

# Skill Acquisition by Instruction Augmentation on Offline Datasets

Anonymous Author(s)

Affiliation

Address

email

1       **Abstract:** In recent years, much progress has been made in learning robotic ma-  
2       nipulation policies that follow natural language instructions. Commonly, such  
3       methods learn from a corpora of robot-language data that was either collected  
4       with specific tasks in mind or expensively re-labelled by humans with rich lan-  
5       guage descriptions in hindsight. Recently, large-scale pretrained vision-language  
6       models like CLIP have been applied to robotics in the form of learning repre-  
7       sentations and planners. Can these pretrained models also be used to cheaply  
8       impart internet-scale knowledge onto offline datasets, providing access to skills  
9       that were not reflected in ground truth labels? To accomplish this, we introduce  
10       **Data-driven Instruction Augmentation for Language-conditioned control (DIAL):**  
11       we utilize semi-supervised language labels leveraging the semantic understanding  
12       of CLIP to propagate knowledge onto large datasets of unlabelled demonstration  
13       data and then train language-conditioned policies on the augmented datasets. This  
14       method enables cheaper acquisition of useful language descriptions compared to  
15       expensive human labels, allowing for more efficient label coverage of large-scale  
16       datasets. We apply DIAL to a challenging real-world robotic manipulation domain,  
17       enabling imitation learning policies to acquire new capabilities and generalize to  
18       60 novel instructions unseen in the original dataset.

## 19   1 Introduction

20   Recent advances in decision making have combined data-driven policies with language models  
21   to enable control policies that respond to natural language instructions, an important capability  
22   for practical adoption of general robots in the real world. A popular method used to accomplish  
23   such language-controlled policies is behavioral cloning (BC) [16, 23, 1], which commonly acquires  
24   language labels in two ways: i) using pre-defined tasks where the task descriptions are provided at the  
25   time of data collection or ii) using cheap unstructured data collect like play data [21, 22] paired with  
26   rich language labels provided by humans in hindsight. Both of these options have major drawbacks,  
27   as pre-defining task instructions prior to data collection may limit data diversity, while hindsight  
28   relabelling is expensive when applied at scale.

29   On the other hand, large-scale pretrained language models (LLMs) and vision-language models  
30   (VLMs) have seen increased adoption due to their ability to leverage internet-scale data to augment  
31   or even replace traditionally engineered parts of robot control systems, such as representation for  
32   perception [27, 31], as task representation for language [16, 20], or as planners [1, 15]. We seek to  
33   apply pretrained VLMs to the datasets themselves: can we use VLMs for *instruction augmentation*,  
34   where we relabel existing offline trajectory datasets with additional language instructions?

35   In this work, we provide an analysis of using instruction augmentation with VLMs to weakly relabel  
36   offline control datasets. We demonstrate this method on a challenging real-world robotic control  
37   domain, showing that instruction augmentation allows policies to acquire understanding of skills not  
38   contained in the original task labels, enabling generalization to 60 novel task instructions. We find  
39   that instruction augmentation with VLMs is especially important for generalizing to skills requiring  
40   understanding of spatial semantic concepts.

41 Our core contributions are as follows:

- 42 • We introduce **Data-driven Instruction Augmentation for Language-conditioned control**  
43 (DIAL) by using CLIP to label offline demonstrations for policy learning
- 44 • We study the sensitivity of policy performance to increasing instruction label noise
- 45 • We show the benefits of combining instruction augmentation predictions with existing labels
- 46 • We demonstrate the scalability of the method to a challenging real-world robotic task

## 47 2 Related Work

48 **Language-instruction following in Robotics** Language-instruction following agents have been  
49 extensively explored in the reinforcement learning setting [19]. Recent advances in deep learning  
50 with large amounts of data has led to works following natural language for robotic manipulations.  
51 Latent Motor Control (LMP) [21] learns hierarchical goal-conditioned policies. Subsequent Language  
52 from Play (LfP) [20] uses language goals provided by large dataset of hindsight human labels on  
53 robotic play data. Similarly, LAVA [22] uses crowd-sourced hindsight labels on diverse play data for  
54 table-top object rearrangements. In contrast, our method does not rely on crowd-sourced language  
55 labels at scale, but instead focuses on collecting just a modest amount of language labels and then  
56 using a learned model to provide weak hindsight labeling of the rest of the data.

57 **Learned Language-conditioned Reward Functions** Prior works have investigated using demon-  
58 strations with language annotations to learn language-conditioned reward functions for utilization in  
59 downstream online [3, 14, 12] or offline RL [26, 8]. The complexity of the language instructions range  
60 from templated language in small-scale environments to crowd-sourced language annotations in real  
61 robotics or open-ended environments such as Minecraft. LOReL [26] learns a reward function from  
62 offline datasets of robot interactions with crowd sourced annotations using a convolutional neural  
63 network trained from scratch combined with a pretrained DistilBERT sentence embedding [30] using  
64 the binary cross entropy. MineCLIP [12] fine-tunes CLIP [29] encoders using a contrastive loss on a  
65 large offline dataset of Minecraft videos and optimizes a language-conditioned control policy through  
66 online RL. While their learned reward function can be used to train agents specifically on novel task  
67 instructions, it requires an expensive and sample-inefficient stage of online RL, which is not tractable  
68 in the real world. A frozen CLIP vision and text encoders has also been used as a baseline method  
69 for imitation learning [24] in the simulated robotic manipulation CALVIN benchmark [25]. Our  
70 approach fine-tunes CLIP on our *real* robot offline dataset and is used for instruction augmentation for  
71 a behavior cloning agent, instead of directly using the CLIP model as a reward model and optimizing  
72 an RL agent.

