb{E}_{\tau^i \sim \mathcal{D} \cup \mathcal{D}'; (\mathbf{s}_i^k, \mathbf{a}_i^k) \sim \tau^k; j \sim q(j|i)} [\log \pi_\phi(\mathbf{a}_i^k | \mathbf{s}_i^k, \mathbf{s}_j^k)] , \quad (4)$$

492 where  $q(j|i)$  is the distribution over future frames that we previously utilized in Equation 3.

493 **Test-time control with  $\pi_\phi$  and  $p_\theta$ .** Once both the goal-reaching policy  $\pi_\phi$  and the image edit-  
 494 ing subgoal generation model  $p_\theta$  are trained, we utilize them together to solve new manipulation  
 495 tasks based on user-specified natural language commands. Given a new scene,  $\mathbf{s}_0^{\text{test}}$ , and a lan-  
 496 guage task description  $l^{\text{test}}$ , SuSIE attempts to solve the task by iteratively generating subgoals  
 497 and commanding the low-level goal-reaching policy with these subgoals. At the start, we sample  
 498 the first subgoal  $\widehat{\mathbf{s}}_+^{\text{test}} \sim p_\theta(\cdot | \mathbf{s}_0^{\text{test}}, l^{\text{test}})$ . Once the subgoal is generated, we then roll out the goal-  
 499 reaching policy  $\pi_\phi$ , conditioned on  $\widehat{\mathbf{s}}_+^{\text{test}}$ , for  $k$  time-steps, such that each action is chosen according  
 500 to  $\mathbf{a}_j^{\text{test}} \sim \pi_\phi(\cdot | \mathbf{s}_j^{\text{test}}, \widehat{\mathbf{s}}_+^{\text{test}})$ . After  $k$  time steps, given the current image  $\mathbf{s}_k^{\text{test}}$ , we refresh the subgoal by  
 501 sampling from the image-editing model again and repeat the process. Note crucially that this recipe  
 502 does not require that the subgoal  $\mathbf{s}^{\text{test}}$  be attained after  $k$  steps, as the generative model effectively  
 503 “replans” a new subgoal based on the current observation. Overall, at test time, we alternate between  
 504 obtaining a new subgoal from  $p_\theta$  and commanding the goal-reaching policy to attain this subgoal,  
 505 until a maximum number of allowed time steps. Pseudocode is provided in Algorithm 2.

## 506 D.3 Implementation Details

507 In Phase (i), we utilize the pre-trained initialization from the InstructPix2Pix model (Brooks et al.,  
 508 2023), trained to perform language-guided image editing and fine-tune it on our robot dataset. Since  
 509 the InstructPix2Pix model utilizes a UNet-based diffusion model architecture, we implement Equa-  
 510 tion 3 using a variational lower bound objective, following the standard recipe for training diffusion  
 511 models. Our image-editing diffusion model operates on images of size  $256 \times 256$ . The language in-  
 512 structions are encoded with a frozen CLIP encoder (Radford et al., 2021). To ensure that this model

---

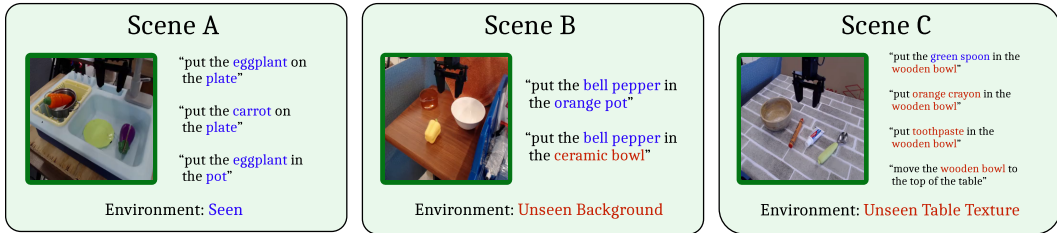
**Algorithm 2** SuSIE: Zero-Shot, Test-Time Execution

---

**Require:** subgoal model  $p_\theta(s_+|s_t, l)$ , policy  $\pi_\phi(\cdot | s_t, s_+)$ , language command  $l^{\text{test}}$ , max episode length  $T$ , goal sampling interval  $K$ , initial state  $s_0^{\text{test}}$

