

# Masked Inverse Reinforcement Learning for Language Conditioned Reward Learning

Anonymous Author(s)

Affiliation

Address

email

1       **Abstract:** Natural language provides a flexible interface for specifying robot  
2 tasks, but language-conditioned reward learning often assumes that instructions are  
3 unambiguous and directly informative. In reality, human language is frequently  
4 ambiguous — and may specify not just what to do, but also what matters in the  
5 environment. In this work, we propose a method that leverages this duality: we use  
6 large language models (LLMs) to extract state feature-level relevance masks from  
7 language and demonstrations, and train a reward function that is both conditioned  
8 on clarified task language and explicitly invariant to irrelevant parts of the state. We  
9 show that this approach improves generalization and sample efficiency in inverse  
10 reinforcement learning, particularly in settings with ambiguous instructions, distrac-  
11 tor objects, or limited data. Our results highlight that disambiguating language with  
12 contextual demonstrations — and using language to guide both goal inference and  
13 state abstraction — enables more robust reward learning from natural instructions.

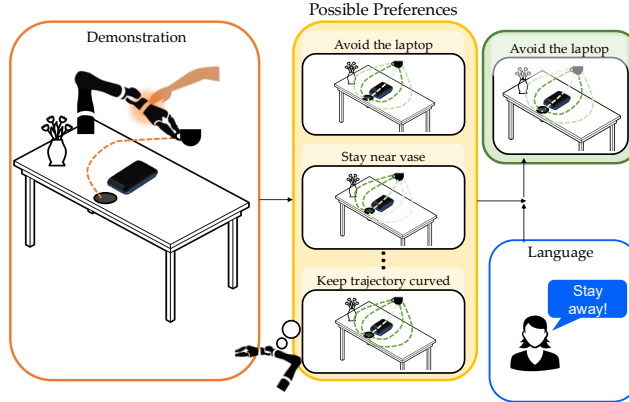
14       **Keywords:** Inverse Reinforcement Learning, Multi-Modal Feedback, Language  
15 Conditioning, Reward Learning

## 16   1 Introduction

17 In robotics, natural language provides a flexible and intuitive interface for specifying the tasks.  
18 However, language-conditioned reward learning typically assumes language instructions are clear  
19 and unambiguous. In practice, human language is often inherently ambiguous – an instruction can  
20 specify not only what the robot should do but also which elements of the environment matter for the  
21 task. Addressing this ambiguity is essential for effective reward learning from limited demonstrations  
22 and generalization to novel tasks or contexts.

23 Language-conditioned reward learning has gained significant interest in recent robotics literature. Fu  
24 *et al.* [1] show that learning a reward model conditioned on language yields behavior that transfers  
25 to novel tasks, whereas directly training a language-conditioned policy was less effective. Poddar  
26 *et al.* [2] learns a latent space that maps language instructions into hidden states to condition the  
27 reward model. Although language is frequently used as an additional modality in robot learning,  
28 existing approaches typically treat language simply as another input to a policy or reward model,  
29 without explicitly structuring the learning around the state features indicated as important by language.  
30 Consequently, these models implicitly infer feature relevance, which can lead to spurious correlations.

31 To address this gap, we propose a method that leverages the duality of language instructions in reward  
32 learning: their ability to specify tasks as well as to indicate relevant environmental state features.  
33 Specifically, our approach uses large language models (LLMs) to extract explicit relevance masks at  
34 the state feature level from language instructions and demonstrations during training. These masks  
35 identify which environmental features are task-relevant according to the provided instruction. Using  
36 these masks, we train a reward function conditioned on language instructions that explicitly ignores



**Figure 1: Overview.** For a robotic task, language can specify not only what to do, but also what matters in the environment. When there are multiple objects (e.g., vase and laptop) and human in the environment, an instruction “stay away from the laptop” states that ‘laptop’ is the only important feature. Even for ambiguous instructions such as “stay away”, when combined with a contextual demonstration (blue trajectory), the instruction can be clarified to include the missing referent, e.g., laptop.

37 irrelevant state features. At inference time, our model can handle not only clear instructions but also  
 38 ambiguous language instructions clarified by a single demonstration per instruction.

39 We introduce *Masked Inverse Reinforcement Learning* (Masked IRL), a framework that integrates  
 40 human demonstrations and language instructions to explicitly guide feature selection for reward  
 41 learning. Masked IRL leverages language-derived masks to dynamically gate relevant features in the  
 42 reward function. For example, given the instruction “stay away from the table” in the scene in Fig. 1,  
 43 our model explicitly ignores laptop-related features and irrelevant end-effector coordinates while  
 44 emphasizing the vertical distance between the robot and table surface. We propose a masking loss  
 45 that penalizes variations in reward predictions resulting from perturbations in state features indicated  
 46 as irrelevant by the language. This builds upon the concept of contextual reliability by Ghosal *et*  
 47 *al.* [3], explicitly training models to identify and ignore spurious or contextually irrelevant features.