73 **Hindsight Relabeling for Goal-conditioned Reinforcement Learning** The relabeling approach  
74 for goal-conditioned reinforcement learning [28] originates from Hindsight Experience Replay (HER)  
75 [2], which relabels the desired goals in a trajectory with achieved goals (hindsight goal) in the same  
76 trajectories to generate positive examples in a sparse reward setting. Relabeling approach has later  
77 been applied to environments where the goals are images [7], task IDs [18], and language instructions  
78 [17, 6, 9]. Early works with templated language goals rely on environment simulators to provide  
79 hindsight labels [17, 6], and more recently [9] uses a learned model. Our work further applies the  
80 relabeling strategy with a learned model that scales to real robot environments.

81 **Semi-supervised Imitation and Offline Reinforcement Learning** Prior works in semi-  
82 supervised imitation learning focuses on labeling missing actions from demonstrations. The approach  
83 of using a small curated dataset to train a model to then label a larger dataset has been explored in  
84 Video PreTraining (VPT) [4]. While VPT uses the small curated dataset to train an inverse dynamics  
85 model (IDM) to label actions, we fine-tuned CLIP [29] on our small dataset with crowd-sourced  
86 natural language annotation in order to relabel the language instructions for a larger dataset of robot  
87 trajectories. While LOReL [26] also applies instruction relabeling to an instruction from another  
88 episode, the relabeling is used to create more *negative* examples for the reward model to train on. In  
89 contrast, our approach creates new *positive* instruction labels for a given trajectory by leveraging an  
90 already fine-tuned VLM, which is used to train a behaviour cloned policy.

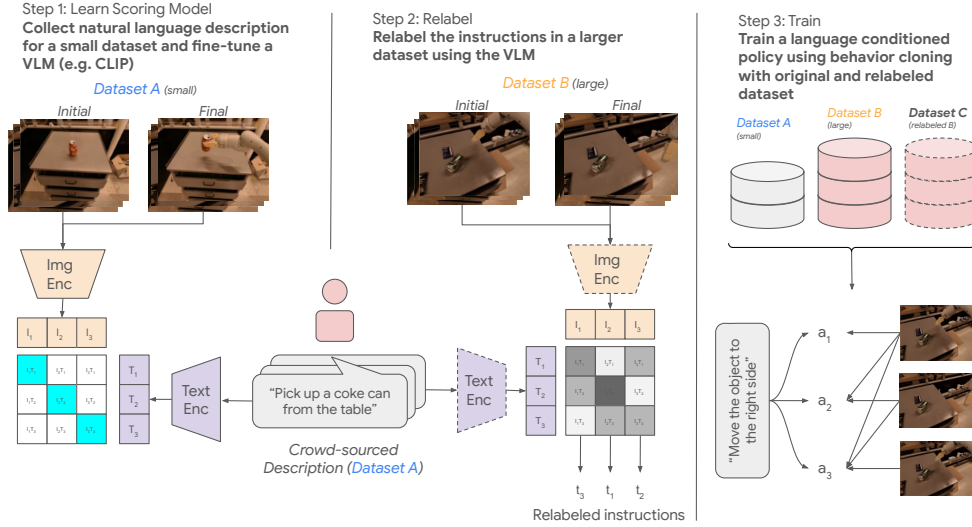


Figure 1: DIAL consists of three steps: (1) Contrastive fine-tuning of a vision-language model (VLM) such as CLIP [29] on small dataset of robot manipulation trajectories with crowd-sourced natural language annotation, (2) using the fine-tuned VLM (in dashed outline) to score and rank the relevance of crowd-sourced annotations against a larger dataset of trajectories to produce novel instruction labels, and (3) training a language-conditioned policy using behavior cloning on the original and relabeled dataset. See Section 3 for more details.

### 91 3 Method

92 In this section, we describe DIAL consisting of three stages: (1) finetuning a VLM’s vision and  
 93 language representation on a small offline dataset of trajectories with crowd sourced episode-level  
 94 natural language description, (2) generating alternative instructions for a larger offline dataset of  
 95 trajectories with the VLM, and (3) learning a language-conditioned policy via behavior-cloning on  
 96 this augmented offline data.

#### 97 3.1 Finetuning Vision-Language Model Representations on Offline Dataset

98 Given an offline dataset of  $N$  trajectories  $[\tau_1, \dots, \tau_N]$ ,  $\tau_n = ((s_0^n, a_0^n), (s_1^n, a_1^n), \dots, (s_T^n))$ , we  
 99 collect a corresponding natural language description  $l^n$  for the  $n$ -th episode describing what the  
 100 robot agent did in the episode via crowd-sourcing. When producing these descriptions, the crowd-  
 101 sourced evaluators observe the first frame,  $s_0$ , and last frame,  $s_T$ , from the agent’s first-person  
 102 view. We refer to these instructions as *hindsight instructions*. Together, we denote the first dataset  
 103  $\mathcal{D}_A = [(\tau_1, l_1), \dots, (\tau_N, l_N)]$  as the paired trajectories and crowd-sourced labels. Our method then  
 104 fine-tunes a vision and language model representation on  $\mathcal{D}_A$ .

