

---

# Concept-based Understanding of Emergent Multi-Agent Behavior

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

1 In this work, we propose novel techniques that improve understanding of emergent  
2 multi-agent behavior via concept-based interpretability. We propose Concept  
3 Bottleneck Policies (CBPs) as a method for learning intrinsically interpretable,  
4 multi-agent policies. By constructing each agent’s policy to be conditioned on its  
5 own estimates of human-understandable concepts, our method also enables novel  
6 forms of behavioral and inter-agent social analysis via concept intervention. For  
7 example, we demonstrate experimentally that our method can reliably detect: (i) the  
8 emergence of coordination, coordination failures (lazy agents), and environmental  
9 demands for coordination in cooperative environments; (ii) the emergence of skill-  
10 based strategy and role assignment in competitive environments; and (iii) emergent  
11 social behaviors such as public resource allocation and exploitation in social  
12 dilemmas. Moreover, CBPs match the performance of standard, non-concept-based  
13 policies; thereby achieving interpretability without sacrificing performance.

## 14 1 Introduction

15 Multi-agent learning continues to play a crucial role in the development of scalable and generally-  
16 capable AI systems. In addition to well-known successes in board games [32, 35, 36] and online  
17 games [6, 44], multi-agent learning has enabled agents to develop a range of capabilities, such as  
18 navigating social dilemmas [20], optimizing traffic flow for autonomous vehicles [43, 45], and even  
19 learning to cooperate with humans [7]. For all its success, interpreting emergent multi-agent behavior  
20 remains an open challenge. In cooperative environments, for example, reward alone is not enough to  
21 understand the nature of a learned coordination strategy, or whether agents have learned to coordinate  
22 at all. In more complex settings, such as social dilemmas, agents often learn nuanced relationships  
23 that include both low-level spatial (e.g., navigation) and high-level social (e.g., exploitation, public  
24 resource allocation) interactions that can only be fully-understood through exhaustive visualization.

25 For this reason, recent work has shifted focus from traditional measures of performance (i.e. reward,  
26 human analysis) to better understanding emergent behaviors [30]. Related methods perform interpretabili-  
27 ty analysis in a *post-hoc* manner, either by measuring basic behavioral statistics [23], or  
28 regressing concepts from policy network activations [27]. An alternative to post-hoc analysis is *intrin-*  
29 *sic interpretability*, where a model’s decisions incorporate human-understandable concepts directly.  
30 In supervised learning settings, for example, it has been shown that it is possible to train a neural  
31 network that is “bottlenecked” by human-understandable concepts, enabling intrinsic interpretability  
32 while maintaining high performance in classification tasks [19].

33 In this work, we introduce an intrinsically-interpretable, concept-based policy architecture for multi-  
34 agent learning. Our proposed architecture, the Concept Bottleneck Policy (CBP), forces an agent to  
35 make decisions by first predicting an intermediate set of human-understandable concepts, then using

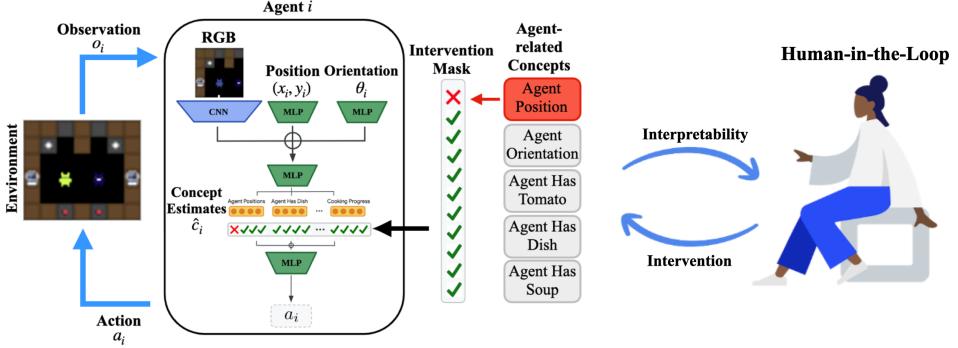


Figure 1: Concept Bottleneck Policies distill agent decisions into human-understandable concepts, yielding interpretable policies for MARL and enabling behavioral analysis via concept intervention.

36 those concepts to select actions. This architecture makes decision-making immediately transparent,  
 37 as agents reveal the environmental and inter-agent factors influencing each of their actions.

38 Beyond intrinsic interpretability on a per-agent basis, our method also **unlocks a novel class of**  
 39 **techniques for understanding emergent multi-agent behavior**. By conditioning each agent’s  
 40 actions on its own concepts from construction, our method supports behavioral analysis via concept  
 41 intervention. Specifically, we can perform interventions over each concept estimate in an agent’s CBP  
 42 and examine its impact on the larger multi-agent system (with respect to rewards, actions, environment  
 43 features, etc). We use this technique to uncover many key aspects of multi-agent behavior.

44 In cooperative environments, our method can reliably detect agents who have learned to coordinate  
 45 from those that act independently, determine which environments require coordination to solve  
 46 vs. those that permit independent solutions, expose inter-agent factors that drive coordination,  
 47 and even diagnose common failure modes such as lazy agents. In competitive environments, our  
 48 method disentangles more sophisticated strategic behavior (e.g., role assignment) as it emerges during  
 49 training. In social dilemmas, we show that it is possible to identify inter-agent social dynamics (e.g.,  
 50 exploitation, public resource sharing) by learning a graph over intervention outcomes. Specifically, we  
 51 learn a sparse graph representing the pairwise relationships of inter-agent intervention effects and show  
 52 that it captures both high-level behavioral patterns (free loader agents vs. public service agents) and  
 53 low-level spatial information (agents who interact). Finally, we show that CBPs perform comparably  
 54 to non-concept based policies, thus achieving interpretability without sacrificing performance.