48 By explicitly leveraging multimodal human feedback, Masked IRL substantially reduces demonstra-  
 49 tion requirements, improves sample efficiency, and enhances generalization by focusing solely on  
 50 relevant task features. We empirically validate our approach using a PyBullet simulation environment  
 51 with a Franka Emika Panda robotic arm. Our experiments highlight Masked IRL’s effectiveness in  
 52 settings with ambiguous language instructions, distractor objects, and limited demonstration data,  
 53 demonstrating improved data efficiency, robustness, and generalization relative to standard IRL  
 54 methods. In summary, our contributions are:

- 55 • A method using large language models (LLMs) to disambiguate language instructions and  
 56 explicitly extract state-feature relevance masks from instructions paired with demonstrations.
- 57 • **Masked IRL**, an IRL framework that conditions rewards on clarified instructions and  
 58 explicitly enforces invariance to irrelevant state features via a novel masking loss.
- 59 • Empirical validation of Masked IRL’s effectiveness on simulated robotic manipulation  
 60 tasks, demonstrating improved generalization, robustness, and data efficiency compared to  
 61 traditional language-conditioned IRL approaches.

## 62 2 Related Work

### 63 2.1 Reward Learning from Human Feedback

64 Inverse reinforcement learning (IRL) learns reward functions from expert demonstrations. Early  
 65 works [4, 5, 6, 7] have shown promising results in robotics but suffer a trade-off between the number  
 66 of expert demonstrations and identifiability [8, 9], i.e., the required amount of demonstrations to

67 identify the true objective function is huge. One fundamental limitation of IRL is that we can only  
68 train one reward function given a set of demonstrations, thereby requiring  $N$  set of demonstrations  
69 and  $N$  training processes to train  $N$  different reward functions. Bobu *et al.* [10] separates feature  
70 learning and reward learning, and uses human trajectory similarity queries to learn a task-agnostic  
71 feature space. However, they still require multiple demonstration sets for different user preferences  
72 and cannot generalize to unseen preferences. Beyond demonstrations alone, incorporating various  
73 modalities of human feedback (e.g., pairwise trajectory comparisons, language) has been shown  
74 to improve reward learning efficiency or reduce human’s cognitive effort. Reinforcement learning  
75 from human feedback (RLHF) methods [11, 12] use pairwise human preferences to guide reward  
76 learning, but these methods often require thousands of human feedback to learn a single reward  
77 function [13]. Previous works [14, 15, 16] leverage API-based LLMs to generate a reward function  
78 as a code or predict weights on sub-rewards. Yu *et al.* [15] use an LLM as a Reward Translator,  
79 mapping high-level instructions into dense reward functions that standard RL can optimize. Recent  
80 works [2, 17] combine pairwise comparisons with language. Poddar *et al.* [2] highlight the need for  
81 personalized reward learning, arguing that aggregating human preferences can obscure individual  
82 human preferences. Their method learns a variational latent user model that personalizes rewards  
83 to individual users. Yang *et al.* [17] incorporates comparative language feedback, where humans  
84 describe which trajectory is better and why. Their model embeds trajectory-language pairs into a  
85 shared space, enabling iterative refinement of the reward function.

## 86 2.2 Language-Conditioned Learning in Robotics

87 Integrating natural language with robot learning has gained significant interest as a way to bridge  
88 human intent with robots. Recent methods leverage language as a conditioning signal in policy  
89 learning and reward modeling. Fu *et al.* [1] propose a language-conditioned reward learning approach  
90 in which IRL is used to ground language commands, showing that the resulting reward functions  
91 transfer better to novel tasks. In parallel, systems like LILAC [18] allow human operators to provide  
92 online language corrections during task execution. While such approaches have shown promising  
93 results, they often use language merely as an auxiliary input without explicit structure for feature  
94 selection. Language has become an essential modality for training robots, as it enables humans to  
95 specify goals, provide feedback, and guide behavior. One prominent approach is to condition policies  
96 or reward functions on language instructions. Ahn *et al.* [19] introduce the Say-Can framework, which  
97 grounds high-level instructions using a large language model (LLM) and constrains execution using  
98 a value function, allowing robots to follow abstract human commands. Huang *et al.* [20] show that  
99 LLMs can serve as zero-shot planners by generating structured action sequences from instructions,  
100 while Huang *et al.* [21] introduce Inner Monologue, a framework that integrates environment feedback  
101 into LLM planning, significantly improving long-horizon task execution. Incorporating LLMs into  
102 robotic control has also gained traction. Liang *et al.* [22] propose Code-as-Policies (CaP), in which  
103 LLMs generate executable code (Python functions) for robotic policies, allowing for interpretable,  
104 structured control. This approach enables robots to generalize to unseen instructions by modifying  
105 their behavior through high-level program synthesis.

106 Beyond LLM-based planning, recent work has explored language-conditioned reward learning.  
107 Yu *et al.* [15] introduce Language to Rewards, where an LLM parses high-level instructions and  
108 outputs a parametric reward function, bridging natural language and robotic reinforcement learning.  
109 Karamcheti *et al.* [23] propose Voltron, a vision-language model for representation learning that  
110 aligns video frames with text descriptions, facilitating language-driven imitation learning. Hwang *et al.*  
111 [24] learn a success detector or a reward function that understands semantic grounding of robot  
112 motions. Other approaches integrate demonstrations with corrective language feedback to directly  
113 gate task-irrelevant features [18]. In such systems, language helps the robot focus on task-relevant  
114 features, thereby reducing the number of demonstrations needed and improving generalization. This  
115 multimodal feedback approach is especially promising in robotics, where safety and efficiency  
116 are paramount. Our work builds on these ideas by combining demonstration data with language  
117 instructions to guide a feature gating mechanism, leading to a reward model that is both data-efficient  
118 and robust.