105 Motivated by promising results of CLIP in robotics in prior works [31, 24], our instantiation of DIAL  
 106 uses CLIP [29] for both instruction augmentation and task representation; nonetheless, other VLMs or  
 107 captioning models could also be used to propose instruction augmentations. Given a batch of  $B$  initial  
 108 state  $s_0$ , final state  $s_T$ , and hindsight instruction  $l$  tuple, the model is trained to predict which of the  
 109  $B^2$  (initial-final state, hindsight instruction) pairs co-occurred. We use CLIP’s Transformer-based text  
 110 encoder  $T_{enc}$  to embed the crowd-sourced instruction to a latent space  $z_l^n = T_{enc}(l^n) / \|T_{enc}(l^n)\| \in$   
 111  $\mathbb{R}^d$  and CLIP’s Vision Transformer-based (ViT) [11] image encoder  $I_{enc}$  to embed the initial and final  
 112 state, and further concatenate these two embeddings and pass through fully connected neural network  
 113  $f_\theta$ , producing the vision embedding  $z_s^n = f_\theta([I_{enc}(s_0^n); I_{enc}(s_T^n)]) / \|f_\theta([I_{enc}(s_0^n); I_{enc}(s_T^n)])\| \in$   
 114  $\mathbb{R}^d$ .  $B^2$  similarity logits are formed by applying dot product across all state-instruction pairs, and a  
 115 symmetric cross entropy loss term is calculated by applying softmax normalization with temperature

116  $\alpha$  across the states and across the text:

$$\mathcal{L}_{CLIP} = - \left[ \sum_{n=1}^B \log \left( \frac{e^{z_l^n \cdot z_s^n / \alpha}}{\sum_{k=1}^B e^{z_l^k \cdot z_s^n / \alpha}} \right) + \sum_{n=1}^B \log \left( \frac{e^{z_l^n \cdot z_s^n / \alpha}}{\sum_{k=1}^B e^{z_l^n \cdot z_s^k / \alpha}} \right) \right] \quad (1)$$

117 **3.2 Instruction Augmentation on Offline Datasets**

118 We are also given a much larger offline dataset of  
 119  $M \gg N$  trajectories  $[\hat{\tau}_1, \dots, \hat{\tau}_M]$ , where  $\hat{\tau}_m =$   
 120  $([(\hat{s}_0^m, \hat{a}_0^m), (\hat{s}_1^m, \hat{a}_1^m), \dots, (\hat{s}_T^m)])$ . These trajec-  
 121 tories may be collected from human teleoperated  
 122 demonstrations on a wide variety of tasks [1], or  
 123 from episodes from unstructured robotic “play”  
 124 collection [21]. In the first scenario, we may  
 125 have access to the original *foresight instructions*,  
 126  $\hat{l}^m$ , given to the human teleoperators to perform  
 127 the  $m$ -th demonstration episode, while in the lat-  
 128 ter case there are no associated instructions with  
 129 the play episodes. Assuming that we do have the  
 130 foresight instructions, we denote this larger of-  
 131 fline dataset as  $\mathcal{D}_B = [(\hat{\tau}_1, \hat{l}_1), \dots, (\hat{\tau}_M, \hat{l}_M)]$ .

132 We use the fine-tuned VLM model to propose al-  
 133 ternative natural language instructions  $\tilde{l}^m$  for the  
 134 trajectory  $\hat{\tau}_m$  to augment the foresight/absent in-  
 135 structions in  $\mathcal{D}_B$ . Our specific instantiation of  
 136 DIAL uses the fine-tuned CLIP text encoder to  
 137 independently embed the crowd-sourced natural  
 138 language instructions from the first stage, i.e.  
 139  $\tilde{l}^m \in L = \{l^1, \dots, l^N\} \sim \mathcal{D}_A$  and store them:

$$\{z_l^1, \dots, z_l^N\} = \{T_{enc}(l^1), \dots, T_{enc}(l^N)\}$$

140 Similarly, we use the fine-tuned CLIP image  
 141 encoder and MLP fusion to embed the initial  
 142 and final observations from the second dataset:

$$\{\hat{z}_s^1, \dots, \hat{z}_s^M\} = \{f_\theta([I_{enc}(\hat{s}_0^1); I_{enc}(\hat{s}_T^1)]), \dots, f_\theta([I_{enc}(\hat{s}_0^M); I_{enc}(\hat{s}_T^M)])\}$$

143 With these embeddings pre-computed, we can retrieve the most likely candidates using  $k$ -Nearest  
 144 Neighbors [13] with cosine similarity between the vision-language embedding pairs  $d(z_l^n, \hat{z}_s^m) =$   
 145  $\frac{z_l^n \cdot \hat{z}_s^m}{\|z_l^n\| \|\hat{z}_s^m\|}$  as the metric. The resulting top- $k$  candidate instructions  $\{\tilde{l}_1^m, \dots, \tilde{l}_k^m\}$  for each trajectory  $\hat{\tau}_m$   
 146 is used to construct the *relabelled* dataset  $\mathcal{D}_C = [(\hat{\tau}_1, \tilde{l}_1^1), \dots, (\hat{\tau}_1, \tilde{l}_k^1), \dots, (\hat{\tau}_M, \tilde{l}_1^M), \dots, (\hat{\tau}_M, \tilde{l}_k^M)]$ .  
 147 Figure 2 visualizes the three datasets generated.

