

# What to Represent, How to Act: Programmatic Feature Induction for Few-Shot Bayesian Imitation

**Anonymous authors**

Paper under double-blind review

## Abstract

1 A natural way to teach an agent a new task is to give it a few expert demonstrations.  
2 But few-shot imitation requires more than reproducing observed state–action pairs: the  
3 agent must infer an underlying policy that generalizes across task instances. We study  
4 this setting as a controlled instance of a broader challenge in big worlds, where gener-  
5 alizing from limited demonstrations requires discovering task-relevant abstractions that  
6 are not given in advance. We propose TwoLIP: Two-level LLM-Guided Induction of  
7 Programmatic Policies. At the first level, TwoLIP uses a large language model (LLM)  
8 to build a library of programmatic state–action features. At the second level, it performs  
9 Bayesian implicit behavioral cloning to learn a programmatic policy that combines the  
10 features. This two-level pipeline separates the problem of what to represent from how to  
11 act, combining a semantic prior over candidate features with a structural prior over pol-  
12 icy composition. We evaluate TwoLIP on Generalization Grid Games, a suite of strategy  
13 tasks where policies must generalize to unseen configurations that differ substantially  
14 from the demonstrations. TwoLIP consistently outperforms behavioral cloning (BC)  
15 baselines, including deep neural models and end-to-end LLM-based policy-generation  
16 approaches. We also find preliminary evidence that induced feature libraries transfer  
17 across related tasks. Together, these results point toward a path for few-shot imitation  
18 in big worlds, where the right abstractions cannot be prespecified.

## 19 1 Introduction

20 In large, open-ended worlds, agents will encounter tasks that were not anticipated when their policy  
21 representations were designed. Consider a household robot learning from a few demonstrations  
22 to organize an unfamiliar kitchen counter: the user may move a mug beside the coffee machine  
23 or place a spoon in a drawer. What transfers is not a trajectory or visual template, but semantic  
24 distinctions such as which object is out of place and where it belongs. This example reflects a  
25 central challenge in modern imitation learning (IL): demonstrations provide a natural interface for  
26 teaching agents, but a few demonstrations rarely specify the policy uniquely. In big worlds, this  
27 ambiguity is amplified, because agents must generalize across changing objects, layouts, and task  
28 instances where the relevant abstractions are not known in advance.

29 Standard BC treats imitation as supervised action prediction: given an observation, predict the  
30 demonstrated action (Pomerleau, 1988). Recent advances in robot learning have scaled this  
31 paradigm to increasingly complex behaviors from demonstration data (Zhao et al., 2023; Chi et al.,  
32 2023; ?, Khazatsky et al., 2024). A related line of work on implicit BC instead scores candidate  
33 state–action pairs and has shown strong results in robot policy learning (Florence et al., 2022).  
34 However, both explicit and implicit BC remain representation-dependent: with very limited data,  
35 they can rely on brittle patterns unless the representation captures the semantic distinctions needed  
36 to choose among actions (Pari et al., 2022).

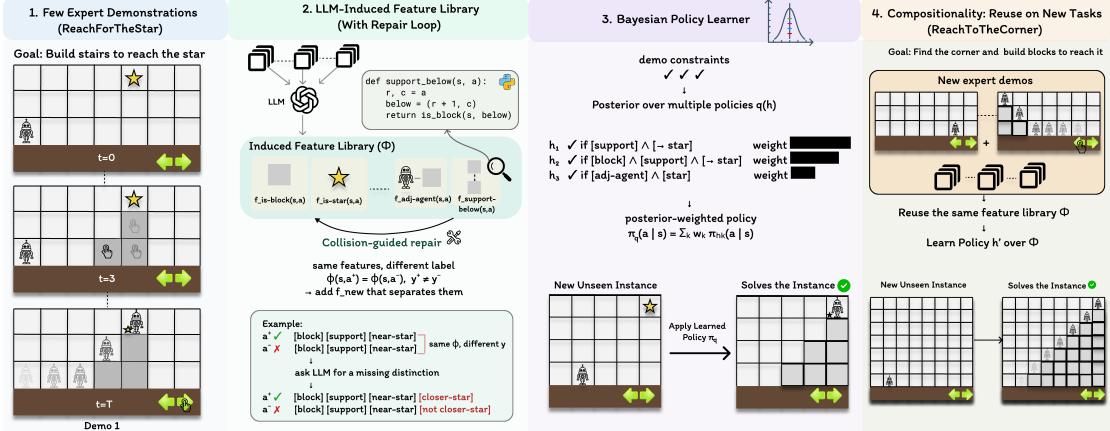


Figure 1: Overview of TwoLIP. From a few expert demonstrations, an LLM induces a reusable feature library, a Bayesian learner composes these features into a logical policy, and the same library can be reused for related tasks.

37 Structured symbolic and programmatic approaches address this representation problem through a  
 38 strong inductive bias, expressing policies in terms of rules or programs (Verma et al., 2018; Silver  
 39 et al., 2020; Xin et al., 2023; Patton et al., 2023; Cui et al., 2025), but typically assume that the  
 40 relevant predicates or domain-specific languages (DSLs) are provided in advance. LLMs offer a  
 41 promising way to relax this assumption: code-as-policies (CAP) and similar lines of work show  
 42 that LLMs can generate executable policy code from language and examples (Liang et al., 2023).  
 43 However, as we show in our experiments, directly asking an LLM to infer a complete policy from a  
 44 few demonstrations can be brittle in complex environments.

45 To address these challenges, we propose TwoLIP—Two-level LLM-Guided Induction of Programmatic  
 46 Policies—which separates the problem of *what to represent* from *how to act*. Rather than  
 47 asking an LLM to generate a complete policy, TwoLIP uses the LLM to induce and refine pro-  
 48 grammatic state–action features (*what to represent*). TwoLIP then performs implicit BC over this  
 49 induced feature space (*how to act*). A further benefit of this two-level pipeline is that it admits a nat-  
 50 ural feedback loop between feature generation and classification. We begin with an initial feature set  
 51 and identify positive and negative state–action pairs that the current features cannot separate. These  
 52 pairs are fed back to the LLM to guide targeted feature generation, and the process repeats until all  
 53 pairs are separable (or a maximum number of rounds is reached). To maintain interpretability and  
 54 quantify uncertainty, we use an implicit BC technique that is programmatic and Bayesian (Silver  
 55 et al., 2020). See Figure 1 for an overview of the full method.

56 We evaluate TwoLIP on Generalization Grid Games (Silver et al., 2020), a suite of strategy tasks  
 57 that test generalization from a few demonstrations to substantially different held-out game instances.  
 58 TwoLIP achieves higher task-completion success than supervised BC and end-to-end LLM policy-  
 59 generation baselines in the low-data regime, while approaching the performance of an oracle that  
 60 uses a hand-designed feature grammar. We also find preliminary evidence that induced feature li-  
 61 braries can be reused across related tasks without retraining, suggesting that the learned represen-  
 62 tations capture reusable task structure rather than demonstration-specific correlations. Together, these  
 63 results suggest that decoupling representation discovery from policy learning is a promising path  
 64 toward few-shot imitation in big worlds, where the right abstractions cannot be specified in advance.