## 119 2.3 Abstractions in Robot Learning

120 Another limitation of IRL comes from the spurious correlations of features. In many robotic tasks,  
121 not all sensory features are relevant for determining the reward. Ghosal *et al.*[3] aim to dynamically  
122 choose which features to rely on based on the current task or context. They have explored conditional  
123 gating mechanisms where a context variable modulates the importance of each input feature. Such  
124 approaches encourage sparsity in the feature set, thereby reducing the effective dimensionality of the  
125 learning problem. In robotics, this is particularly valuable since different tasks may require attention  
126 to different subsets of sensor modalities or object attributes. By integrating contextual reliability, one  
127 can obtain a more robust and interpretable model that adapts to the nuances of each task. Feature  
128 relevance varies with context, making contextual feature selection an essential component of robust  
129 learning. Unlike static feature selection, contextual feature selection dynamically selects which  
130 features are relevant based on auxiliary information such as the task or environment. In robotics,  
131 contextual feature selection is crucial for multi-task learning. Some skill learning frameworks enable  
132 robots to dynamically select relevant object attributes for different tasks, reducing learning complexity  
133 and improving generalization. Peng *et al.* [8] deals with overparameterization of reward by iteratively  
134 generating features and learning a reward on top of the current feature set. Peng *et al.* [8] uses  
135 language-guided contrastive explanations to iteratively extract and validate semantically meaningful  
136 features for the reward function. [25] uses background knowledge of language models to build state  
137 representations for unseen tasks. We aim to learn which state features matter under different user  
138 preferences, thereby improving sample efficiency and interpretability.

## 139 3 Problem Formulation

140 We consider the problem of learning reward functions that capture the unknown preferences held by  
141 a human given a small number of user demonstrations and language.

### 142 3.1 Preliminaries

143 We model our problem as a Markov Decision Process  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R} \rangle$  with states  $s \in \mathcal{S}$ , actions  
144  $a \in \mathcal{A}$ , transition probability  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ , and rewards  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . A solution  
145 to the MDP is a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  that specifies what actions the robot should take in different  
146 states. The reward function is typically parameterized (e.g. a neural network)  $\mathcal{R}_\theta(s)$ , and is intended  
147 to capture the human’s preference for how the robot should perform the task. To optimize task  
148 performance, the robot seeks a trajectory  $\tau = \{s^0, \dots, s^T\}$  that maximizes the cumulative reward  
149  $\mathcal{R}_\theta(\tau) = \sum_{s^t \in \tau} \mathcal{R}_\theta(s^t)$  and executes the corresponding actions.

### 150 3.2 Maximum Entropy Inverse RL (MaxEnt IRL)

151 In practice, the reward function  $\mathcal{R}_\theta$  is typically unknown to the robot or very challenging to manually  
152 specify. Thus, in IRL the robot’s goal is to *learn* this reward function from human feedback, such as  
153 demonstrations. Given a dataset of human-demonstrated trajectories  $\mathcal{D} = \{\tau_i\}_{i=1}^N$ , the robot treats  
154 them as evidence of the human’s preferred behavior and attempts to infer the reward parameters  $\theta$  that  
155 explain the underlying objective. We adopt the maximum entropy (MaxEnt) framework for modeling  
156 human decision-making [7, 5], where the human is assumed to be a noisily optimal agent who selects  
157 trajectories with probability proportional to their exponentiated reward:

$$p(\tau | \theta) = \frac{e^{\mathcal{R}_\theta(\tau)}}{\int_{\bar{\tau}} e^{\mathcal{R}_\theta(\bar{\tau})} d\bar{\tau}} \propto \exp(\mathcal{R}_\theta(\tau)) \quad (1)$$

158 This model captures the intuition that while humans generally act optimally, suboptimal trajectories  
159 are still possible, but occur with exponentially lower probability as their reward decreases [7]. To  
160 recover the reward parameters, we maximize the log-likelihood of the demonstrations:

$$\theta^* = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \sum_{\tau \in \mathcal{D}} \log p(\tau | \theta) . \quad (2)$$

161 Since the partition function in the denominator is intractable to compute exactly, we follow prior  
 162 work [5, 26] and use importance sampling to approximate it. Once the reward is learned, the  
 163 robot can act according to the policy that optimizes it. While MaxEnt IRL provides a principled  
 164 framework for inferring rewards from demonstrations, learning a flexible reward function directly  
 165 from high-dimensional states typically demands thousands of demonstrations per task [27, 28, 29],  
 166 which is costly and impractical to scale. With limited data, learned rewards often capture spurious  
 167 correlations between state features that accidentally co-occur with task success rather than reflecting  
 168 true human intent. This fundamentally limits generalization, particularly in environments with  
 169 distractors, ambiguous cues, or structural variations.

170 To address this, we propose leveraging natural language as an additional, structured form of supervi-  
 171 sion. Our key insight is that language plays a dual role in reward learning: 1) it *conveys information*  
 172 *about the human’s intent*, enabling a shared reward model to generalize across tasks via language  
 173 conditioning; and 2) it *implicitly communicates which aspects of the state are task-relevant*, providing  
 174 a signal for filtering out irrelevant environmental variation. By exploiting this natural duality, we learn  
 175 a language-conditioned reward function that both shares structure across tasks and ignores spurious  
 176 correlations, resulting in more generalizable rewards from significantly fewer demonstrations.