148 The hyperparameter  $k$  trades off precision and recall of the relabelled dataset. A smaller  $k$  will  
 149 return mostly relevant candidate instructions, while a larger  $k$  value can recall a broader spectrum of  
 150 potential hindsight descriptions for the episode at the expense of introducing irrelevant instructions.  
 151 We will investigate the effects of  $k$  in Section 5 on the downstream policy performance.

152 **3.3 Learning Language Conditioned Policies with Behaviour Cloning**

153 Given a dataset of robot trajectories and corresponding augmented language instructions, we can train  
 154 a language-conditioned control policy with Behavior Cloning (BC). While instruction augmented  
 155 offline datasets can be used by any downstream language-conditioned policy learning method such as  
 156 offline RL or BC, we limit our work to the conceptually simpler BC in order to focus our analysis on  
 157 the importance of instruction augmentation.

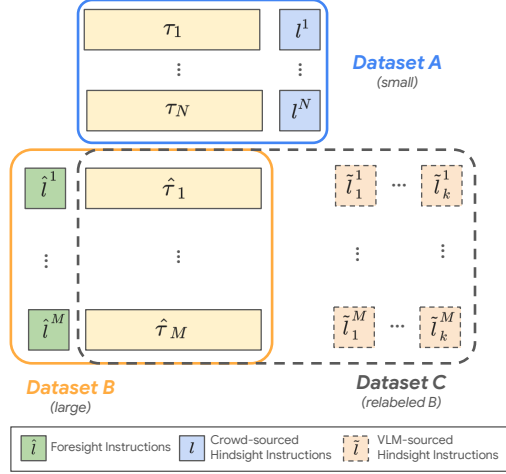


Figure 2: The construction of datasets: Dataset A ( $\mathcal{D}_A$ ) (blue) consists of the  $N$  trajectories  $\{\tau_n\}_{n=1}^N$  labeled with crowd-sourced hindsight instructions  $\{l^n\}_{n=1}^N$  describing what the robot agent performed in the episode. Dataset B ( $\mathcal{D}_B$ ) (yellow) consists of a much larger set of trajectories,  $\{\hat{\tau}_m\}_{m=1}^M$  generated by foresight instructions  $\{\hat{l}^m\}_{m=1}^M$  without hindsight labels. Dataset C ( $\mathcal{D}_C$ ) (black, dashed) contains Dataset B trajectories relabelled with VLM-sourced hindsight instruction(s)  $\{\tilde{l}_1^m, \dots, \tilde{l}_k^m\}_{m=1}^M$ .



Figure 3: (a) A mobile manipulator robot performs a variety of manipulation tasks with various objects and cabinet drawers in an office kitchen environment. (b) An example of some of the kitchen objects found in the demonstration dataset. (c) The mobile manipulator robot receives RGB images from an over-the-shoulder camera and uses a 7 DoF arm with parallel-jaw grippers.

## 158 4 Experimental Setup

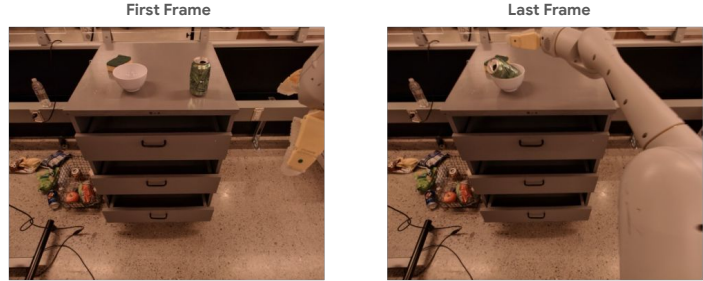
### 159 4.1 Environment, Robot, and Datasets

160 We implement DIAL in a challenging real-world robotic manipulation setting in a kitchen environment  
 161 similar to SayCan [1]. We focus on the practically-motivated setting where a dataset of teleoperated  
 162 demonstrations is available, collected for downstream imitation learning [1, 16]. An [Everyday Robots](#)  
 163 robot [33], a mobile manipulator with RGB observations, is placed in an office kitchen to interact  
 164 with common objects using concurrent [34] continuous closed-loop control from pixels, as shown in  
 165 Figure 3. The robot uses parallel-jaw grippers, an over-the-shoulder RGB camera, and a 7 DoF arm.  
 166 We collect a large-scale dataset of over 80,000 robot trajectories via human teleoperation ( $\mathcal{D}_B$  in  
 167 Section 3.2), where teleoperators perform 551 unique tasks motivated by common manipulation skills  
 168 and objects in a kitchen environment [1]. Afterwards, we leverage crowd-sourced human annotators  
 169 to label 2,800 robot trajectories with two possible hindsight instructions each, resulting in a total of  
 170 5,600 unique episodes with crowdsourced captions ( $\mathcal{D}_A$  in Section 3.1). Human annotators are shown  
 171 the first and last frame of the episode and asked to provide a free-form text description describing  
 172 how a robot should be commanded to go from the start to the end.

### 173 4.2 Instruction Augmentation and Policy Training

174 After finetuning a CLIP model on 5,600 annotated episodes using the procedure described Section  
 175 3.1, we then perform instruction augmentation on the 80,000 demonstrations which do not  
 176 contain hindsight instructions ( $\mathcal{D}_C$  as in Section 3.2). By increasing the number  $k$  of instruction  
 177 augmentations applied to each episode, we produce three instructed augmented datasets: 80,000  
 178 relabeled demonstrations ( $k = 1$ ), 240,000 relabeled demonstrations ( $k = 3$ ), and 800,000 relabeled  
 179 demonstrations ( $k = 10$ ).