- 1:  $t \leftarrow 0$
- 2: **while**  $t \leq T$  **do**
- 3:   Sample  $s_+^{\text{test}} \sim p_\theta(s_+|s_t^{\text{test}}, l^{\text{test}})$  ▷ Sample a new subgoal every  $K$  steps
- 4:   **for**  $j = 1$  to  $k$  **do**
- 5:     Sample  $\mathbf{a}_t \sim \pi_\phi(\cdot | s_t^{\text{test}}, s_+^{\text{test}})$  ▷ Predict the action from current state and subgoal
- 6:     Execute  $\mathbf{a}_t$
- 7:      $t \leftarrow t + 1$
- 8:   **end for**
- 9: **end while**

---



Unseen in BridgeData / Seen in BridgeData

Figure 3: **Real-world experimental setup.** We evaluate our method in 3 real-world scenes. The scenes become progressively more difficult from left to right, due to both an increasing visual departure from the robot training data and an increasingly confounding mixture of both seen and unseen objects.

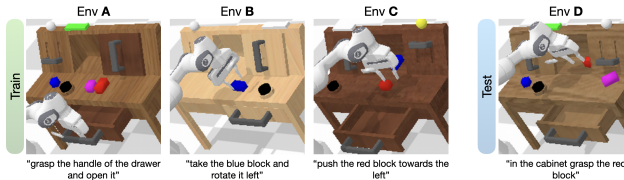
513 pays attention to the input state and the language command it is conditioned on, we apply classifier-  
514 free guidance (Ho & Salimans, 2022) separately to both the language and the image, similarly to  
515 InstructPix2Pix. To obtain a robust goal-reaching policy in Phase (ii), we follow the implementation  
516 details prescribed by Walke et al. (2023). More details about the training hyperparameters and the  
517 architecture of this goal-reaching policy are provided in Appendix G.1.1.

## 518 E Experimental Evaluation

519 The goal of our experiments is to evaluate the efficacy of SuSIE at improving generalization and  
520 motor control in open-world robotic manipulation tasks. To this end, our experiments aim to study  
521 the following questions: **(1)** Can SuSIE generate plausible subgoals for novel tasks, objects and  
522 environments, even those that lie outside of the robot training distribution? **(2)** Are the generated  
523 subgoals useful for solving a task specified by a novel language command, in zero-shot?, **(3)** Does  
524 SuSIE exhibit an elevated level of precision and dexterity compared to other approaches that do  
525 not use subgoals?, and **(4)** How crucial is pre-training on Internet data for attaining zero-shot gen-  
526 eralization? To answer these questions, our experiments compare SuSIE to several prior methods  
527 including state-of-the-art approaches for training language-conditioned policies that leverage pre-  
528 trained vision-language models in a variety of ways.

### 529 E.1 Experimental Scenarios and Comparisons

530 **Real-world experimental setup and datasets.** We conduct our real-robot experiments on a Wid-  
531 owX250 robot platform. Our robot dataset is BridgeData V2 (Walke et al., 2023), a large and diverse  
532 dataset of robotic manipulation behaviors designed for evaluating open-vocabulary instructions. The  
533 dataset contains over 60k trajectories, 24 environments, 13 skills, and hundreds of objects. Our eval-  
534 uations present three different scenarios 3, designed specifically to test the ability of various methods  
535 at different levels of open-world generalization: **Scene A:** this scene includes an environment and  
536 objects that are well-represented in BridgeData V2; **Scene B:** this scene is situated in an environ-  
537 ment with a seen tabletop but a novel background and distractors, where the robot must move a seen  
538 object (bell pepper) into a choice of seen container (orange pot) or unseen container (ceramic bowl);  
539 and **Scene C:** this scene includes a table texture unlike anything in BridgeData V2 and requires ma-  
540 nipulating both seen and unseen objects. We expect Scene C to be the hardest since the robot needs  
541 to carefully ground the language command to identify the correct object while resisting its affinity  
542 for an object that is well-represented in the data (the spoon).



543 **Simulation tasks.** We run our simulation  
 544 experiments in CALVIN (Mees  
 545 et al., 2022b), a benchmark for long  
 546 horizon, language-conditioned manipulation.  
 547 CALVIN consists of four  
 548 simulated environments, A, B, C, D,

549 and each environment comes with a dataset of human-collected play trajectories. Approximately  
 550 35% of these rollouts are annotated with language. Each environment consists of a Franka Emika  
 551 Panda robot arm positioned next to a desk with various manipulatable objects, including a drawer,  
 552 sliding cabinet, light switch, and various colored blocks. Environments are differentiated by the  
 553 positions of these objects and their textures. With this benchmark, we study the most challenging  
 554 zero-shot multi-environment scenario: training on A, B, and C, and testing on D. We follow the  
 555 evaluation protocol from Mees et al. (2022b). During evaluation, a policy given a fixed number of  
 556 timesteps (default 360) to complete a chain of five language instructions.