## 177 4 Method

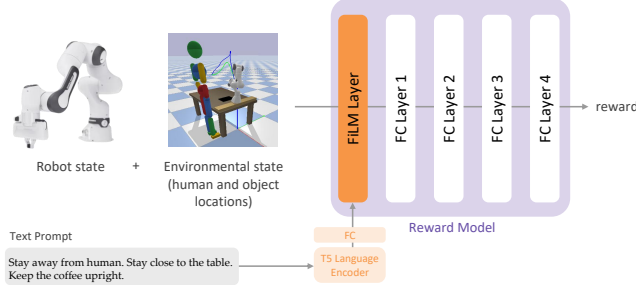
178 We present Masked Inverse Reinforcement Learning for Language Conditioned Reward Learning  
 179 (Masked IRL), a method which leverages demonstrations paired with human language instructions to  
 180 learn a language-conditioned reward function. Our approach exploits language’s two distinct signals:  
 181 language commands condition the preference captured by the reward model, and a language-informed  
 182 masking loss is used to enforce invariance to task-irrelevant state aspects. We generate this mask  
 183 directly from language commands and implement a masking loss that forces the reward function  
 184 to ignore spurious state elements. By combining this masking loss with a language-conditioned  
 185 architecture, Masked IRL achieves improved sample efficiency, requiring fewer demonstrations to  
 186 learn generalizable rewards.

### 187 4.1 Preliminaries

188 We assume the human maintains a set of ground truth state features  $\phi(s)$  which are only known  
 189 to the human, not the observing agent. We assume the ground truth reward for preference  $i$  is a  
 190 function of these features,  $\mathcal{R}_i^*(\phi(s))$ , where  $\mathcal{R}_i^*(\tau) = \sum_{s^t \in \tau} \mathcal{R}_i^*(\phi(s^t))$ . Given a set of training  
 191 preferences  $\mathcal{P}_{train} = \{1, 2, \dots, N\}$ , we collect a training dataset  $\mathcal{D} = \{\tau_i, \ell_i\}_{i=1}^N$ , where each paired  
 192 demonstration  $\tau_i$  and language command  $\ell_i$  correspond to preference  $i \in \mathcal{P}_i$ . We aim to learn a  
 193 general reward function  $\mathcal{R}_\theta(s|\ell_j)$  that captures the ground truth reward for a new preference  $j$  where  
 194  $j \notin \mathcal{P}_{train}$ . Our goal is to learn a reward function that can generalize to unseen preferences given just  
 195 a single language command  $\ell_j$ . Since we lack access to the ground truth state features, our inferred  
 196 reward is state-based  $\mathcal{R}_\theta(\tau|\ell_j) = \sum_{s^t \in \tau} \mathcal{R}_\theta(s^t|\ell_j)$ . We assume that all ground truth training and  
 197 test preferences are functions of the same set of ground truth human features, representing a consistent  
 198 intermediate representation unknown to the agent. We use language commands in our training dataset  
 199 in two distinct ways. First we condition our model on these language inputs, following established  
 200 practices in prior methods. Novel to our approach is our second usage – we convert language  
 201 commands into state-based masks that inform a specialized training loss, promoting invariance to  
 202 irrelevant state elements.

### 203 4.2 Language for State Masking

204 We extract state relevance from language by translating language commands into binary feature  
 205 masks. For each demonstration-language pair  $\{\tau, \ell\} \in \mathcal{D}$  we use language command  $\ell$  to generate a  
 206 binary mask  $m \in \{0, 1\}^d$ , where  $d$  is the dimension of the input state  $s$ . Each mask element is 1 for  
 207 state indices relevant to the specified preference, and 0 otherwise. We augment our dataset with these



**Figure 2: Network Architecture.** We condition the reward model on language instructions using FiLM layers. The conditioned reward model infers a scalar reward of a robot’s 9-dim state.

208 language-generated masks to create  $\mathcal{D}' = \{\tau_i, \ell_i, m_i\}_{i=1}^N$ . These masks are produced by leveraging  
 209 large language models ..

210 To ensure that the reward model is invariant to features deemed irrelevant by the language command,  
 211 we introduce a masking loss. Let  $s^{(j)}$  denote a perturbed version of state  $s \in \tau$ , where element  $j$  is  
 212 modified (such as through the addition of Gaussian noise) and all other elements remain unchanged.  
 213 The masking loss becomes

$$\mathcal{L}_{\text{mask}}(\theta) = \sum_{\tau, \ell, m \in \mathcal{D}'} \sum_{s \in \tau} \sum_{j=1}^d (1 - m_j) \left| R_{\theta}(s^{(j)} | \ell) - R_{\theta}(s | \ell) \right|, \quad (3)$$

214 where  $m_j$  represents the  $j^{\text{th}}$  element of  $m$ . This loss term penalizes changes in the reward when  
 215 irrelevant features are perturbed, forcing the reward model to ignore these features.

216 The final training loss becomes

$$\mathcal{J}(\theta) = \mathcal{L}_{\text{IRL}}(\theta) + \lambda \mathcal{L}_{\text{mask}}(\theta), \quad (4)$$

217 where  $\lambda > 0$  is a hyperparameter controlling the trade-off between fitting the demonstrations and  
 218 enforcing invariance to irrelevant state elements.