180 When increasing  $k$ , the augmented datasets become larger but the proposed instructions may become  
 181 increasingly irrelevant or inaccurate. To measure how instruction augmentation accuracy changes as  
 182 we increase  $k$ , we ask human labelers to rate whether the proposed captions are factually accurate  
 183 descriptions of a given episode. We show an example of predicted instruction augmentations in  
 184 Figure 4 and measure the accuracy of predicted instructions in Table ??.



| Instruction Augmentation Prediction by CLIP  | Accurate? |
|--|-----------|
| #1: pick up the green can and place it in the bowl which is at the left side of the table  | ✓         |
| #2: lift green can from table and place it in white cup                                    | ✓         |
| #3: pick up the green can which is close to the water bottle and place it in the bowl      | ✗         |
| #4: place green can into the plastic white bowl  | ✓         |
| #5: pick the green can from the bottom right of the table and place it into the white bowl | ✓         |
| #6: pick up the silver can and place it in the white bowl                                  | ✗         |
| #7: bring the blue can and place it into white paper bowl                                  | ✗         |
| #8: pick up the green can from the bottom left side of the table                           | ✗         |
| #9: pick up the green can from the bottom side of the table and drop it into bowl          | ✓         |
| #10: pick up the red bull can and drop it in the white bowl                                | ✗         |

Figure 4: The top 10 proposed instruction augmentations for a single episode with original foresight instruction place green can in white bowl. In some cases, the predicted captions provide additional semantic information such as describing the location of the can or the material of the bowl.

| Category  | Instruction Samples   |
|-----------|---|
| Spatial   | ['knock down the right soda', 'raise the left most can', 'raise bottle which is to the left of the can']                                |
| Rephrased | ['pick up the apple fruit', 'liftt the fruit' [sic], 'lift the yellow rectangle']   |
| Semantic  | ['move the lonely object to the others', 'push blue chip bag to the left side of the table', 'move the green bag away from the others'] |

Table 1: Sample novel instructions in each evaluation category. Spatial tasks focus on tasks involving Spatial relationships, Rephrased tasks contain tasks that directly map to a foresight skill, and Semantic tasks describe semantic concepts not contained in the relabeled or original datasets. In total, there are 60 instructions across the three categories.

185 Using these various instruction augmented datasets, we train vision-based language-conditioned  
 186 behavior cloning policies similar to the formulation in BC-Z [16], as described in Section 3.3.  
 187 Compared to BC-Z, we use a larger Transformer [32] based backbone instead of ResNet18 and  
 188 use a CLIP language encoder instead of a Universal Sentence Encoder [5]. Nonetheless, we treat  
 189 the behavior cloning policy as an independent component of our method and focus on studying  
 190 instruction augmentation methods; we do not explore different policy architectures or losses in this  
 191 work.

### 192 4.3 Evaluation

193 In contrast to prior works [16] on instruction following, we focus our evaluation only on *novel*  
 194 *instructions unseen during training*. To source these novel instructions, we crowd-source instructions  
 195 from a different set of humans than the original dataset labelers and filter out any instructions already  
 196 contained in either the instruction augmentation process in Section 3.2 or in the original set of  
 197 551 foresight tasks in Section 4.1; in total, we sample 60 novel evaluation instructions. While  
 198 these evaluation instructions were not curated with specific properties in mind, after sourcing these  
 199 instructions we organize them into various semantic categories to allow for more detailed analysis of  
 200 qualitative policy performance; some examples are shown in Table 1.

- 201 1. **Spatial**: 40 tasks focusing on instructions involving reasoning about spatial relationships.  
 202 For example, this includes specifying an object’s initial position relative to other objects in  
 203 the scene.
- 204 2. **Rephrased**: 10 tasks which are linguistic re-phrasings of the original 551 foresight tasks.  
 205 For example, this includes referring to sodas and chips by their colors instead of their brand  
 206 name.
- 207 3. **Semantic**: 10 tasks which encompass skills not contained in the original dataset. For  
 208 example, this includes the instruction of moving objects away from all other objects, since  
 209 the original dataset only contains trajectories of moving objects towards other objects.

## 210 5 Experimental Results

### 211 5.1 Does using DIAL improve policy performance on unseen tasks?

212 We investigate to what extent a behavior-cloned policy can be successfully learned from instruction  
 213 augmented datasets, even when some relabeled instructions are potentially inaccurate. We use *all*  
 214 available datasets containing foresight labels (FS), ground-truth hindsight labels (GT), and instruction  
 215 augmentation (IA). We vary the amount of instruction augmentation by setting the hyperparameter  
 216  $k = \{1, 3, 10\}$ , resulting in additional 80k to 800k trajectory-instruction pairs. As baselines, we also  
 217 consider training policies *without* instruction augmentation, i.e. only on FS, and on (FS + GT).

218 Table 2 summarises the evaluation results across three categories of novel tasks. Additional baselines  
 219 we consider in Table 5 include methods that perform instruction augmentation without visual context.  
 220 We find that only instruction augmentation using CLIP is able to perform well at novel “Spatial” tasks  
 221 that require visual understanding and “Semantic” tasks that introduce generalizing to semantic skills  
 222 not contained in the original foresight instructions.