557 **Comparisons.** Our experiments cover methods that utilize pre-trained models of vision and language  
 558 in language-guided robot control in a variety of ways. While there are several prior methods  
 559 that tackle language-based robotic control as we discuss in Section B, in our experiments, we choose  
 560 to compare to a representative subset of these prior methods to maximally cover the possible set  
 561 of comparisons. We compare to (a) RT-2 (Brohan et al., 2023a) which is one of the most recent  
 562 works utilizing a pre-trained VLM for initializing the robot policy (specifically, RT-2-X (Anonymous),  
 563 which was also trained and evaluated on BridgeData V2), generalizing prior work (Shridhar  
 564 et al., 2022); (b) MOO (Stone et al., 2023), which utilizes pre-trained object detectors to obtain  
 565 bounding box information for the policy and then trains a language-conditioned behavioral cloning  
 566 policy (denoted as “LCBC/MOO”); and (c) UniPi (Du et al., 2023), which trains an entire language-  
 567 conditioned video prediction model starting from a pre-trained video initialization. Since the original  
 568 UniPi model utilized proprietary pre-trained initializations that are not available publicly, we  
 569 re-implemented this method using our own video model following the guidelines in Du et al. (2023),  
 570 but were unable to obtain high-quality generations. In our simulation experiments though, we also  
 571 evaluate another reimplementation of UniPi, using a video diffusion model trained by concurrent  
 572 work (Ajay et al., 2023b). We present details for MOO and UniPi baselines in Appendix G.2.  
 573 Finally, we remark that while we did try to apply GenAug (Chen et al., 2023) as a representative  
 574 semantic augmentation approach in our real-world experiments, we were not able to obtain 3D mesh  
 575 predictions for objects in Bridgedata V2, needed for this approach.

576 We also compare to language-conditioned behavioral cloning (“LCBC”) (Walke et al., 2023), trained  
 577 to produce actions conditioned on an embedding of the natural language task description (Walke  
 578 et al., 2023); and an oracle goal-conditioned behavioral cloning (“GCBC oracle”) approach for tasks  
 579 that require manipulating objects previously seen in the robot data. We observed that in Scene A,  
 580 simple LCBC outperforms MOO. However, in Scenes B and C, which include tasks with unseen objects,  
 581 MOO is crucial for achieving non-zero success. Hence, we report LCBC in Scene A and MOO  
 582 in Scenes B and C. In simulation, we also compare to additional methods previously studied on the  
 583 CALVIN benchmark. These include methods that explicitly tackle long-horizon language-based  
 584 control on CALVIN such as multi-context imitation (MCIL) (Lynch & Sermanet, 2020), hierarchical  
 585 universal language-conditioned policy (HULC) (Mees et al., 2022a), and improved variants of  
 586 HULC (Ge et al., 2023). We also compare to other state-of-the-art methods from Ge et al. (2023) that  
 587 employ an identical training and evaluation protocol as our experiments, namely MdetrLC (Kamath  
 588 et al., 2021), and AugLC (Pashkevich et al., 2019).

## 589 E.2 Can SuSIE Generate Plausible and Meaningful Subgoals?

590 To answer question (1), we start by presenting qualitative examples of intermediate subgoals generated  
 591 by the SuSIE image-editing model in Figure 4. Even on previously unseen trajectories and  
 592 language commands, the model is able to produce visually high-quality and useful subgoals involving  
 593 the gripper grasping and moving objects. This is nontrivial since it requires the model to have  
 594 not only the *semantic knowledge* to detect which pixels in the image correspond to a given object,  
 595 but also an understanding of *dynamics* to predict how to move and rotate the gripper to grasp it.