### 219 4.3 Masked IRL for Language Conditioned Reward Learning

220 We pair our masking loss with a language-conditioned architecture to additionally leverage the  
 221 intent captured by language instructions. Specifically, we apply Feature-wise Linear Modulation  
 222 (FiLM) [30] to the first fully connected (FC) layer of the reward model (see Fig. 2) to condition  
 223 the reward model based on the language inputs. This FiLM layer applies language-dependent affine  
 224 transformations to intermediate network features, allowing language commands to dynamically  
 225 modulate reward components directly. As opposed to simply concatenating the language command  
 226 with the input state, FiLM targets conditioning input to explicitly modulate intermediate network  
 227 features, providing the ability to scale features, negate them, or shut them off entirely. This method  
 228 enables using language for a dual purpose: both as a gating mechanism that filters out irrelevant  
 229 state aspects, and as an adaptation function that adjusts intermediate feature weights based on the  
 230 preference captured in language. Algorithm 1 shows the training procedure for Masked IRL.

### 231 4.4 Clarifying Ambiguous Language Instructions

232 For ambiguous language instructions, we systematically generate the instructions within two types:  
 233 (1) referent omitted and (2) expression omitted. Referent omitted instructions do not include the  
 234 object that the user actually cares about, and only include instructions such as “stay away”, “stay  
 235 close”, and “carry it upright”. Expression omitted instructions have the information about what object  
 236 the user wants to refer to, but does not mention how the user wants the relationship between the  
 237 robot and the object to be. For instance, “table”, “laptop”, or “human” can be expression omitted  
 238 instructions. To generate state masks from ambiguous instructions, we provide the information of  
 239 a demonstration trajectory as tabular data in text, along with the instruction to an LLM. We use



240 Chain-of-thought reasoning to let LLM generate the response step-by-step, including its reasoning  
 241 process to generate the clarified language instruction. For instance, given the instruction ‘stay away’  
 242 and a demonstration where the robot moves away from the table, the LLM might reason: ‘The robot  
 243 avoids the table. Therefore, the instruction likely refers to avoiding the table.’ The clarified instruction  
 244 becomes ‘stay away from the table’, which is then mapped to a binary mask emphasizing the end  
 245 effector’s z position and de-emphasizing human or laptop locations. Then, we query the LLM again  
 246 to convert the clarified language instruction into a 9-dim state mask that represents the importance of  
 247 each state dimension.

---

**Algorithm 1:** Masked IRL with Language Conditioning

---

**Input:** Demonstrations  $\{(\tau_i, \theta_i)\}_{i=1}^N$ , training trajectories  $\mathcal{T}$ , language encoder  $E$ , reward network  $R$ , learning rate  $\eta$ , iterations  $I$ , batch size  $B$ , masked loss weight  $\lambda$ , noise scale  $\sigma$ .

**for** *epoch* 1 to  $I$  **do**

    Shuffle demo and training indices

**for** *each minibatch*  $b$  of size  $B$  **do**

        Form demo inputs:  $X_d^b = \{(\bar{s}_i, c_i)\}$  and compute cost  $C_d^b = R(X_d^b)$

        Form training inputs  $X_t^b = \{(\bar{s}_j, c_j)\}$  with cost  $C_t^b = R(X_t^b)$

        Compute maxent loss:

$$248 \quad \mathcal{L}_{\text{IRL}}^b = \text{mean}(C_d^b) + \log\left(\text{mean}(\exp(-C_t^b))\right)$$

        Perturb demo states:  $\bar{s}'_i = \bar{s}_i + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$  (only in dimensions where  $\Pi(\theta_i) = 0$ );

        Compute perturbed cost  $C_d'^b = R(\{\bar{s}'_i, c_i\})$  and masked loss:

$$\mathcal{L}_{\text{mask}}^b = \text{mean}\left(|C_d^b - C_d'^b|\right)$$

        Update parameters:  $\theta \leftarrow \theta - \eta \nabla\left(\mathcal{L}_{\text{IRL}}^b + \lambda \mathcal{L}_{\text{mask}}^b\right)$

**end**

**end**

**return**  $\theta$ .

---

## 249 5 Experiments

250 We evaluate our method on a robotic task to move a coffee from a start to a goal location in a PyBullet  
 251 simulator, where there is a human, a table, and a laptop in the environment. Each state consists of  
 252 the position and rotation of the robot’s end effector, objects (table and laptop), and a human in the  
 253 environment. In each task, only a subset of features is relevant to the reward. Human instructions  
 254 (e.g., “stay away from the laptop”) are provided to guide the feature gating.

### 255 5.1 Dataset.

256 We generate a dataset of 20 object configurations and 10 start-goal pairs per configuration for a task  
 257 of moving a coffee mug, each with 5 robot trajectories, in PyBullet simulator. We also generate  
 258 242 language instructions that are mapped into ground truth reward functions that define human  
 259 preferences. For clear language instructions, we construct the dataset with 50 train instructions and  
 260 30 test instructions. Each instruction has a corresponding 5-dim theta value that describes human’s  
 261 ground truth reward function. We use GPT-4o API to infer the state mask only from each clear  
 262 instruction, without any information about the ground truth reward. For inferring state masks from  
 263 ambiguous instructions, we pair each instruction with its corresponding expert demonstration and pass  
 264 the information of the language instruction and demonstration to GPT-4o as described in [Section 4.4](#).  
 265 We train each model with 10 demonstrations per human preference.