## 65 2 Problem Setting

66 We consider IL in a finite-horizon Markov decision process  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, \rho_0, \mathcal{G}, H)$ , where  $\mathcal{S}$  is  
 67 the state space,  $\mathcal{A}$  is the action space,  $P(s' | s, a)$  is the transition kernel,  $\rho_0$  is a distribution over  
 68 initial states,  $\mathcal{G} \subseteq \mathcal{S}$  is the set of goal states, and  $H$  is the horizon. The learner is given a small set of

69 expert demonstrations  $\mathcal{D} = \{\tau_i\}_{i=1}^N$ , where  $\tau_i = (s_0^i, a_0^i, s_1^i, a_1^i, \dots, s_{T_i}^i)$ ,  $s_t^i \in \mathcal{S}$ ,  $a_t^i \in \mathcal{A}$ , and  $N$  is  
 70 small (e.g.,  $N \leq 10$  in our experiments). The demonstrations are generated by an expert policy

$$\pi^* : \mathcal{S} \rightarrow \mathcal{A}, \quad a_t^i = \pi^*(s_t^i),$$

71 which induces successful behavior in  $\mathcal{M}$ . We assume demonstrations come from a consistent expert  
 72 with a canonical action for each demonstrated state, though not necessarily a uniquely optimal one:  
 73 if multiple actions are equally optimal, each is valid. The goal is to learn a policy  $\pi \in \Pi$ , with  
 74  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ , from  $\mathcal{D}$  that imitates the expert and generalizes beyond the observed trajectories.  
 75 We evaluate a learned policy by rolling it out from held-out initial states  $s_0 \sim \rho_{\text{test}}$  and measuring  
 76 task-completion success:

$$J(\pi) = \Pr_{\substack{s_0 \sim \rho_{\text{test}}, \\ a_t \sim \pi(\cdot | s_t), \\ s_{t+1} \sim P(\cdot | s_t, a_t)}} [\exists t \leq H \text{ such that } s_t \in \mathcal{G}].$$

## 77 3 Method

78 **Overview.** TwoLIP separates the few-shot imitation problem into two levels. First, it constructs  
 79 a library of semantic state–action features: reusable Boolean functions of a candidate action in the  
 80 current state. Second, it performs Bayesian inference over structured policies that combine these  
 81 features to decide which actions are admissible. The remainder of this section describes the two  
 82 levels of TwoLIP. We first define the behavior-learning problem by converting demonstrations into  
 83 a contrastive state–action dataset, which provides positive and negative examples for the learning  
 84 task (Section 3.1). We then introduce the induced policy class and Bayesian inference procedure  
 85 (Sections 3.2 and 3.3). Finally, we describe how TwoLIP constructs the feature library from demon-  
 86 strations and uses feature collisions to diagnose and repair missing abstractions (Sections 3.4, 3.5,  
 87 and 3.6). Algorithm 1 summarizes the pipeline, and Figure A.1 visualizes a learned policy.

### 88 3.1 Contrastive State–Action Imitation Data

89 We convert demonstrations into a labeled contrastive dataset over state–action pairs. Each demon-  
 90 strated action  $a_t^i$  at state  $s_t^i$  gives a positive example  $(s_t^i, a_t^i, 1)$ . To construct negatives, we select  
 91 alternative actions from a candidate set  $\mathcal{C}^-(s_t^i, a_t^i) \subseteq \mathcal{A} \setminus \{a_t^i\}$ . This set may contain all non-expert  
 92 actions or only a sampled subset. The resulting dataset is  $\mathcal{D}_{\text{sa}} = \{(s, a, y)\}$ , with  $y = 1$  for expert  
 93 actions and  $y = 0$  for alternatives. This turns imitation into a state–action discrimination problem:  
 94 the learner compares the expert action against other actions available in the same state, assigning  
 95 high admissibility, or low energy, to expert actions and lower admissibility to alternatives (Florence  
 96 et al., 2022; LeCun et al., 2006).

### 97 3.2 Feature-Induced Implicit Policies

98 A feature is a Boolean function over state–action pairs. Given a feature library  $\Phi = \{\phi_1, \dots, \phi_m\}$ ,  
 99 each state–action pair is mapped to a Boolean vector  $x_\Phi(s, a) = (\phi_1(s, a), \dots, \phi_m(s, a)) \in$   
 100  $\{0, 1\}^m$ . TwoLIP learns a structured classifier  $h : \{0, 1\}^m \rightarrow \{0, 1\}$ , where  $h(x_\Phi(s, a)) = 1$   
 101 indicates that  $a$  is consistent with the expert at state  $s$  and  $h(x_\Phi(s, a)) = 0$  indicates that it is  
 102 not. This defines an implicit policy over candidate actions: rather than predicting an action directly  
 103 from the state,  $h$  evaluates each candidate state–action pair through the induced feature representa-  
 104 tion. The test-time decision rule, which uses the posterior over classifiers, is given in Section 3.3.  
 105 To keep the learned policy human-interpretable, we restrict  $h$  to logical rules over the induced  
 106 features. Concretely, policies are represented as disjunctions of conjunctions of feature literals,  
 107  $h(x_\Phi(s, a)) = \bigvee_r \bigwedge_{j \in I_r} \ell_{rj}(s, a)$ , where  $I_r \subseteq \{1, \dots, m\}$  and  $\ell_{rj}(s, a) \in \{\phi_j(s, a), \neg\phi_j(s, a)\}$ .  
 108 Each conjunction gives one sufficient condition for an action to be admissible, while the disjunction  
 109 combines multiple such reasons. In practice, we instantiate this policy class using decision trees over  
 110 Boolean features (Breiman et al., 1984): each path to a positive leaf corresponds to one conjunction,  
 111 and the set of positive leaves defines the disjunction.

### 112 3.3 Bayesian Inference over Structured Policies

113 The policy-composition layer builds on the Bayesian logical policy-learning formulation of Log-  
 114 ical Programmatic Policies (LPP) (Silver et al., 2020). This Bayesian view is useful in our few-  
 115 demonstration setting because the data often underdetermines the policy: many classifiers may ex-  
 116 plain the observed state–action labels, and committing to a single one can overfit to accidental pat-  
 117 terns in the demonstrations. We therefore maintain an approximate posterior over structured policies.  
 118 Let  $\mathcal{H}_\Phi$  denote the hypothesis class of structured classifiers over the induced feature representation  
 119  $x_\Phi$ . Each candidate policy  $h \in \mathcal{H}_\Phi$  receives a posterior score that combines its agreement with the  
 120 contrastive dataset  $\mathcal{D}_{\text{sa}}$  and a structural prior favoring simpler policies. In particular, policies that  
 121 mark expert actions as admissible and alternatives as inadmissible receive higher data likelihood,  
 122 while more complex policies receive lower prior probability.

123 Exact inference over  $\mathcal{H}_\Phi$  is intractable, so we approximate the posterior with a weighted collec-  
 124 tion of policy particles  $q(h) = \sum_{k=1}^K w_k \delta_{h_k}(h)$ , where each  $h_k$  is a structured classifier over  
 125 the induced feature space and  $w_k$  is proportional to its posterior score. At test time, the posterior  
 126 mixture assigns each candidate action a soft admissibility score  $S_q(s, a) = \mathbb{E}_{h \sim q}[h(x_\Phi(s, a))] =$   
 127  $\sum_{k=1}^K w_k h_k(x_\Phi(s, a))$ . The policy then selects an action with maximum posterior expected admis-  
 128 sibility,  $\pi_q(s) \in \arg \max_{a \in \mathcal{A}} S_q(s, a)$ . This posterior mixture acts as a soft energy-like state–action  
 129 scoring model and is especially useful in the few-demonstration setting, where several distinct poli-  
 130 cies may be consistent with the observed behavior. Additional details are in Appendix A.

### 131 3.4 Demo-Guided State–Action Feature Induction

132 The policy learner above assumes a feature library  $\Phi$ , but in real-world problems the relevant ab-  
 133 stractions are often not known in advance. We therefore do not assume access to an oracle set of  
 134 human-designed features, DSL primitives, or grammar rules tailored to the task. Instead, TwoLIP  
 135 constructs this library automatically from demonstrations using an LLM. The prompt template used  
 136 for feature induction is provided in Appendix E.1.

137 Rather than asking the LLM to synthesize the full policy, we ask it to generate feature-template  
 138 generators: small Python functions that define families of state–action features. For example, in a  
 139 grid task, a template may check whether the candidate action selects an object of type  $c$ , or whether  
 140 it moves toward an object with property  $p$ . These templates are expanded into candidate Boolean  
 141 features; Before adding candidates to the library, we evaluate them on  $\mathcal{D}_{\text{sa}}$  and apply a *feature-*  
 142 *filtering* phase. This removes features that fail to execute, violate the required type signature, are  
 143 constant on the dataset, or duplicate the behavior of an existing feature on the observed examples.  
 144 The surviving features form the initial library  $\Phi^{(0)} = \{\phi_1, \dots, \phi_m\}$ , where the superscript 0 denotes  
 145 the library before any repair rounds.