## 596 E.3 Is the Synthesized Subgoal Useful for Completing New Commands?

597 **Simulation results.** We present performance for SuSIE and other comparisons in Table 1, in terms  
 598 of success rates (out of 1.0) for completing each language instruction in the chain. Observe that

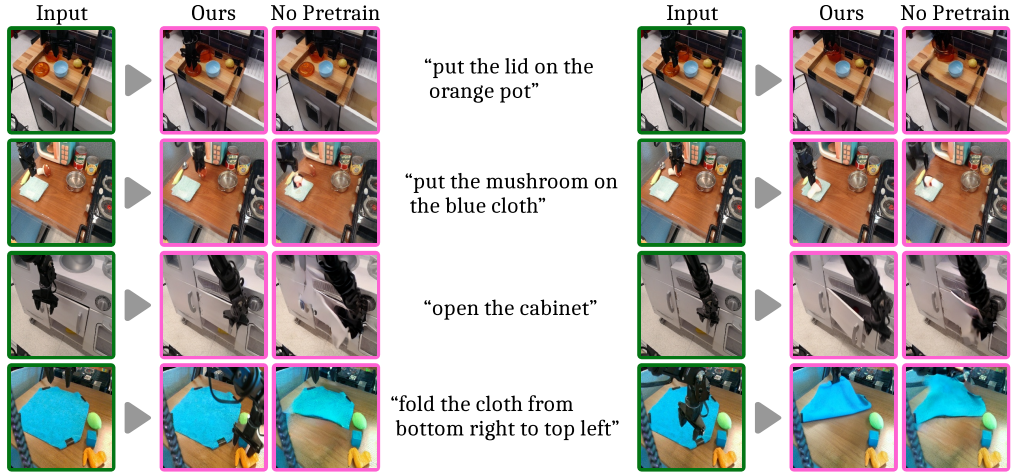


Figure 4: **Examples of subgoals synthesized by SuSIE.** A Comparison between the pre-trained diffusion model initialized from InstructPix2Pix (**Ours**) and random initialization on BridgeData. Each row is a trajectory from a holdout in-distribution validation set, where the objects and environments are all seen but the particular trajectory and language label are not. The fine-tuned model consistently generates better subgoals.

599 SuSIE is able to complete instructions with a significantly higher success rate than LCBC, outper-  
 600 forming prior methods on this benchmark, including both the reimplementations of the closest prior  
 601 approach, UniPi. Concretely, we observe more than about 20% improvement in the success rates for  
 602 completing the first and second language tasks in the chain, and approximately 10% improvement  
 603 for the remaining tasks. This indicates that SuSIE is able to produce useful subgoals that enable the  
 604 low-level policy to accomplish tasks in this novel environment.

605 **Real-world results.** We present performance of real-world evaluations in Table 2. Observe that in  
 606 Scene A, SuSIE achieves the highest success rate on all three tasks, attaining an average success  
 607 rate of 73% which improves over RT-2-X by **69%**. In Scene B, SuSIE again outperforms other prior  
 608 approaches on the two tasks, successfully grounding both the novel ceramic bowl and the previously  
 609 seen orange pot. In the most challenging Scene C (unseen domain, unseen objects), SuSIE attains  
 610 a success rate of 45%, outperforming MOO by about **260%**. However, RT-2-X outperforms SuSIE  
 611 in this scene. We believe that the superior performance of RT-2-X compared to SuSIE in Scene C is  
 612 because it is a much larger 55B parameter model, initializes from a proprietary VLM, and is trained  
 613 on much more data — including BridgeData V2, but also a vast quantity of additional tabletop  
 614 manipulation. These differences in the amount of data and parameters put our method, which only  
 615 utilizes BridgeData V2, at quite an unfair advantage against RT-2-X. Nevertheless, SuSIE is still  
 616 able to recognize the novel objects and attain a high absolute success rate of 45%.

#### 617 E.4 Does SuSIE Improve Precision and Low-Level Skill Execution?

618 Our real-world and simulated results clearly demonstrate the efficacy of SuSIE in executing novel  
 619 language commands in a variety of scenarios. In this section, we visualize some evaluation rollouts  
 620 from our experiments in Scene A to understand if SuSIE works merely because it enhances the  
 621 generalization of the policy to semantic changes in the visual observation or if it actually does  
 622 improve the precision of the low-level control by commanding meaningful subgoals. Observe in  
 623 Figure 5 that the RT-2-X policy often produces actions that fail to precisely orient the gripper around  
 624 the target object or close the gripper early. In contrast, policy executions obtained via SuSIE are more  
 625 precise, and execute actions that attempt to match the gripper and object positions to the generated  
 626 subgoal, allowing the policy to succeed at the task.