266 **5.2 Baselines**

267 Traditional IRL learns a single reward function from demonstrations, without contextual modulation.  
268 This often results in a reward model that uses all features indiscriminately, making it vulnerable to  
269 spurious correlations when demonstrations cover multiple tasks or environments. To demonstrate  
270 the effectiveness of masking loss, we compare Masked IRL and MaxEnt IRL on two different types  
271 of reward model - single model and multiple model. We refer to ‘single model’ as a language-  
272 conditioned reward model, regardless of the usage of masking loss. ‘Multiple model’ refers to a  
273 set of language-unconditioned reward models, where each element of the set is a reward model that  
274 corresponds to a specific language instruction, i.e., user preference. For multiple model methods,  
275 we only evaluate on seen human preferences, since unseen human preferences do not have any  
276 corresponding trained reward model. However, for single model methods, we evaluate on both seen  
277 and unseen human preferences, i.e., language instructions. Since we focus on the single model  
278 experiments in this section, experiment details and results for multiple models are in the appendix.

279 For single model approaches, we compare:

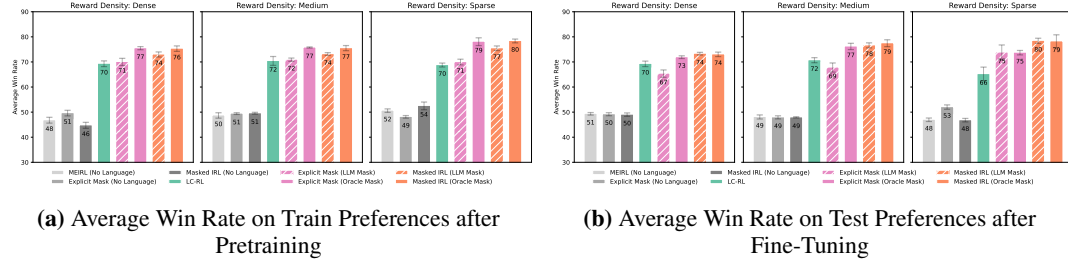
- 280 • **LC-RL [1] (Language-conditioned)**. A 4-layer MLP reward model is conditioned on  
281 language embedding using FiLM. We use the standard maximum entropy loss function to  
282 train this model.
- 283 • **Masked IRL (Language-conditioned, Oracle Mask)**. Same architecture as the LC-RL  
284 baseline but uses the weighted masking loss in addition to the maximum entropy loss for  
285 training. Oracle state mask is used.
- 286 • **Masked IRL (Language-conditioned, LLM Mask)**. Same architecture as the LC-RL  
287 baseline but uses the weighted masking loss in addition to the maximum entropy loss for  
288 training. LLM generated state mask is used.
- 289 • **Explicit Mask (Language-conditioned, Oracle Mask)**. Same architecture and training  
290 loss as the LC-RL baseline but uses oracle state mask to mask out irrelevant state dimensions  
291 given a language instruction.
- 292 • **Explicit Mask (Language-conditioned, LLM Mask)**. Same architecture and training loss  
293 as the LC-RL baseline but uses LLM generated state mask to mask out irrelevant state  
294 dimensions given a language instruction.
- 295 • **MaxEnt IRL (No Language)**. We use the standard maximum entropy loss function to train  
296 a 4-layer MLP reward model. This model is not conditioned on language.
- 297 • **MaxEnt IRL (No Language)**. Same architecture as the MaxEnt IRL baseline but uses the  
298 weighted masking loss in addition to the maximum entropy loss for training. Oracle state  
299 mask is used.
- 300 • **Explicit Mask (No Language)**. Same architecture and training loss as the MaxEnt IRL  
301 baseline but uses oracle state mask to mask out irrelevant state dimensions given a language  
302 instruction.

303 For simplicity, we omit ‘language-conditioned’ when we refer to language-conditioned single model  
304 baselines.

305 **5.3 Evaluation Metrics.**

306 We evaluate all models by calculating the average win rate, where the average win rate measures  
307 how often our learned reward model correctly prefers better trajectories compared to ground-truth  
308 preferences. We measure the average win rate on three different reward densities: sparse, medium,  
309 and dense, where sparser reward density implies less features are important in the environment. The  
310 sparsity of the ground truth reward model is chosen based on the number of valid features from 1 to 5  
311 (sparse: 1, 2, medium: 3, dense: 4, 5). We also We run all experiments with 5 different random seeds  
312 (12345, 23451, 34512, 45123, and 51234) and show the average and standard error across seeds.





**Figure 3: Experiment Results.** (a) and (b) show the average win rate of single model methods on different reward densities after pretraining on 40 train preferences for 1k epochs and fine-tuning on 30 test preferences for 100 epochs, respectively. All models are trained with 10 demonstrations per user preference and evaluated with unseen trajectories with novel object configurations. Error bars show the standard error across 5 different seeds.