### 146 3.5 Detecting Feature Insufficiency with Collisions

147 The induced feature library determines what distinctions the policy learner can express. To deter-  
 148 mine whether the current feature set can separate expert actions from alternatives, we introduce the  
 149 concept of **feature collisions**. A collision occurs when two labeled examples have the same feature  
 150 representation but different labels:  $x_\Phi(s, a) = x_\Phi(s', a')$  and  $y \neq y'$ .

151 When this happens, under the consistency assumption from Section 2, no deterministic classifier  
 152 whose prediction depends only on  $x_\Phi$  can correctly classify both examples. Thus, mixed-label  
 153 collisions identify an irreducible source of empirical error for any policy expressed using the current  
 154 feature library. To measure this effect, we group examples in  $\mathcal{D}_{\text{sa}}$  into *feature buckets*, where each  
 155 bucket contains all examples with the same feature vector  $x_\Phi(s, a)$ . A bucket is mixed-label if it  
 156 contains both expert actions and alternatives. We quantify collision severity by counting either the  
 157 number of mixed-label buckets or the number of positive–negative pairs within them.

158 This diagnostic is one of the key distinctions of our framework. Because feature induction and policy  
 159 learning are separated, the feature library is an explicit intermediate object that can be inspected,  
 160 filtered, and repaired before learning the final policy.

### 161 3.6 Collision-Guided Feature Repair

162 TwoLIP uses feature collisions as targeted feedback for improving the library. At repair round  $r$ ,  
 163 after constructing the current feature set  $\Phi^{(r)}$ , we identify mixed-label buckets in  $\mathcal{D}_{\text{sa}}$  and sample  
 164 representative<sup>1</sup> positive–negative collision pairs from them. These pairs are given to the LLM as  
 165 counterexamples showing where the current abstraction fails, and the LLM is prompted to propose  
 166 additional predicates that distinguish the expert action from the alternative. Appendix E.2 provides  
 167 the full prompt and implementation details.

168 The newly proposed features are instantiated and optionally filtered. We accept them only if they  
 169 reduce the collision-induced empirical-error bound; otherwise, we discard them and keep the library  
 170 unchanged. Accepted features are added as  $\Phi^{(r+1)} = \Phi^{(r)} \cup \Delta\Phi^{(r)}$ . We then recompute collisions  
 171 and repeat until the repair budget is exhausted or no collisions remain. The final library is passed to  
 172 the policy learner. The next proposition formalizes the effect of an accepted repair step.

173 **Proposition 1** (Collision-added features strictly reduce irreducible empirical error). *Let  $\Phi$  be a*  
 174 *feature library and let  $\Psi$  be a set of additional features added by an accepted repair step. Define*

$$R^*(\Phi) = \min_{h: \{0,1\}^{|\Phi|} \rightarrow \{0,1\}} \frac{1}{|\mathcal{D}_{\text{sa}}|} \sum_{(s,a,y) \in \mathcal{D}_{\text{sa}}} \mathbf{1}\{h(x_\Phi(s,a)) \neq y\},$$

175 *the minimum empirical classification error achievable by any deterministic classifier whose predic-*  
 176 *tion depends only on  $x_\Phi(s,a)$ . Then*

$$R^*(\Phi \cup \Psi) < R^*(\Phi).$$

177 *A proof is given in Appendix B.*

178 Collision-freeness should be interpreted carefully. If no collisions remain, then  $\mathcal{D}_{\text{sa}}$  is separable at  
 179 the representation level: there exists some deterministic classifier over the observed feature vectors  
 180 that can fit the labels. This does not guarantee generalization, nor does it guarantee that a restricted  
 181 or regularized policy learner will choose a perfect classifier. However, the presence of collisions is  
 182 a concrete certificate that the current feature library is insufficient for fitting the labeled examples.  
 183 More discussion about the expressiveness vs generalizability trade-off is written in Appendix G.

## 184 4 Experiments and Results

185 We evaluate TwoLIP along three main axes: (i) generalization to unseen task instances, (ii) sample  
 186 efficiency with respect to the number of demonstrations, and (iii) composability of the learned fea-  
 187 ture library across tasks. We compare against a set of baselines, analyze quantitative performance,  
 188 and conduct ablation studies to isolate the contribution of individual components of TwoLIP .

### 189 4.1 Experimental Setup

190 **Environments and Tasks** We evaluate TwoLIP on Generalization Grid Games (Silver et al.,  
 191 2020), a suite of strategy tasks with a shared grid-based state and action space. Each state is a  
 192 variable-sized  $h \times w$  grid with cell values from a finite set  $V$ , and each action selects a grid cell,  
 193  $a \in \{1, \dots, h\} \times \{1, \dots, w\}$ . We consider five tasks: *TwoPileNim*, *CheckmateTactic*, *StopTheFall*,  
 194 *ReachForTheStar*, and *Chase*<sup>2</sup>, each with different objectives and strategic requirements. For each

<sup>1</sup>We rank mixed-label buckets by positive–negative pair count and sample from the top-ranked buckets.

<sup>2</sup>Short forms: nim, ct, stf, rfts, and ec, respectively.

## Algorithm 1: TwoLIP

**Require:** demos  $\mathcal{D}$ , negative sampler  $\mathcal{C}^-$ , query budget  $Q$ , rounds  $R$ , particles  $K$   
**Ensure:** feature library  $\Phi$ , posterior  $q$   
1:  $\mathcal{D}_{sa} \leftarrow \text{BUILD DATA}(\mathcal{D}, \mathcal{C}^-)$   
2:  $\Phi \leftarrow \text{INDUCE FEATURES}(\mathcal{D}_{sa}; Q)$   
3: **for**  $r = 1, \dots, R$  **do**  
4:    $(X, y) \leftarrow \text{FEATURIZE}(\mathcal{D}_{sa}, \Phi)$   
5:    $\mathcal{B} \leftarrow \text{FIND COLLISIONS}(X, y)$   
6:   **if**  $\mathcal{B} = \emptyset$  **then break**  
7:   **end if**  
8:    $\Phi \leftarrow \Phi \cup \text{REPAIR FEATURES}(\mathcal{B})$   
9: **end for**  
10:  $(X, y) \leftarrow \text{FEATURIZE}(\mathcal{D}_{sa}, \Phi)$   
11:  $q \leftarrow \text{INFER POSTERIOR}(X, y, \Phi; K)$   
12: **return**  $\Phi, q$

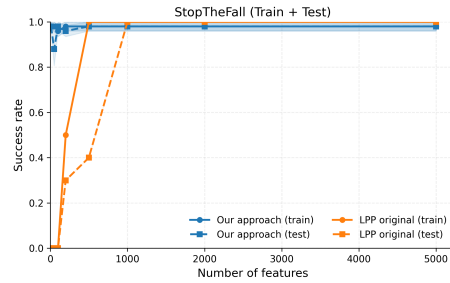


Figure 2: Sample efficiency on STF.

195 task, instances vary in grid size, layout, and object configurations, requiring generalization beyond  
196 the training demonstrations. We use 10 training instances and 10 held-out test instances per task,  
197 with a 200-step episode horizon. Our primary metric is test success rate: the fraction of held-out  
198 instances solved within this horizon. Unless otherwise stated, all methods that rely on LLM calls  
199 use GPT-4.1 (OpenAI, 2025) with temperature 0.0 and top-p 1.0.

200 **4.2 Baselines**

201 We compare TwoLIP against baselines for IL and policy representation.

202 **Fully Convolutional Network (FCN).** We train an FCN (Long et al., 2015) using the same state  
203 and action representations as TwoLIP: eight  $3 \times 3$  convolutional layers with stride 1, padding 1, and  
204 ReLU nonlinearities, exploiting local spatial structure through parameter sharing.

205 **CAP+Demo.** We evaluate an end-to-end LLM baseline inspired by Code as Policies (Liang et al.,  
206 2023), prompting the model to synthesize a Python policy from demonstrations. Unlike TwoLIP, it  
207 infers task structure directly from demonstrations without an explicit feature representation (Austin  
208 et al., 2021; Chen et al., 2021).

209 **VLM-Imitation.** We prompt a vision-language model with demonstration image sequences to  
210 predict actions for new visual states, testing whether multimodal models can perform IL directly  
211 without structured features.