627 To understand the contribution of the subgoal prediction towards improved precision, we also eval-  
 628 uate an oracle GCBC policy on a subset of tasks. This policy is trained on identical robot data as  
 629 SuSIE; however, we at test time we command the policy with a real image of the completed task,  
 630 which our method does not require. Observe that even then this GCBC oracle fails to accomplish the  
 631 task due to issues with imprecise object localization and untimely gripper closing. Corroborated by  
 632 numerical results in Table 3, these experiments validate our claim that utilizing subgoal prediction  
 633 is crucial for enabling precise low-level skill execution and control.

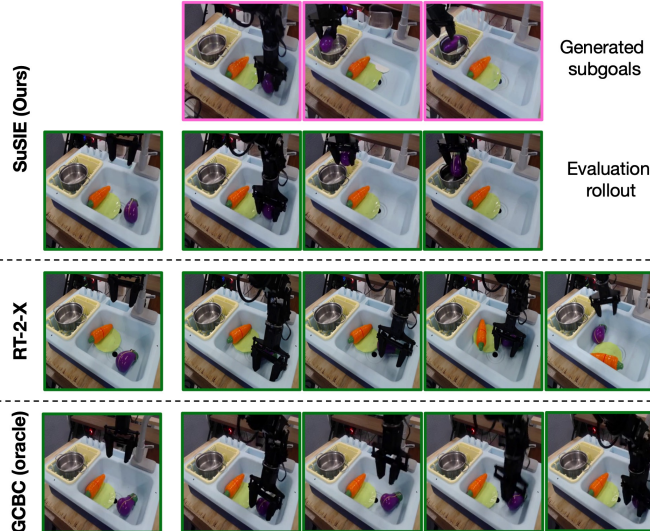


Figure 5: Visualizing rollouts from SuSIE, RT-2-X, and oracle GCBC. While RT-2-X and oracle GCBC often fail to precisely localize or grasp the object, generated subgoals from the image editing model in SuSIE guide the low-level controller precisely, improving low-level skill execution with novel language commands.

## 634 E.5 Is Pre-Training on Internet Data Crucial for Zero-Shot Generalization?

635 Finally, we conduct an experiment to understand if pre-training is crucial for generating meaningful  
 636 subgoals. We train a second image editing model without InstructPix2Pix initialization, but using  
 637 the same UNet architecture, image autoencoder, and text encoder as InstructPix2Pix. Observe in  
 638 Figure 4 that the pre-trained model consistently generates superior subgoals.

## 639 F Discussion and Future Work

640 We presented a method for robotic control from language instructions that uses pre-training to gener-  
 641 ate subgoals to guide low-level goal-conditioned policy, which is unaware of language. The subgoals  
 642 are generated by an image-editing diffusion model fine-tuned on robot data. This system improves  
 643 both zero-shot generalization to new objects, and the precision of the overall policy, because the  
 644 subgoal model incorporates semantic benefits from pre-training and commands the low-level pol-  
 645 icy to reach more meaningful subgoals. Our experiments show that SuSIE improves over prior  
 646 techniques on the CALVIN benchmark and attains good performance in three different scenes for  
 647 a real-world manipulation task, outperforming language-conditioned behavioral cloning, and often  
 648 outperforming the state-of-the-art, instruction-following approach, RT-2-X, that is trained on more  
 649 than an order of magnitude more robot data.

650 Our method is simple and provides good performance, but it does have limitations that suggest  
 651 promising directions for future work. For instance, the diffusion model and the low-level policy are  
 652 trained separately indicating that the diffusion model itself is also unaware of the capabilities of the  
 653 low-level policy — it is trained on the same dataset, but assumes that anything that is reachable in  
 654 the dataset can also be reached by the policy. We hypothesize that performance can be improved  
 655 by making the diffusion model aware of the low-level policy’s capabilities. More broadly, we found  
 656 the performance of our method to often be bottlenecked by the performance of the low-level policy,  
 657 suggesting that addressing either of these limitations might lead to a more performant method for  
 658 importing Internet-scale knowledge into robotic manipulation.