### 223 5.2 Does using DIAL for *unlabeled* datasets improve policy performance on unseen tasks?

224 Starting with a dataset of 5,600 trajectories with crowd-sourced hindsight labels, we apply different  
 225 amounts of instruction augmentation onto a dataset of 80,000 trajectories that do not have any  
 226 hindsight language labels. This experiment emulates the practical setting of when a large amount of  
 227 unstructured trajectory data is available but hindsight labels are expensive to collect, such as robot  
 228 play data [10, 21, 22]. We find that training on the instruction augmented trajectories increases  
 229 performance on a set of 60 sampled novel instructions not seen in the original hindsight label set, as  
 230 shown in Table 3. However, overall performance suffers when increasing the number of augmented  
 231 instructions from  $k = 3$  to  $k = 10$ , suggesting there is some limit to how much label inaccuracy the  
 232 language-conditioned policies can tolerate.

| Instruction Augmented Dataset Properties |                |                   | Evaluation on Novel Instructions |              |              |              |
|--|----------------|-------------------|----------------------------------|--------------|--------------|--------------|
| Episodes w/ FS                           | Episodes w/ GT | Episodes w/ IA    | Spatial                          | Rephrased    | Semantic     | Overall      |
| 80k                                      | 0              | 0                 | 33.3%                            | 62.5%        | 10.0%        | 35.0%        |
| 80k                                      | 5600           | 0                 | 45.2%                            | <b>87.5%</b> | 0.0%         | 43.3%        |
| 80k                                      | 5600           | 80k ( $k = 1$ )   | 59.5%                            | 75.0%        | 30.0%        | <b>56.7%</b> |
| 80k                                      | 5600           | 240k ( $k = 3$ )  | <b>64.3%</b>                     | 50.0%        | 30.0%        | 55.0%        |
| 80k                                      | 5600           | 800k ( $k = 10$ ) | 35.7%                            | 50.0%        | <b>40.0%</b> | 35.0%        |

Table 2: Combining episodes with foresight labels of the structured tasks attempted during data collection (FS) with groundtruth crowd-sourced hindsight instructions (GT) with an increasing amount  $k$  of instruction augmentation (IA). DIAL performs the best at challenging “Spatial” tasks.

| Instruction Augmented Dataset Properties |                   |             | Evaluation on Novel Instructions |              |              |              |
|--|-------------------|-------------|----------------------------------|--------------|--------------|--------------|
| Episodes w/ GT                           | Episodes w/ IA    | IA Accuracy | Spatial                          | Rephrased    | Semantic     | Overall      |
| 5600                                     | 0                 | N/A         | 23.8%                            | 37.5%        | 0.0%         | 21.7%        |
| 5600                                     | 80k ( $k = 1$ )   | 68.0%       | 50.0%                            | <b>75.0%</b> | 0.0%         | 45.0%        |
| 5600                                     | 240k ( $k = 3$ )  | 65.3%       | <b>52.4%</b>                     | 50.0%        | <b>20.0%</b> | <b>46.7%</b> |
| 5600                                     | 800k ( $k = 10$ ) | 57.0%       | 38.1%                            | 62.5%        | 10.0%        | 36.7%        |

Table 3: Training on groundtruth crowd-sourced hindsight instructions (GT) compared with utilizing increasing the amount  $k$  of instruction augmentation on unlabeled data (IA), with a corresponding decrease in label accuracy. Instruction Augmentation up to  $k = 3$  significantly improves overall novel instruction performance, especially on ‘‘Spatial’’ tasks requiring visual reasoning.

| Model                 | Task Instruction Encoder | Evaluation on Novel Instructions |              |              |              |
|-----------------------|--------------------------|----------------------------------|--------------|--------------|--------------|
|                       |                          | Spatial                          | Rephrased    | Semantic     | Overall      |
| GT Only               | USE                      | 16.7%                            | 33.3%        | 0.0%         | 18.6%        |
| GT Only               | FT CLIP                  | 23.8%                            | 37.5%        | 0.0%         | 21.7%        |
| FS + GT               | Pretrained CLIP          | 42.9%                            | <b>75.0%</b> | 0.00%        | 40.0%        |
| FS + GT               | FT CLIP                  | 42.9%                            | <b>75.0%</b> | 20.0%        | 41.7%        |
| FS + GT + IA, $k = 1$ | USE                      | 47.6%                            | 50.0%        | 10.0%        | 43.3%        |
| FS + GT + IA, $k = 1$ | FT CLIP                  | <b>59.5%</b>                     | <b>75.0%</b> | <b>30.0%</b> | <b>56.7%</b> |

Table 4: Comparing downstream policy performance when improving the task representation from USE [5] to Pretrained CLIP [29] to Finetuned CLIP (FT CLIP), as described in Section 3.1. We find that the FT CLIP representation is the best task representation in all dataset settings.

### 233 5.3 Is a VLM good at relabeling also a good task representation?

234 We study whether a VLM fine-tuned for instruction augmentation can also act as a better task  
235 representation for policy learning in the form of a more powerful language embedding. Across the  
236 various groundtruth and relabeled datasets we focus on, we find that Finetuned CLIP is the most  
237 effective task representation, as seen in Table 4. Finetuned CLIP is a good representation not only for  
238 freeform language instructions like those contained in the finetuning dataset in Section 4.2, but also  
239 for structured foresight commands like those contained in Section 4.1.