212 **Oracle DSL Baseline.** We include a structured oracle baseline based on the LPP setup (Silver  
213 et al., 2020). This baseline uses a hand-designed DSL tailored to the tasks, so the policy learner  
214 already has access to task-relevant symbolic features. TwoLIP does not assume accessing this DSL;  
215 it has to induce its feature library from demonstrations. We use this comparison to see how close  
216 learned features can get to a structured policy learner that is given the right abstractions in advance.

217 **4.3 Generalization Performance**

218 We first evaluate whether TwoLIP generalizes from few demonstrations. Figure A.2 reports test  
219 success versus the number of training demonstrations on 10 held-out task instances, excluded from  
220 both feature generation and policy learning. Shaded regions show standard error over 10 random  
221 seeds; the training curves are in Appendix C.

222 Across tasks, TwoLIP substantially outperforms the end-to-end baselines and is often competitive  
223 with the oracle DSL baseline. In some tasks, it matches or exceeds the oracle despite not using  
224 the hand-designed DSL. The FCN baseline improves slightly with more demonstrations but  
225 remains well below TwoLIP, suggesting that direct supervised action prediction needs more data to

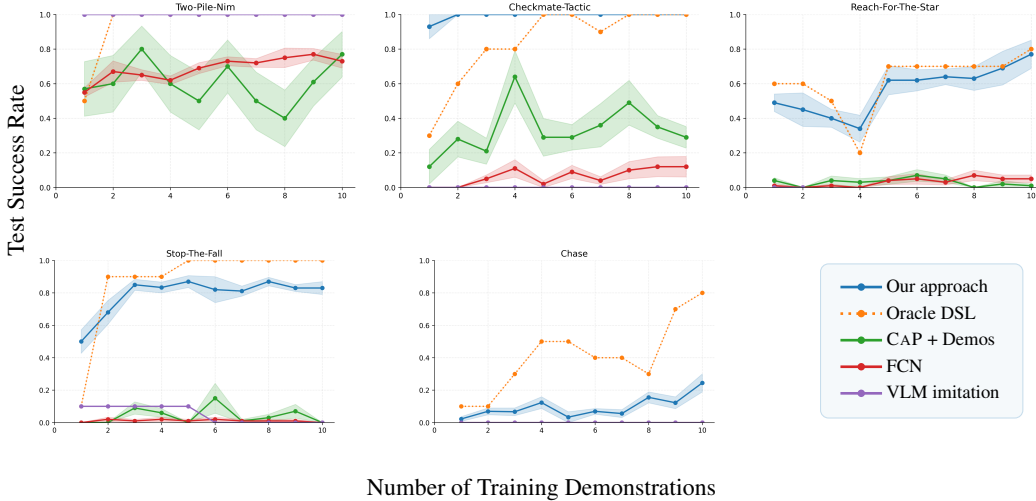
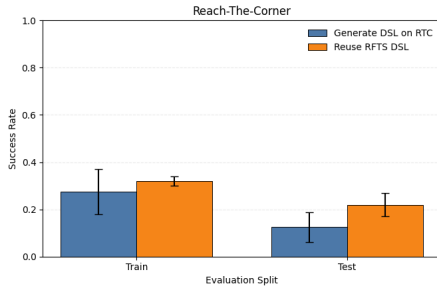
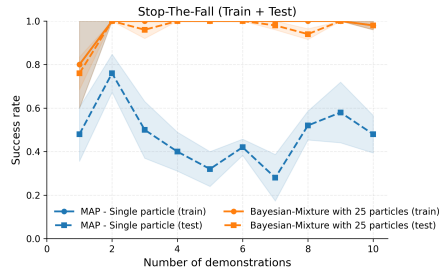


Figure 3: Performance on held-out test task instances as a function of the number of training demonstrations. Each subplot reports mean task-completion success over 10 training seeds.



(a) Transferability on RTC with 10 demos.



(b) Performance of MAP vs. mixture policy.

226 learn the relevant spatial patterns. The LLM-based baselines also perform poorly on most tasks,  
 227 likely because direct policy synthesis requires inferring task logic and reliable spatial reasoning  
 228 from text prompts. The main exception is *TwoPileNim*, whose state structure is simpler than the  
 229 grid-navigation tasks. In one harder domain, TwoLIP remains below the oracle, likely due to the  
 230 task difficulty and the expert strategy, but still outperforms baselines that achieve near-zero success.

231 Overall, these results show that separating feature induction from policy learning lets TwoLIP re-  
 232 cover task structure strong enough to match, and sometimes exceed, a policy learner given a hand-  
 233 designed DSL, while substantially outperforming end-to-end imitation baselines.

#### 234 4.4 Sample Efficiency

235 We next compare how much representation DSL size is needed to reach strong performance. The oracle  
 236 DSL baseline enumerates features from a fixed hand-designed DSL, independent of the demon-  
 237 strations. In contrast, TwoLIP induces features conditioned on the demonstrated behavior, so the  
 238 feature library can focus on distinctions that are useful for the current task. Figure 2 shows that  
 239 TwoLIP reaches performance comparable to the oracle DSL baseline while using far fewer features,  
 240 often by one to two orders of magnitude depending on the environment. This compactness advan-  
 241 tage holds for both training and held-out test instances, suggesting that TwoLIP induces a compact,  
 242 task-relevant feature library rather than relying on a large hand-enumerated representation.

#### 243 4.5 Composability of Learned Features

244 We next test whether induced feature libraries transfer across related tasks. As a preliminary experi-  
 245 ment, we introduce Reach-The-Corner (RTC), a new task in the same environment as RFTS; details  
 246 are in Appendix H. We compare an RTC-specific library to reusing the fixed RFTS library. Figure 4a  
 247 shows that the reused library performs close to the RTC-specific one on both train and held-out test  
 248 instances, suggesting that induced features capture reusable environment structure rather than only  
 249 task-specific correlations. Here we transfer only the initial library without repair; future work will  
 250 study when transferred libraries need repair.

#### 251 4.6 Ablation Studies

252 We ablate two components contributing to TwoLIP’s performance. First, we compare the full poste-  
 253 rior ensemble to the MAP policy. The mixture improves robustness, especially with few demonstra-  
 254 tions, by averaging over multiple plausible policies; Figure 4b shows this result. Second, we study  
 255 collision repair: increasing the repair budget improves train/test accuracy while reducing positive-  
 256 negative collisions, indicating that the initial feature library misses task-relevant distinctions. Perfor-  
 257 mance becomes near-perfect before all collisions are removed, suggesting that repair mainly resolves  
 258 the ambiguities most relevant to the downstream learner. Figure A.3 reports this ablation.

### 259 5 Related Work

260 **IL for robot policies.** Recent robot learning has increasingly moved toward large-scale IL, where  
 261 policies are trained directly from demonstrations. Diffusion policies model multimodal action dis-  
 262 tributions for visuomotor manipulation (Chi et al., 2025), while generalist robot policies and vision-  
 263 language-action (VLA) models leverage broad robot datasets and language/vision pretraining to  
 264 generalize across tasks and embodiments (Open X-Embodiment Collaboration et al., 2024; Octo  
 265 Model Team et al., 2024; Kim et al., 2024; Black et al., 2025). These works show the benefits of  
 266 scaling data, model capacity, and pretraining. In contrast, TwoLIP induces compact state-action  
 267 features and learns a structured policy over them from few demonstrations.

268 **Data-efficient IL.** Prior work on data-efficient IL includes apprenticeship and inverse reinforce-  
 269 ment learning, which infer rewards or costs from demonstrations before solving for a policy (Abbeel  
 270 and Ng, 2004), and meta-imitation, which learns across task distributions to adapt from one or a few  
 271 demonstrations at test time (Duan et al., 2017; Finn et al., 2017a;b). While effective, these methods  
 272 typically rely on many related tasks and demonstrations during meta-training. We instead study  
 273 few-shot target-task imitation, where abstractions must be inferred directly from demonstrations.

274 **Structured and programmatic policy learning.** Another way to improve generalization is to add  
 275 symbolic, relational, or programmatic structure to the policy class. Relational RL and IL use logical  
 276 representations to generalize across objects and states (Džeroski et al., 2001; Natarajan et al., 2011),  
 277 while programmatic policy methods search over interpretable programs (Verma et al., 2018). LPP  
 278 shows that Bayesian inference over symbolic policy structure can support few-shot generalization  
 279 when a task-relevant DSL is available (Silver et al., 2020). These approaches show the value of  
 280 structure, but usually assume that the relevant DSL or features are provided by a human designer.  
 281 TwoLIP keeps the structured policy bias, but induces the feature library from demonstrations.