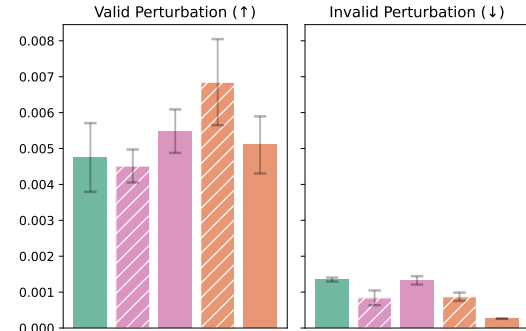
### 5.4 Results

**The effectiveness of Masking Loss on Single Model.** Fig. 3 shows the average win rate across reward densities (dense, medium, sparse) for both (a) train and (b) test preferences. Across all reward densities and both train and test preferences, Masked IRL consistently outperforms the LC-RL baseline that lacks masking loss. Explicit Mask that uses oracle state mask also outperforms LC-RL but shows significant performance decline when LLM generated masks are used. In contrast, Masked IRL outperforms LC-RL with both oracle and LLM generated masks, demonstrating its robustness to the quality of the state masks. This implies that the masking loss effectively reduces spurious correlations by enforcing invariance to irrelevant features, thereby enhancing the stability and efficacy of reward learning. All language unconditioned baselines show poor performance compared to language conditioned models, which shows the effectiveness of training a language conditioned reward model for multiple preferences. Fig. 4 shows the reward variance when valid (state mask element is 1) and invalid (state mask element is 0) state dimensions are perturbed in test trajectories. With oracle masks, Masked IRL shows the strongest invariance when invalid state dimensions are perturbed.

#### Performance on ambiguous instructions.

When we use our Masked IRL single model trained with 10 demonstrations per human preference to evaluate trajectories given a single ambiguous language instruction and an expert demonstration to disambiguate language, we get an average win rate of 63.1% on the instructions. The lower performance compared to the performance on clear test instructions may be due to the inaccuracy of clarifying ambiguous instructions to clear instructions using LLMs.

**Future Work and Limitations** Although our Masked IRL framework effectively improves generalization and sample efficiency, several limitations remain. First, our reliance on LLMs introduces potential inaccuracies in generating relevance masks, particularly when instructions are ambiguous or nuanced, which can affect the overall robustness of the reward model. Future work could explore methods for refining mask accuracy through interactive human feedback or advanced prompting strategies. Additionally, our current evaluations focus on relatively constrained robotic tasks; extending the approach to more complex, dynamic, or multi-agent environments could further validate the generality of Masked IRL. Lastly, investigating ways to integrate explicit uncertainty estimation in the masking process could enhance the reliability of our approach in real-world deployments.

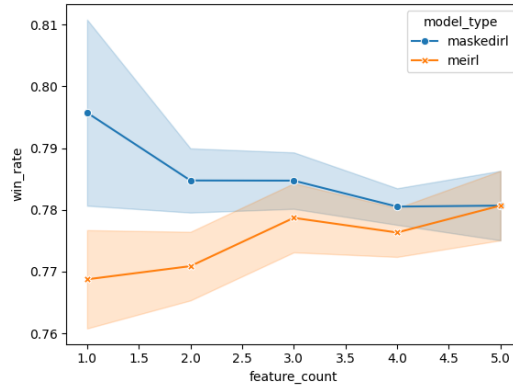


**Figure 4: Reward Variance when Valid or Invalid State Dimensions are Perturbed.** The plots show reward variance when state dimensions that are valid or invalid given language instructions are perturbed. Higher variance when valid dimensions are perturbed and lower variance when invalid dimensions are perturbed imply the reward model is more sensitive to valid changes and less sensitive to invalid changes in the environment, respectively.

## References

- [1] J. Fu et al. From language to goals: Inverse reinforcement learning for vision-based instruction following. In *International Conference on Learning Representations (ICLR)*, 2019.
- [2] S. Poddar, Y. Wan, H. Ivison, A. Gupta, and N. Jaques. Personalizing reinforcement learning from human feedback with variational preference learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [3] G. Ghosal, A. Setlur, D. S. Brown, A. D. Dragan, and A. Raghunathan. Contextual reliability: When different features matter in different contexts. In *Proceedings of ICML*, 2023.
- [4] A. Y. Ng and S. Russell. Algorithms for inverse reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, 2000.
- [5] C. Finn, S. Levine, and P. Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. In *International conference on machine learning*, pages 49–58. PMLR, 2016.
- [6] P. Abbeel and A. Y. Ng. Apprenticeship Learning via Inverse Reinforcement Learning. In *Proceedings of the International Conference on Machine learning*, 2004.
- [7] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey, et al. Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- [8] A. Peng, B. Z. Li, I. Sucholutsky, N. Kumar, J. A. Shah, J. Andreas, and A. Bobu. Adaptive language-guided abstraction from contrastive explanations. *Conference on Robot Learning (CoRL)*, 2024.
- [9] T. Summers, R. Hawkins, M. K. Ho, T. Griffiths, and D. Hadfield-Menell. How to talk so ai will learn: Instructions, descriptions, and autonomy. *Advances in Neural Information Processing Systems*, 35:34762–34775, 2022.
- [10] A. Bobu, Y. Liu, R. Shah, D. S. Brown, and A. D. Dragan. Sirl: Similarity-based implicit representation learning. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot Interaction*, pages 565–574, 2023.
- [11] P. F. Christiano, J. Leike, T. B. Brown, M. Martic, S. Legg, and D. Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4299–4307, 2017.
- [12] S. Casper, X. Davies, C. Shi, T. K. Gilbert, J. Scheurer, J. Rando, R. Freedman, T. Korbak, D. Lindner, P. Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- [13] M. Hwang, G. Lee, H. Kee, C. W. Kim, K. Lee, and S. Oh. Sequential preference ranking for efficient reinforcement learning from human feedback. *Advances in Neural Information Processing Systems*, 36:49088–49099, 2023.
- [14] Y. J. Ma, W. Liang, G. Wang, D.-A. Huang, O. Bastani, D. Jayaraman, Y. Zhu, L. Fan, and A. Anandkumar. Eureka: Human-level reward design via coding large language models. *arXiv preprint arXiv:2310.12931*, 2023.
- [15] W. Yu, N. Gileadi, C. Fu, and et al. Language to rewards for robotic skill synthesis. In *Conference on Robot Learning (CoRL)*, 2023.
- [16] M. Hwang, L. Weihs, C. Park, K. Lee, A. Kembhavi, and K. Ehsani. Promptable behaviors: Personalizing multi-objective rewards from human preferences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16216–16226, 2024.
- [17] Z. Yang, M. Jun, J. Tien, S. J. Russell, A. D. Dragan, and E. Bıyık. Trajectory improvement and reward learning from comparative language feedback. 2024.