## 240 6 Conclusion

241 In this work, we introduced DIAL, a method that uses VLMs to label offline datasets for language-  
242 conditioned policy learning. We show that control policies are able to utilize relabeled demonstrations  
243 even when some labels are inaccurate, suggesting that DIAL is able to provide a cheap and automated  
244 option to extract additional semantic knowledge from offline control datasets. As the performance of  
245 internet-scale VLMs improve, we expect that DIAL might work increasingly better on even richer  
246 control settings.



## References

- 248 [1] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, K. Gopalakrishnan,  
249 K. Hausman, A. Herzog, D. Ho, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, E. Jang, R. J. Ruano,  
250 K. Jeffrey, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, K.-H. Lee, S. Levine,  
251 Y. Lu, L. Luu, C. Parada, P. Pastor, J. Quiambao, K. Rao, J. Rettinghouse, D. Reyes, P. Sermanet,  
252 N. Sievers, C. Tan, A. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, S. Xu, and M. Yan.  
253 Do as i can and not as i say: Grounding language in robotic affordances. In *arXiv preprint*  
254 *arXiv:2204.01691*, 2022.
- 255 [2] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,  
256 O. Pieter Abbeel, and W. Zaremba. Hindsight experience replay. *Advances in neural information*  
257 *processing systems*, 30, 2017.
- 258 [3] D. Bahdanau, F. Hill, J. Leike, E. Hughes, A. Hosseini, P. Kohli, and E. Grefenstette. Learning  
259 to understand goal specifications by modelling reward. *arXiv preprint arXiv:1806.01946*, 2018.
- 260 [4] B. Baker, I. Akkaya, P. Zhokhov, J. Huizinga, J. Tang, A. Ecoffet, B. Houghton, R. Sampedro,  
261 and J. Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos.  
262 *arXiv preprint arXiv:2206.11795*, 2022.
- 263 [5] D. Cer, Y. Yang, S.-y. Kong, N. Hua, N. Limtiaco, R. S. John, N. Constant, M. Guajardo-  
264 Cespedes, S. Yuan, C. Tar, et al. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*,  
265 2018.
- 266 [6] H. Chan, Y. Wu, J. Kiros, S. Fidler, and J. Ba. Actrce: Augmenting experience via teacher’s  
267 advice for multi-goal reinforcement learning. *arXiv preprint arXiv:1902.04546*, 2019.
- 268 [7] Y. Chebotar, K. Hausman, Y. Lu, T. Xiao, D. Kalashnikov, J. Varley, A. Irpan, B. Eysenbach,  
269 R. Julian, C. Finn, et al. Actionable models: Unsupervised offline reinforcement learning of  
270 robotic skills. *arXiv preprint arXiv:2104.07749*, 2021.
- 271 [8] A. S. Chen, S. Nair, and C. Finn. Learning generalizable robotic reward functions from"  
272 in-the-wild" human videos. *arXiv preprint arXiv:2103.16817*, 2021.
- 273 [9] G. Cideron, M. Seurin, F. Strub, and O. Pietquin. Self-educated language agent with hindsight  
274 experience replay for instruction following. 2019.
- 275 [10] Z. J. Cui, Y. Wang, N. Muhammad, L. Pinto, et al. From play to policy: Conditional behavior  
276 generation from uncurated robot data. *arXiv preprint arXiv:2210.10047*, 2022.
- 277 [11] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,  
278 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for  
279 image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 280 [12] L. Fan, G. Wang, Y. Jiang, A. Mandlkar, Y. Yang, H. Zhu, A. Tang, D.-A. Huang, Y. Zhu,  
281 and A. Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale  
282 knowledge. *arXiv preprint arXiv:2206.08853*, 2022.
- 283 [13] E. Fix and J. L. Hodges. Discriminatory analysis. nonparametric discrimination: Consistency  
284 properties. *International Statistical Review/Revue Internationale de Statistique*, 57(3):238–247,  
285 1989.
- 286 [14] J. Fu, A. Korattikara, S. Levine, and S. Guadarrama. From language to goals: Inverse rein-  
287 forcement learning for vision-based instruction following. *arXiv preprint arXiv:1902.07742*,  
288 2019.
- 289 [15] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners:  
290 Extracting actionable knowledge for embodied agents. *arXiv preprint arXiv:2201.07207*, 2022.
- 291 [16] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn. Bc-z:  
292 Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*,  
293 pages 991–1002. PMLR, 2022.
- 294 [17] Y. Jiang, S. S. Gu, K. P. Murphy, and C. Finn. Language as an abstraction for hierarchical deep  
295 reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- 296 [18] D. Kalashnikov, J. Varley, Y. Chebotar, B. Swanson, R. Jonschkowski, C. Finn, S. Levine, and  
297 K. Hausman. Mt-opt: Continuous multi-task robotic reinforcement learning at scale. *arXiv*  
298 *preprint arXiv:2104.08212*, 2021.