282 **Language models and executable policy structure.** Recent work uses language models to gener-  
 283 ate executable code for embodied control, most notably CAP (Liang et al., 2023). These methods  
 284 show that LLMs can produce structured control logic, but typically ask the model to generate the  
 285 policy directly. In contrast, we use the LLM to propose candidate state-action features, which are  
 286 executed, filtered, repaired, and passed to a separate policy learner. This separation lets us test and  
 287 improve the representation before committing to a policy.<sup>3</sup>

<sup>3</sup>Discussion and future work are in Appendix F.

288 **References**

- 289 Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning. In  
290 *Proceedings of the Twenty-first International Conference on Machine Learning*, page 1. ACM,  
291 2004. doi: 10.1145/1015330.1015430.
- 292 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
293 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large  
294 language models. *arXiv preprint arXiv:2108.07732*, 2021.
- 295 Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail,  
296 Michael Robert Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya Ghosh, Lachy  
297 Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin  
298 LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z.  
299 Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tan-  
300 ner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky.  $\pi_{0.5}$ :  
301 A vision-language-action model with open-world generalization. In *Proceedings of The 9th Con-  
302 ference on Robot Learning*, volume 305 of *Proceedings of Machine Learning Research*, pages  
303 17–40, 2025.
- 304 Leo Breiman, Jerome H. Friedman, Richard A. Olshen, and Charles J. Stone. *Classification and  
305 Regression Trees*. Wadsworth International Group, 1984.
- 306 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
307 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
308 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 309 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran  
310 Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Proceedings of  
311 Robotics: Science and Systems*, 2023.
- 312 Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ  
313 Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffu-  
314 sion. *The International Journal of Robotics Research*, 44(10–11):1480–1501, 2025. doi:  
315 10.1177/02783649241273668.
- 316 Guofeng Cui, Yuning Wang, Wensen Mao, Yuanlin Duan, and He Zhu. Abstraction refinement-  
317 guided program synthesis for robot learning from demonstrations. *Proceedings of the ACM on  
318 Programming Languages*, 9(OOPSLA2):555–583, 2025. doi: 10.1145/3763070.
- 319 Yan Duan, Marcin Andrychowicz, Bradly Stadie, OpenAI Jonathan Ho, Jonas Schneider, Ilya  
320 Sutskever, Pieter Abbeel, and Wojciech Zaremba. One-shot imitation learning. In *Advances  
321 in Neural Information Processing Systems*, volume 30, 2017.
- 322 Sašo Džeroski, Luc De Raedt, and Kurt Driessens. Relational reinforcement learning. *Machine  
323 Learning*, 43(1–2):7–52, 2001. doi: 10.1023/A:1007694015589.
- 324 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation  
325 of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*,  
326 volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR, 2017a.
- 327 Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imita-  
328 tion learning via meta-learning. In *Proceedings of the 1st Annual Conference on Robot Learning*,  
329 volume 78 of *Proceedings of Machine Learning Research*, pages 357–368. PMLR, 2017b.
- 330 Pete Florence, Corey Lynch, Andy Zeng, Oscar A. Ramirez, Ayzaan Wahid, Laura Downs, Adrian  
331 Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In  
332 *Proceedings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine  
333 Learning Research*, pages 158–168. PMLR, 2022.

- 334 Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth  
335 Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty El-  
336 lis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint*  
337 *arXiv:2403.12945*, 2024.
- 338 Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair,  
339 Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Ben-  
340 jamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn.  
341 Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- 342 Yann LeCun, Sumit Chopra, Raia Hadsell, Marc’Aurelio Ranzato, and Fu Jie Huang. A tutorial on  
343 energy-based learning. In Gökhan Bakir, Thomas Hofmann, Bernhard Schölkopf, Alexander J.  
344 Smola, Ben Taskar, and S. V. N. Vishwanathan, editors, *Predicting Structured Data*. MIT Press,  
345 2006.
- 346 Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and  
347 Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE*  
348 *International Conference on Robotics and Automation (ICRA)*, pages 9493–9500, 2023. doi:  
349 10.1109/ICRA48891.2023.10160591.
- 350 Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic seg-  
351 mentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
352 pages 3431–3440, 2015. doi: 10.1109/CVPR.2015.7298965.
- 353 Sriraam Natarajan, Saket Joshi, Prasad Tadepalli, Kristian Kersting, and Jude Shavlik. Imitation  
354 learning in relational domains: A functional-gradient boosting approach. In *Proceedings of the*  
355 *Twenty-Second International Joint Conference on Artificial Intelligence*, pages 1414–1420, 2011.  
356 doi: 10.5591/978-1-57735-516-8/IJCAI11-239.
- 357 Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep  
358 Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, Jianlan Luo, You Liang Tan, Pannag Sanketi,  
359 Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source  
360 generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- 361 Open X-Embodiment Collaboration, Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek  
362 Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya  
363 Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, An-  
364 chit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, et al. Open x-embodiment: Robotic  
365 learning datasets and rt-x models. In *2024 IEEE International Conference on Robotics and Au-*  
366 *tomation (ICRA)*, pages 6892–6903, 2024. doi: 10.1109/ICRA57147.2024.10611477.
- 367 OpenAI. Introducing gpt-4.1 in the api. <https://openai.com/index/gpt-4-1/>, April  
368 2025. Accessed: 2026-05-12.
- 369 Jyothish Pari, Nur Muhammad Mahi Shafiullah, Sridhar Arunachalam, and Lerrel Pinto. The sur-  
370 prising effectiveness of representation learning for visual imitation. In *Robotics: Science and*  
371 *Systems*, 2022.
- 372 Noah Patton, Kia Rahmani, Meghana Missula, Joydeep Biswas, and Isil Dillig. Programming-by-  
373 demonstration for long-horizon robot tasks. *arXiv preprint arXiv:2305.03129*, 2023.
- 374 Dean A. Pomerleau. ALVINN: An autonomous land vehicle in a neural network. In *Advances in*  
375 *Neural Information Processing Systems*, volume 1, pages 305–313, 1988.
- 376 Stephane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. A reduction of imitation learning and  
377 structured prediction to no-regret online learning. In *Proceedings of the Fourteenth International*  
378 *Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learn-*  
379 *ing Research*, pages 627–635. PMLR, 2011.

- 380 Tom Silver, Kelsey R. Allen, Alex K. Lew, Leslie Pack Kaelbling, and Josh Tenenbaum. Few-shot  
381 bayesian imitation learning with logical program policies. In *Proceedings of the AAAI Conference*  
382 *on Artificial Intelligence*, volume 34, pages 10251–10258, 2020. doi: 10.1609/aaai.v34i06.6587.
- 383 Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri.  
384 Programmatically interpretable reinforcement learning. In *Proceedings of the 35th International*  
385 *Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*,  
386 pages 5052–5061. PMLR, 2018.
- 387 Jimmy Xin, Linus Zheng, Kia Rahmani, Jiayi Wei, Jarrett Holtz, Isil Dillig, and Joydeep Biswas.  
388 Programmatic imitation learning from unlabeled and noisy demonstrations. *arXiv preprint*  
389 *arXiv:2303.01440*, 2023.
- 390 Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual  
391 manipulation with low-cost hardware. In *Proceedings of Robotics: Science and Systems*, Daegu,  
392 Republic of Korea, July 2023. doi: 10.15607/RSS.2023.XIX.016.

393 **6 Appendix**394 **A Bayesian Logical Policy Learning Details**

395 This appendix gives the details of the Bayesian logical policy-learning component used in TwoLIP.

396 **A.1 Posterior over structured classifiers**

397 Given a feature library  $\Phi = \{\phi_1, \dots, \phi_m\}$ , each state–action pair  $(s, a)$  is represented by  $x_\Phi(s, a) \in$   
 398  $\{0, 1\}^m$ . Let  $\mathcal{H}_\Phi$  denote the class of structured classifiers over this representation. Each classifier  
 399  $h \in \mathcal{H}_\Phi$  predicts whether an action should be treated as admissible in a state, with  $h(x_\Phi(s, a)) = 1$   
 400 indicating admissibility and  $h(x_\Phi(s, a)) = 0$  indicating inadmissibility.