## 659 G Further Details of Implementation

660 We provide implementation details for SuSIE and the baselines.



## 661 G.1 SuSIE implementation details

### 662 G.1.1 Goal-reaching policy

663 We use a diffusion model for our goal-reaching policy since recent work has shown that diffusion-  
664 based policies can better capture multi-modality in robot data (Chi et al., 2023; Hansen-Estruch  
665 et al., 2023), leading to improved performance across a variety of tasks. In our implementation  
666 (which follows Walke et al. (2023)), the observation and goal image are stacked channel-wise before  
667 being passed into a ResNet-50 image encoder. This image encoding is used to condition a diffusion  
668 process that models the action distribution. We use the DDPM (Denoising Diffusion Probabilistic  
669 Models) objective as introduced by Ho et al. (2020). The diffusion process uses an MLP with 3  
670 256-unit layers and residual connections. Following Chi et al. (2023), rather than predicting a single  
671 action, we predict a sequence of  $k$  actions to encourage temporal consistency. We use an action  
672 sequence length of  $k = 4$ . We use the Adam optimizer (Kingma & Ba, 2015) with a learning  
673 rate of  $3e-4$  and a linear warmup schedule with 2000 steps. We augment the observation and goal  
674 with random crops, random resizing, and color jitter. During training, the goal associated with an  
675 observation is selected by uniformly sampling an observation from a window of future timesteps in  
676 the trajectory. Specifically, we sample a goal from 0-20 steps in the future.

677 At test time, we have several options for how to predict and execute action sequences. Chi et al.  
678 (2023) use receding horizon control, sampling  $k$ -length action sequences and only executing some  
679 of the actions before sampling a new sequence. This strategy can make the policy more reactive.  
680 However we found that the robot behavior was quite jerky as the policy switched between different  
681 modes in the action distribution with each sample. Instead, we use a temporal ensembling strategy  
682 similar to Zhao et al. (2023). We predict a new  $k$ -length action sequence at each timestep and execute  
683 a weighted average of the last  $k$  predictions.

## 684 G.2 Baseline implementation details

### 685 G.2.1 Language-conditioned behavior cloning (LCBC)

686 We use the language-conditioned behavior cloning method from Walke et al. (2023) and Myers et al.  
687 (2023). The instruction is encoded using the MUSE sentence embedding Yang et al. (2019), then the  
688 image observation is encoded using a ResNet-50 with FiLM conditioning on the language encoding  
689 Perez et al. (2017). The output is passed into a fully connected policy network with 3 256-unit layers  
690 to produce the action. We use the Adam optimizer Kingma & Ba (2015) with a learning rate of  $3e-4$   
691 and a linear warmup schedule with 2000 steps. We augment the observation and goal with random  
692 crops, random resizing, and color jitter.

### 693 G.2.2 UniPi

694 UniPi (Du et al., 2023) trains a video diffusion model,  $p_{\theta}(\tau|s_0, l)$  to generate a sequence of frames  
695 given a language command and an initial frame. The original paper employs the model architecture  
696 from Imagen Video (Ho et al., 2022a,b). To achieve higher resolution and longer videos for their  
697 real-world results, the authors leverage a 1.7B 3D U-Net and four pre-trained super-resolution mod-  
698 els from Imagen Video, with 1.7B, 1.7B, 1.4B, and 1.2B parameters, respectively. Since the original  
699 models and codes are not publicly available, we tried to replicate their approach in two different  
700 ways.

701 **UniPi (ours).** We implemented a 3D U-Net video diffusion model, following Ho et al. (2022b,a),  
702 combining UniPi’s first-frame conditioning. Due to limited computes, we did not train spa-  
703 tial/temporal super-resolution models; instead, we trained a 3D U-Net-based diffusion model  
704 to directly generate images with a resolution of  $128 \times 128$ . The model includes 4 residual  
705 blocks, with (input channels, output channels) as follows: (64, 64), (64, 128), (128, 256), and  
706 (256, 640). The model is trained to produce the trajectory with a fixed horizon of 10 frames  
707  $\tau_t = \{s_t, s_{t+1}, \dots, s_{t+9}\}$  conditioned on the current frame  $s_t$  and language command. We used  
708 a frozen pre-trained CLIP (Radford et al., 2021) encoder to obtain the language embeddings.