- 397 [18] Y. Cui, S. Karamcheti, et al. “no, to the right” – online language corrections for robotic  
398 manipulation via shared autonomy. In *Proceedings of HRI 2023*, 2023.
- 399 [19] M. Ahn, A. Brohan, N. Brown, and et al. Do as i can, not as i say: Grounding language in  
400 robotic affordances. In *Conference on Robot Learning (CoRL)*, 2022.
- 401 [20] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners:  
402 Extracting actionable knowledge for embodied agents. In *International Conference on Machine  
403 Learning (ICML)*, 2022.
- 404 [21] W. Huang, F. Xia, T. Xiao, and et al. Inner monologue: Embodied reasoning through planning  
405 with language models. In *Conference on Robot Learning (CoRL)*, 2022.
- 406 [22] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as  
407 policies: Language model programs for embodied control. In *IEEE International Conference  
408 on Robotics and Automation (ICRA)*, 2023.
- 409 [23] S. Karamcheti, S. Nair, A. S. Chen, T. Kollar, C. Finn, D. Sadigh, and P. Liang. Language-driven  
410 representation learning for robotics. In *Robotics: Science and Systems (RSS)*, 2023.
- 411 [24] M. Hwang, J. Hejna, D. Sadigh, and Y. Bisk. Motif: Motion instruction fine-tuning. *IEEE  
412 Robotics and Automation Letters*, 2025.
- 413 [25] A. Peng, I. Sucholutsky, B. Li, T. Summers, T. Griffiths, J. Andreas, and J. Shah. Learning with  
414 language-guided state abstractions. In *International Conference on Learning Representations*,  
415 2024.
- 416 [26] A. Bobu, M. Wiggert, C. Tomlin, and A. D. Dragan. Inducing structure in reward learning by  
417 learning features. *The International Journal of Robotics Research*, 41(5):497–518, 2022.
- 418 [27] T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel. Deep imitation  
419 learning for complex manipulation tasks from virtual reality teleoperation. In *2018 IEEE  
420 International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May  
421 21-25, 2018*, pages 1–8. IEEE, 2018. doi:10.1109/ICRA.2018.8461249. URL [https://doi.  
422 org/10.1109/ICRA.2018.8461249](https://doi.org/10.1109/ICRA.2018.8461249).
- 423 [28] R. Rahmatizadeh, P. Abolghasemi, L. Bölöni, and S. Levine. Vision-based multi-task ma-  
424 nipulation for inexpensive robots using end-to-end learning from demonstration. In *2018  
425 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia,  
426 May 21-25, 2018*, pages 3758–3765. IEEE, 2018. doi:10.1109/ICRA.2018.8461076. URL  
427 <https://doi.org/10.1109/ICRA.2018.8461076>.
- 428 [29] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine.  
429 Learning complex dexterous manipulation with deep reinforcement learning and demonstrations.  
430 In H. Kress-Gazit, S. S. Srinivasa, T. Howard, and N. Atanasov, editors, *Robotics: Science and  
431 Systems XIV, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA, June 26-30, 2018*,  
432 2018. doi:10.15607/RSS.2018.XIV.049. URL [http://www.roboticsproceedings.org/  
433 rss14/p49.html](http://www.roboticsproceedings.org/rss14/p49.html).
- 434 [30] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with a  
435 general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*,  
436 volume 32, 2018.



**Figure 5: Multiple model performance per feature counts.** Comparing Masked IRL to MaxEnt IRL for multiple model baselines, the average win rate is improved the most when there are least number of valid features in the ground truth reward of the simulated human. As the number of valid features increase from 1 to 5, the performance gap between Masked IRL and MaxEnt IRL decreases.

## 437 A Additional Experiments

### 438 A.1 Multiple Model Experiments

439 For multiple model approaches, we compare:

- 440 • **MaxEnt IRL (No Language, multiple model).** We train a 3-layer MLP that inputs a  
441 9-dimensional state and outputs a scalar reward value for each state. We train this baseline  
442 with standard maximum entropy loss.
- 443 • **Masked IRL (No Language, multiple model).** Same architecture as the baseline but uses  
444 the weighted masking loss in addition to the maximum entropy loss for training.

445 **The effectiveness of Masking Loss on Multiple Model.** Fig. 5 shows the effect of having masking  
446 loss in multiple model methods. Interestingly, the performance improvement by using masking loss  
447 is maximized when the number of valid features for the ground truth reward of the simulated human  
448 is minimized to 1. As the number of valid features increases, the gap between Masked IRL and  
449 MaxEnt IRL decreases. This is a desired behavior because when all features are valid, i.e., all state  
450 dimensions are relevant to the instruction,