401 We define a posterior over classifiers,

$$p(h \mid \mathcal{D}_{\text{sa}}, \Phi) \propto p(\mathcal{D}_{\text{sa}} \mid h, \Phi) p(h \mid \Phi),$$

402 where  $\mathcal{D}_{\text{sa}} = \{(s_i, a_i, y_i)\}_{i=1}^n$  is the contrastive state–action dataset. Positive examples correspond  
 403 to demonstrated actions,  $y_i = 1$ , while negative examples correspond to alternative actions,  $y_i = 0$ .  
 404 The likelihood rewards classifiers that agree with these labels, while the prior favors simpler logical  
 405 policies.

406 **A.2 Soft likelihood with allow/reject penalties**

407 For each labeled example  $(s_i, a_i, y_i)$ , let  $\hat{y}_i = h(x_\Phi(s_i, a_i))$ . A deterministic classifier would either  
 408 agree or disagree with the label. However, because the demonstrations are few and the negative set  
 409 may contain actions that are not necessarily impossible, we use a softened likelihood that penalizes  
 410 different types of mistakes rather than assigning zero probability to any imperfect policy.

411 One convenient form is an exponential penalty model:

$$p(\mathcal{D}_{\text{sa}} \mid h, \Phi) \propto \exp \left( - \sum_{i=1}^n \omega_i [\lambda_{\text{fn}} \mathbf{1}\{y_i = 1, \hat{y}_i = 0\} + \lambda_{\text{fp}} \mathbf{1}\{y_i = 0, \hat{y}_i = 1\}] \right).$$

412 Here,  $\omega_i$  is an optional example weight,  $\lambda_{\text{fn}}$  penalizes rejecting a demonstrated action, and  $\lambda_{\text{fp}}$   
 413 penalizes allowing a negative action. In our setting, false negatives are usually more severe: if  $h$   
 414 rejects the demonstrated action, then the induced policy cannot reproduce the expert choice in that  
 415 state. False positives are also penalized, but they may be treated more softly because some alternative  
 416 actions can be behaviorally plausible, especially when only a few demonstrations are available.

417 Equivalently, this likelihood can be written as a product of per-example probabilities,

$$p(\mathcal{D}_{\text{sa}} \mid h, \Phi) = \prod_{i=1}^n p(y_i \mid h, x_\Phi(s_i, a_i))^{\omega_i},$$

418 where

$$p(y_i = 1 \mid h, x_i) = \begin{cases} 1 - \epsilon_{\text{fn}}, & h(x_i) = 1, \\ \epsilon_{\text{fn}}, & h(x_i) = 0, \end{cases} \quad p(y_i = 0 \mid h, x_i) = \begin{cases} 1 - \epsilon_{\text{fp}}, & h(x_i) = 0, \\ \epsilon_{\text{fp}}, & h(x_i) = 1. \end{cases}$$

419 The parameters  $\epsilon_{\text{fn}}$  and  $\epsilon_{\text{fp}}$  are the softened error probabilities corresponding to the false-negative  
 420 and false-positive penalties. This form avoids making the posterior collapse when no candidate  
 421 classifier perfectly separates the contrastive dataset.

422 **A.3 Structural prior**

423 The prior  $p(h \mid \Phi)$  encodes a preference for compact logical explanations. Since each classifier is  
 424 represented as a logical rule over Boolean features, we penalize structural complexity, such as the  
 425 number of clauses, the number of feature literals, or the depth of the corresponding decision tree. A  
 426 generic form is

$$p(h \mid \Phi) \propto \exp(-\alpha \text{complexity}(h)),$$

427 where  $\alpha \geq 0$  controls the strength of the simplicity bias. This prior discourages policies that fit  
 428 the demonstrations using overly specific conjunctions of features, which are more likely to encode  
 429 accidental correlations in the small demonstration set.

430 In our implementation, the structured classifiers are instantiated using decision trees over the induced  
 431 Boolean features. A path to a positive leaf corresponds to a conjunction of feature literals, and the  
 432 set of positive leaves forms a disjunction. Thus, decision-tree complexity provides a natural proxy  
 433 for logical policy complexity.

434 **A.4 Approximate posterior inference**

435 Exact inference over  $\mathcal{H}_\Phi$  is intractable because the number of logical formulas grows combinato-  
 436 rially with the number of features. We therefore approximate the posterior with a weighted set of  
 437 particles,

$$q(h) = \sum_{k=1}^K w_k \delta_{h_k}(h),$$

438 where each  $h_k$  is a candidate structured classifier and  $w_k$  is its normalized posterior weight.

439 Candidate particles are generated by repeatedly training logical learners on subsets of the induced  
 440 feature library. For a feature subset  $\Phi_\ell \subseteq \Phi$ , we construct the corresponding feature matrix

$$X_\ell = \{x_{\Phi_\ell}(s_i, a_i) : (s_i, a_i, y_i) \in \mathcal{D}_{\text{sa}}\}.$$

441 A decision-tree learner is then fit to  $(X_\ell, y)$ , producing one or more candidate classifiers. Each  
 442 candidate classifier is scored using the posterior objective,

$$\tilde{w}_k = p(\mathcal{D}_{\text{sa}} \mid h_k, \Phi) p(h_k \mid \Phi),$$

443 and the weights are normalized as

$$w_k = \frac{\tilde{w}_k}{\sum_{j=1}^K \tilde{w}_j}.$$

444 This particle approximation allows TwoLIP to maintain multiple plausible policies instead of com-  
 445 mitting to a single explanation of the demonstrations.

446 **A.5 Posterior mixture policy**

447 Each classifier  $h_k$  induces an admissibility prediction for every candidate action. The posterior  
 448 mixture combines these predictions into a soft admissibility score,

$$S_q(s, a) = \mathbb{E}_{h \sim q}[h(x_\Phi(s, a))] = \sum_{k=1}^K w_k h_k(x_\Phi(s, a)).$$

449 At test time, the policy selects an action with maximum posterior expected admissibility,

$$\pi_q(s) \in \arg \max_{a \in \mathcal{A}} S_q(s, a).$$

450 This mixture can be interpreted as a soft state–action scoring model: actions marked admissible by  
 451 many high-posterior policies receive higher scores, while actions rejected by most policies receive  
 452 lower scores. When the few demonstrations underdetermine the correct policy, this mixture reduces  
 453 sensitivity to any single classifier and allows several plausible explanations to contribute to action  
 454 selection.

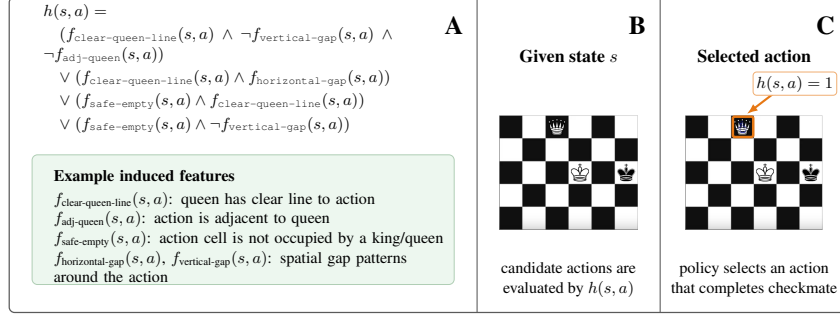


Figure A.1: Example learned policy for Checkmate Tactic. (A) The learned logical policy composes LLM-induced features into a readable rule over state–action pairs. (B) Given a board state, candidate actions are evaluated by the policy. (C) The selected action satisfies the learned rule and solves the instance.