709 **UniPi (HIP, Ajay et al. (2023b))** For the second approach, we followed the UniPi replication in  
710 Ajay et al. (2023b). We trained a latent video diffusion model from PVDM (Yu et al., 2023a),  
711 building upon the codebase <https://github.com/sihyun-yu/PVDM> where we added first

712 frame conditioning. We first trained the video autoencoder to project video of size  $16 \times 128 \times 128$   
713 into latent representation, followed by training a PVDM-L model that uses a 2D U-Net architecture.  
714 We used a Flan-T5- Base (Chung et al., 2022) encoder to obtain the language embeddings.

715 **Data and training details.** To incorporate knowledge from internet data into video models, we  
716 utilize Ego4D (Grauman et al., 2022), a large-scale human egocentric video dataset with language  
717 annotations. For UniPi (ours), we first pre-trained the video model on Ego4D for 270K steps, and  
718 fine-tuned it on the robotics dataset, CALVIN for the simulation and BridgeData v2 for the real  
719 world, for additional 200K steps. We use a batch size of 4 during the training. For UniPi (HIP), we  
720 jointly trained a single model on all Ego4D, BridgeData v2, and CALVIN dataset at the same time.  
721 The autoencoder was trained for 85K steps, and the PVDM-L model was trained for 200K steps.  
722 We use a batch size of 8 during the training.

723 **Inverse model and test time control.** To extract actions from generated videos, we trained an in-  
724 verse dynamics model  $\pi_\phi(\cdot | s_t, s_{t+1})$  to predict the action from two adjacent frames. We employed  
725 the same architecture as our GCBC policy described in Section D.2 and set the goal horizon  $k$  to 1.  
726 During test time, given the current observation  $s_t$  and the language command  $l$ , we synthesize  $H$   
727 image frames from the video model and apply the inverse dynamics model to obtain the correspond-  
728 ing  $H - 1$  actions. The predicted actions are executed, and we generate a new video from  $s_{t+H-1}$   
729 and repeat the process until it reaches the maximum episode step.

730 **Generated videos.** While the quality of the video model trained on the simulation dataset is good  
731 enough for solving the tasks on the CALVIN benchmark as shown in Table 1, we found that it is  
732 nontrivial to obtain a high-quality generation for the real-world dataset. Additionally, sampling the  
733 video model of UniPi to rollout a real robot is extremely time-consuming. Therefore, we evaluated  
734 UniPi only in simulations.

### 735 G.2.3 MOO

736 MOO (Stone et al., 2023) utilizes a mask to represent the target objects and incorporates it as an  
737 additional channel in the observation. Specifically, they train a language-conditioned policy that  
738 takes a 4-channel image and a language command as inputs. To acquire the mask for target objects,  
739 the Owl-ViT (Minderer et al., 2022) detector is employed. This detector is an open-vocabulary  
740 object detection model, pre-trained on internet-scale datasets, and it is used to extract the bounding  
741 boxes of the objects of interest from the image. For tasks like "move X to Y," MOO calculates the  
742 bounding box for X, representing the object of interest, and Y, indicating the target place. A mask  
743 is then created where the pixel at the center of the predicted bounding box is assigned a value of 1.0  
744 for X and 0.5 for Y.

745 **Extracting object entities from BridgeData V2 language annotations.** In order to obtain the  
746 mask, it is necessary to extract the entities corresponding to the object of interest, denoted as X, and  
747 the target place, Y, from the language command. In MOO’s original paper, the authors assume that  
748 the language in their dataset is structured in a way that facilitates the easy separation of X and Y.  
749 Specifically, they employ a dataset that exclusively consists of language annotations such as "pick  
750 X," "move X near Y," "knock X over," "place X upright," and "place X into Y."

751 Given that the language annotations in BridgeData v2 are diverse and unstructured, it is challenging  
752 to naively extract X and Y. We utilized the the API of OpenAI’s `gpt-3.5-turbo-instruct`  
753 model to extract the object of interest and the target place (if any) from the language annotations,  
754 and input them into Owl-ViT to create masks. We then train a mask conditioned LCBC policy using  
755 the same architecture as described in Section G.2.1. Following the original work, we removed X and  
756 Y from the prompt and replaced the word X with "object of interest" and the word Y with "target  
757 place". For example, given a language prompt "put the eggplant in the pot", we use a modified  
758 prompt "put object of interest in target place" as the input to the policy during both training and test  
759 time.

760 **Test time.** During test time, we use oracle masks annotated from the initial camera observations of  
761 each test trial. To enable this, we build a simple interface on the robot machine, allowing the tester  
762 to create the masks by clicking on the initial camera image at the beginning of each test trial.