- 299 [19] J. Luketina, N. Nardelli, G. Farquhar, J. Foerster, J. Andreas, E. Grefenstette, S. Whiteson,  
300 and T. Rocktäschel. A survey of reinforcement learning informed by natural language. *arXiv*  
301 *preprint arXiv:1906.03926*, 2019.
- 302 [20] C. Lynch and P. Sermanet. Language conditioned imitation learning over unstructured data.  
303 *arXiv preprint arXiv:2005.07648*, 2020.
- 304 [21] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet. Learning  
305 latent plans from play. In *Conference on robot learning*, pages 1113–1132. PMLR, 2020.
- 306 [22] C. Lynch, A. Wahid, J. Tompson, T. Ding, J. Betker, R. Baruch, T. Armstrong, and P. Florence.  
307 Interactive language: Talking to robots in real time. *arXiv preprint arXiv:2210.06407*, 2022.
- 308 [23] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese,  
309 Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations for  
310 robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.
- 311 [24] O. Mees, L. Hermann, and W. Burgard. What matters in language conditioned robotic imitation  
312 learning. *arXiv preprint arXiv:2204.06252*, 2022.
- 313 [25] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard. Calvin: A benchmark for language-  
314 conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and*  
315 *Automation Letters (RA-L)*, 7(3):7327–7334, 2022.
- 316 [26] S. Nair, E. Mitchell, K. Chen, S. Savarese, C. Finn, et al. Learning language-conditioned robot  
317 behavior from offline data and crowd-sourced annotation. In *Conference on Robot Learning*,  
318 pages 1303–1315. PMLR, 2022.
- 319 [27] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation  
320 for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- 321 [28] M. Plappert, M. Andrychowicz, A. Ray, B. McGrew, B. Baker, G. Powell, J. Schneider, J. Tobin,  
322 M. Chociej, P. Welinder, et al. Multi-goal reinforcement learning: Challenging robotics  
323 environments and request for research. *arXiv preprint arXiv:1802.09464*, 2018.
- 324 [29] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,  
325 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision.  
326 In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- 327 [30] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. Distilbert, a distilled version of bert: smaller,  
328 faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 329 [31] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipula-  
330 tion. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.
- 331 [32] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and  
332 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,  
333 30, 2017.
- 334 [33] X Development, LLC. Everyday Robots. <http://www.everydayrobots.com>, 2022. Ac-  
335 cessed: 2022-06-15.
- 336 [34] T. Xiao, E. Jang, D. Kalashnikov, S. Levine, J. Ibarz, K. Hausman, and A. Herzog. Think-  
337 ing while moving: Deep reinforcement learning with concurrent control. *arXiv preprint*  
338 *arXiv:2004.06089*, 2020.

339 **A Appendix**

340 **A.1 Instruction Augmentation Accuracy**

341 As described in Section 4.3, instruction prediction ac-  
 342 curacy may decrease when increasing the number  $k$  of  
 343 instruction augmentations. In Figure 5, we sample 50  
 344 episodes and ask human labelers to assess the predicted  
 345 instruction accuracy as we increase the number of pre-  
 346 dictions produced by CLIP. While the initial predictions  
 347 are correct often, the later predictions are often factually  
 348 inaccurate. The top-20-th instruction prediction is only  
 349 factually accurate 20.0% of the time. An example of the  
 350 top 10 predictions of an episode is shown in Figure 4.

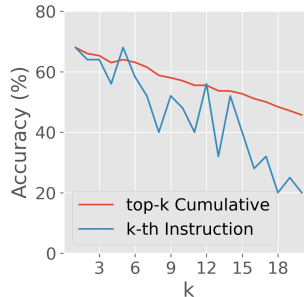


Figure 5: The accuracy of the top 20 instruction augmentation predictions of a sample of 50 episodes that have been relabeled by a Finetuned CLIP model in Section 4.2.

351 **A.2 Additional Experiments**

352 While our proposed method utilizes instruction augmenta-  
 353 tion with pretrained visual-language models, we can also  
 354 attempt to increase the diversity of task instructions with  
 355 non-visual methods. Two potential methods to do this are  
 356 madlibs-style augmentations that replace words in the foresight instructions with synonyms and  
 357 with Gaussian Noise augmentations that add noise with variance=0.05 to the text embeddings of  
 358 foresight instructions. In Table 5, we compare relabeling methods in a setting similar to Section 5.1,  
 359 where we apply relabeling to ground-truth labels from 80,000 episodes with foresight tasks and 5,600  
 360 episodes with groundtruth tasks. We note that while our dataset allows the baseline methods to relabel  
 361 starting from the ground-truth foresight labels, “IA with CLIP” is able to relabel potentially unlabeled  
 362 episodes, a setting that is not possible for the baseline methods.

| Relabeling Method         | Evaluation on Novel Instructions |              |              |              |
|---------------------------|----------------------------------|--------------|--------------|--------------|
|                           | Spatial                          | Rephrased    | Semantic     | Overall      |
| No relabeling             | 33.3%                            | 62.5%        | 10.0%        | 35.0%        |
| Madlibs Text Augmentation | 31.0%                            | <b>87.5%</b> | 20.0%        | 35.0%        |
| Gaussian Noise            | 31.4%                            | 75.0%        | 0.0%         | 30.0%        |
| IA with CLIP, $k = 1$     | 59.5%                            | 75.0%        | 30.0%        | <b>56.7%</b> |
| IA with CLIP, $k = 3$     | <b>64.3%</b>                     | 50.0%        | <b>30.3%</b> | 55.0%        |

Table 5: Comparing instruction augmentation with CLIP (IA) with non-visually grounded ways of relabeling the foresight tasks. We try Madlibs-style text augmentation as well as adding task embedding Gaussian noise. Policies train on foresight labels, groundtruth hindsight labels, and the additional relabeled episodes. While these improve performance on “Rephrased” tasks, they fail to improve performance on other task categories.