## 455 B Proof of Propositions 1

456 *Proof.* Fix a feature library  $\Phi$ . The representation  $x_\Phi$  partitions  $\mathcal{D}_{\text{sa}}$  into feature buckets, where two  
 457 examples belong to the same bucket if they have the same feature vector. Let  $\mathcal{B}_\Phi$  denote the set of  
 458 buckets induced by  $\Phi$ . For each bucket  $B \in \mathcal{B}_\Phi$ , define

$$n_B^+ = |\{(s, a, y) \in B : y = 1\}|, \quad n_B^- = |\{(s, a, y) \in B : y = 0\}|.$$

459 Any deterministic classifier whose prediction depends only on  $x_\Phi(s, a)$  must assign the same label  
 460 to all examples in a bucket  $B$ . Therefore, the best possible prediction on  $B$  is the majority label, and  
 461 the minimum number of mistakes on that bucket is

$$\min(n_B^+, n_B^-).$$

462 Thus,

$$R^*(\Phi) = \frac{1}{|\mathcal{D}_{\text{sa}}|} \sum_{B \in \mathcal{B}_\Phi} \min(n_B^+, n_B^-).$$

463 Now consider the augmented feature library  $\Phi \cup \Psi$ . Adding features can only refine the original  
 464 partition: each bucket  $B \in \mathcal{B}_\Phi$  is split into sub-buckets  $B_1, \dots, B_m$  under the representation  $x_{\Phi \cup \Psi}$ .  
 465 For these sub-buckets,

$$\sum_{j=1}^m n_{B_j}^+ = n_B^+, \quad \sum_{j=1}^m n_{B_j}^- = n_B^-.$$

466 The contribution of the refined buckets to the optimal empirical error is

$$\sum_{j=1}^m \min(n_{B_j}^+, n_{B_j}^-).$$

467 For nonnegative numbers  $a_j, b_j$ , we have

$$\sum_{j=1}^m \min(a_j, b_j) \leq \min\left(\sum_{j=1}^m a_j, \sum_{j=1}^m b_j\right).$$

468 Applying this with  $a_j = n_{B_j}^+$  and  $b_j = n_{B_j}^-$  gives

$$\sum_{j=1}^m \min(n_{B_j}^+, n_{B_j}^-) \leq \min(n_B^+, n_B^-).$$

469 Therefore, each original bucket contributes no more error after adding  $\Psi$ . Summing over all original  
 470 buckets yields

$$R^*(\Phi \cup \Psi) \leq R^*(\Phi).$$

471 If the repair step is accepted because it reduces the collision-induced empirical-error bound, then for  
 472 at least one original bucket  $B$ ,

$$\sum_{j=1}^m \min(n_{B_j}^+, n_{B_j}^-) < \min(n_B^+, n_B^-),$$

473 while all other buckets contribute no more error than before. Hence the total optimal empirical error  
 474 strictly decreases:

$$R^*(\Phi \cup \Psi) < R^*(\Phi).$$

475

□

## 476 C Additional Experimental Results

477 This appendix includes additional plots complementing the main experiments, including the training  
 478 success curves corresponding to the held-out test results reported in the main text.



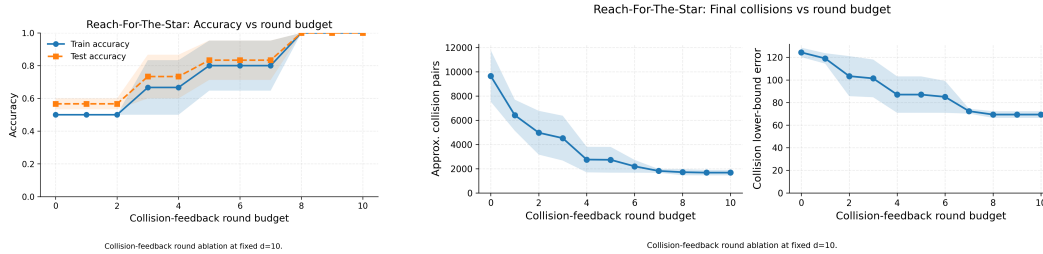
Figure A.2: Training performance as a function of the number of demonstrations. Each subplot reports mean task-completion success on training task instances over 10 training seeds.

## 479 D LLM Prompt Structure

480 TwoLIP uses LLM calls in two parts of the feature-induction pipeline: initial feature-template gen-  
 481 eration and collision-guided feature repair. We do not rely on the LLM to directly output actions or  
 482 policies. Instead, the LLM proposes executable Boolean state-action feature templates, which are  
 483 then validated, filtered, and used by the downstream Bayesian policy learner.

### 484 D.1 Feature-Template Generation Prompt

485 We use the following prompt template to ask the LLM to propose executable Python feature tem-  
 486 plates from expert demonstrations. The full prompt includes task-specific demonstrations and envi-  
 487 ronment summaries; here we show only the high-level structure.



(a) Accuracy vs. repair rounds

(b) Collision metrics vs. repair rounds

Figure A.3: Effect of collision-guided repair on *RFTS* with  $d = 10$ .

## 488 E Prompt Templates

489 We show simplified prompt templates used by TwoLIP . Task-specific demonstrations, collision  
 490 examples, and implementation details are replaced by placeholders.

### 491 E.1 Feature-Template Generation Prompt

#### Feature-Template Generation Prompt

##### SYSTEM

You are an expert feature-library designer for program synthesis and Logical Programmatic Policies in grid-based games.

Your task is to read expert demonstrations  $(s, a, s')$  and propose a compact, expressive library of Boolean feature templates  $f_i(s, a)$ .

The goal is not to enumerate many variants, but to introduce the minimal set of structurally distinct predicate families necessary to express the expert’s decision rule.

##### ENVIRONMENT

{description of the grid state, clicked-cell action, transition dynamics, and feature-function interface}

##### TEMPLATE CONSTRAINTS

{allowed placeholders, token-role rules, action validation logic, and examples of forbidden patterns such as hard-coded tokens or duplicate complements}

##### DESIGN PRINCIPLES

{features should be structurally expressive, minimal, general, scale-invariant, action-grounded, and non-redundant}

##### INPUT: DEMONSTRATIONS

{expert trajectories summarized with board structure, local views, ray features, distances, and transition effects}

##### TASK

Infer the structural distinctions that separate expert actions from plausible alternatives, generate feature templates that encode those distinctions, and stop once semantic coverage is sufficient.

##### OUTPUT FORMAT

Return only valid JSON:

```
{
  "features": [
    { "id": "...", "name": "...", "source": "..." }
  ]
}
```

492

493 **E.2 Collision-Repair Prompt****Collision-Repair Prompt**

494

**SYSTEM**

You are an expert feature-library repair agent for Logical Programmatic Policies in grid-based games. Your task is to propose new Boolean feature functions  $f_i(s, a)$  that repair feature collisions. A collision occurs when positive and negative examples have identical feature vectors under the current feature library, but different labels.

**OBJECTIVE**

Introduce genuinely new structural distinctions that separate the colliding examples and generalize beyond the provided collision bucket.

**ENVIRONMENT**

{description of the grid state, clicked-cell action, allowed token constants, and feature-function interface}

**IMPLEMENTATION CONSTRAINTS**

{allowed tokens, action validation logic, no imports, no raw strings, no unsupported helpers}

**ANTI-OVERFITTING CONSTRAINTS**

{no hard-coded coordinates, no board-size assumptions, no dataset-specific thresholds, no scan-order anchors, no memorization}

**PREFERRED STRUCTURAL FEATURES**

{ray visibility, proximity, alignment, reachability, boundary relations, gaps, between-ness, or other structural relations}

**COLLISION EVIDENCE**

{one or more collision buckets containing positive examples, negative examples, local/ray/global summaries, and difference hints}

**BUCKET HANDLING**

Analyze each collision bucket independently. If multiple buckets reveal the same missing structural distinction, propose a feature that resolves them jointly.

**TASK**

Compare positive and negative examples, identify missing structural invariants, and propose new non-redundant features that break the collision while generalizing to new boards and placements.

**OUTPUT FORMAT**

Return only valid JSON:

```
{
  "features": [
    {
      "id": "...",
      "name": "...",
      "description": "...",
      "source": "..."
    }
  ]
}
```

495

496 **F Discussion/Future Work**

497 This work focuses on the role of inductive bias in few-demonstration IL. Our goal is not to argue for  
 498 a particular optimization algorithm, but to show that, in settings where demonstrations are limited,  
 499 it can be useful to first learn a task-relevant feature representation and then learn a policy over that  
 500 representation. This gives the learner a more structured hypothesis space than end-to-end action  
 501 prediction, while still avoiding the need for a manually designed feature set. At the same time, this  
 502 creates an important tradeoff: the feature-generation process must be expressive enough to make the

503 task learnable, but constrained enough to avoid simply memorizing the demonstrations. We include  
504 additional practical details on how we balance this tradeoff in Appendix G.

505 A common challenge in IL is compounding error: small mistakes made by the learned policy can  
506 move the agent into states that were not covered by the expert demonstrations, where further mis-  
507 takes become more likely (Ross et al., 2011). This was not a major issue in the grid-game domains  
508 we study here, since the policy chooses among discrete candidate actions and the action space is  
509 relatively controlled. However, we expect this issue to become more important in future continuous-  
510 action settings. In continuous control, small action errors can accumulate over time and push the  
511 agent farther from the demonstrated trajectories. Extending TwoLIP to these settings may therefore  
512 require interactive data collection, recovery demonstrations, or planners that can help correct local  
513 mistakes.

514 A second direction is to move beyond generating task-specific features from scratch. In this work,  
515 the LLM proposes a feature library for each task, and we then repair that library using collision  
516 feedback. A natural next step is to treat the feature library itself as a reusable object. Features  
517 learned for one task could initialize the library for a related task, and new demonstrations could be  
518 used to add, revise, or remove features over time. This suggests a more lifelong or continual-learning  
519 version of TwoLIP, where the learner gradually builds a library of reusable abstractions rather than  
520 restarting feature induction for every new task.

521 Another interesting direction is to use the learned policies as components inside a planner. The poli-  
522 cies learned by TwoLIP are not just black-box action predictors; they define interpretable conditions  
523 under which actions are admissible. This makes them potentially useful as skills, action filters, or  
524 proposal mechanisms for planning. For example, a planner could use a learned policy to restrict  
525 the action space to semantically plausible actions, while still searching over longer-horizon conse-  
526 quences. This could combine the generalization benefits of structured imitation with the robustness  
527 of planning.

528 Finally, our results suggest that LLMs may be especially useful not as direct policy generators, but  
529 as tools for proposing and revising abstractions. Instead of asking the model to solve the entire con-  
530 trol problem end-to-end, TwoLIP asks it to generate candidate features, tests those features against  
531 demonstrations, and uses failures as feedback for repair. This separation gives the learning algo-  
532 rithm a way to inspect, filter, and improve the representation before committing to a policy. We  
533 believe this is an important direction for using language models in IL: not as standalone agents, but  
534 as sources of structured hypotheses that can be grounded and tested through interaction with data.

## 535 G Expressiveness vs. Generalization

536 The proposed feature-generation pipeline introduces a tradeoff between *expressiveness* and *gener-*  
537 *alization*. These two forces correspond to different goals in the learning process.

538 **Expressiveness.** Feature collisions indicate that the current representation is insufficient. Recall  
539 that a collision occurs when two state–action pairs have the same feature vector but different labels:

$$x(s, a) = x(s', a') \quad \text{but} \quad y(s, a) \neq y(s', a').$$

540 In this case, no deterministic classifier

$$\pi : \{0, 1\}^m \rightarrow \{0, 1\}$$

541 can perfectly fit the data. Collision detection therefore identifies violations of the condition

$$\exists \pi \in \mathcal{H} \quad \text{such that} \quad \pi(x_i) = y_i \quad \forall i.$$

542 To resolve collisions, the system generates additional features. If  $\phi_m$  is the current feature map and  
543  $\phi_{m+1}$  is the augmented map after adding a new feature, then collision repair expands the effective

544 hypothesis class:

$$\mathcal{H}_m \subseteq \mathcal{H}_{m+1}.$$

545 Thus, collision-driven feature generation improves learnability: once collisions are removed, the  
546 training data can in principle be separated by some classifier in the expanded feature space.

547 **Generalization.** However, eliminating collisions does not guarantee generalization. If the feature  
548 map becomes too expressive, the hypothesis class can memorize the training set. For example, once  
549 every positive example has a unique feature vector, one can construct the DNF

$$\pi(x) = \bigvee_{i: y_i=1} \left( \bigwedge_{j: x_{ij}=1} f_j \wedge \bigwedge_{j: x_{ij}=0} \neg f_j \right),$$

550 which fits each positive training example individually. In this case, the empirical risk may go to  
551 zero,

$$\hat{R}_{\text{train}}(\pi) \rightarrow 0,$$

552 while the true risk may still be high:

$$R(\pi) = \mathbb{E}_{(s,a) \sim \mathcal{D}_{\text{test}}} [\mathbf{1}[\pi(s, a) \neq y(s, a)]].$$

553 This can happen when newly generated features capture false shortcut correlations that are predictive  
554 in the demonstrations but unstable across test configurations.

555 **Controlling the tradeoff.** A natural direction for future work is to make this expressiveness–  
556 generalization tradeoff explicit in the learner. One approach is to detect and reject false shortcut  
557 features, such as predicates that rely on overly specific coordinates, object counts, or demonstration-  
558 specific artifacts. Another is to assign complexity costs to features and use these costs as part of the  
559 structural prior, biasing inference toward more stable and reusable predicates.

560 A second direction is to filter features across demonstrations. After generating a feature library, one  
561 can score features by how consistently they are useful across different demonstrations, and remove  
562 features that only activate for a small number of examples. This would discourage features that  
563 explain isolated demonstrations but do not capture the broader task structure.

564 Finally, the hypothesis class itself can be restricted to reduce memorization. For decision-tree learn-  
565 ers, this can be done through standard capacity controls such as maximum depth, minimum leaf size,  
566 or cost-complexity pruning. Similarly, collision repair can be stopped early instead of continuing  
567 until all training collisions are removed, since fully eliminating collisions may over-specialize the  
568 feature space. Together, these mechanisms provide ways to balance the need for expressive features  
569 with the need for policies that generalize.

## 570 H Task Details

571 All tasks use the Generalization Grid Games interface: each state is an  $h \times w$  grid with discrete  
572 object types, and each action selects a grid cell  $a \in \{1, \dots, h\} \times \{1, \dots, w\}$ . The transition is  
573 deterministic, and an episode is successful if the task goal is reached within the horizon.

574 **TwoPileNim.** *TwoPileNim* is a spatial variant of Nim. The grid contains token cells arranged into  
575 two piles, and an action removes tokens according to the game dynamics. The goal is to select  
576 moves that follow the expert Nim strategy and reach a winning terminal configuration.

577 **CheckmateTactic.** *CheckmateTactic* is a chess-tactics task represented on a grid. The state con-  
578 tains chess pieces, and the agent must select the move or target cell that progresses toward a check-  
579 mating configuration. Success requires recognizing the relevant spatial relations between pieces,  
580 such as attacks, alignment, and king safety.

581 **StopTheFall.** *StopTheFall* contains gravity-like dynamics in which objects may fall unless the  
582 agent intervenes. The goal is to click cells that prevent the falling object from reaching an unsafe  
583 configuration, for example by placing or selecting supporting cells. The task requires reasoning  
584 about vertical structure, support, and future motion.

585 **ReachForTheStar.** *ReachForTheStar* is a grid-construction/navigation task in which the agent  
586 must reach a star. Actions modify the grid so that the agent can make progress toward the star, often  
587 requiring intermediate structures rather than simply selecting the goal cell. The task tests whether  
588 the learner can infer reusable spatial concepts such as support, adjacency, connectivity, and progress  
589 toward the target.

590 **Chase.** *Chase* requires the agent to catch or reach a moving target while respecting the grid layout.  
591 The expert behavior often involves intermediate control decisions before directly pursuing the target.  
592 This makes the task challenging because successful behavior depends not only on distance to the  
593 target, but also on the semantics of action cells, obstacles, and the current stage of the strategy.

594 **Reach-The-Corner.** *Reach-The-Corner* (RTC) is a new task introduced for the transfer experi-  
595 ment. It uses the same environment family as *ReachForTheStar*, but changes the goal from reaching  
596 a star to reaching a designated corner of the grid. This preserves much of the underlying environment  
597 structure while changing the task objective, making it useful for testing whether features induced on  
598 *ReachForTheStar* transfer to a related task.