55 In sum, our contributions are as follows: (i) We introduce Concept Bottleneck Policies as an  
 56 interpretable, concept-based architecture for MARL; (ii) We introduce a suite of novel techniques for  
 57 multi-agent behavioral understanding that leverage concept intervention and graph learning; (iii) We  
 58 show that CBPs can match the performance of non-concept-based network architectures.

## 59 2 Related Work

60 Our work lies at the intersection of interpretability and MARL. In interpretability, computer vision  
 61 works have used saliency maps to help provide pixel-level explanations for model predictions  
 62 [37, 38, 42, 39]. Interpretability using grounded concepts has also been explored to provide more  
 63 meaningful, human-understandable explanations of model decisions [18, 14, 5, 13, 8, 48, 40]. As  
 64 described earlier, concept bottleneck models [19] constrain the networks through such grounded  
 65 concepts, enabling analysis of the model via intervention.

66 RL interpretability techniques include those that either increase the transparency of agent decision-  
 67 making directly [4, 25]; or analyze agent behaviors from a post-hoc perspective [47, 2, 34, 24,  
 68 17, 30, 27, 11]. Approaches that directly increase transparency include those that combine model  
 69 compression and decision trees to yield more interpretable agent policies [4], or rely upon causal  
 70 decision trees with limited depth to explore causality of agent decisions [25]. However, these  
 71 approaches have not yet been extended to RL settings with complex observations (e.g., images),  
 72 where function approximators are typically needed. Post-hoc methods have used traditional saliency-  
 73 mapping techniques to attribute decisions to image observation regions [47], highlighting states  
 74 that lead to major differences in agent behaviors [2], and summarized key agent behaviors using

75 measures such as action uncertainty [34]. Increasingly, natural language has been used to aid human  
76 understanding of agents through instructions [1] or explanations [15].

77 Interpretability of MARL agents has received limited attention in prior works, with primary investigations  
78 gauging agents’ internal representations by making predictions about future outcomes [24],  
79 visualizing latent clusterings of agents’ neural activation vectors [17, 26], or using offline behavioral  
80 analysis to learn behavioral spaces over agents [30]. Recent work has also analyzed decision-making  
81 in two-player games such as Chess [27], and Hex [11]. In contrast to our work, these approaches focus  
82 on post-hoc interpretability, whereas ours makes agent decisions directly transparent and enables  
83 novel forms of downstream analysis. To the best of our knowledge, ours is the first work to identify  
84 coordination, lazy agents, emergent skill learning, and social behaviors via intervention.

### 85 3 Background

86 **Markov Games** A partially-observable Markov game [22] of  $N$  agents is defined by the tuple  
87  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O})$  where  $\mathcal{S}$ ,  $\Omega$ , and  $\mathcal{A}$  are the game’s global state-space, joint observation-  
88 space, and joint action-space, respectively. In each state, each agent  $i$  selects an action  $a_i \in \mathcal{A}_i$ ,  
89 yielding a joint action  $\mathbf{a} = (a_1, \dots, a_N)$  for all agents. Following action selection, the environment  
90 transitions to a new state according to the transition function  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ , produces new  
91 observations for each agent following the observation function  $\mathcal{O} : \mathcal{S} \times \mathcal{A} \rightarrow \Omega$ , and emits a reward  
92 defined by the reward function  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^N$ . The actions of each agent are dictated by a policy  
93  $\pi_i(a_i|o_i)$  and the collection of all individual policies  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_N)$  is called the joint policy.

94 **Concept Bottleneck Models** CBMs are a class of intrinsically-interpretable neural network archi-  
95 tecture for supervised learning [19]. CBMs make predictions by first estimating a set of human-  
96 understandable concepts, then producing an output based on those concepts—i.e., the network is  
97 “bottlenecked” by concepts. Formally, given a dataset  $\{(x^{(j)}, y^{(j)}, c^{(j)})\}_{j=1}^n$  consisting of inputs  
98  $x \in \mathbb{R}^d$ , outputs  $y \in \mathbb{R}$ , and human-understandable concepts  $c \in \mathbb{R}^k$  (where  $k$  is the number of unique  
99 concepts), a concept bottleneck model learns two mappings:  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  from input-space to concept-  
100 space, and  $g: \mathbb{R}^k \rightarrow \mathbb{R}$  from concept-space to output-space. The model can then make predictions  
101  $\hat{y} = f(g(x))$  as a composition of those mappings. For example, in an application of arthritis detection  
102 from MRI images, the model first predicts related concepts such as knee joint spacing and the  
103 presence of bone spurs, then uses those concepts to classify arthritis severity.

### 104 4 Concept Bottleneck Policies for MARL

105 Here we introduce Concept Bottleneck Policies (CBPs) as an intrinsically interpretable, concept-based  
106 policy learning method for MARL (see Figure 1). First, we extend the Markov game formalism to  
107 support concept-based learning. Then we introduce the CBP architecture and show how it fits within  
108 any policy learning scheme. Finally, we demonstrate our method’s potential for interpreting emergent  
109 multi-agent behavior, introducing multiple behavioral analysis techniques that use intervention to  
110 detect a variety of multi-agent social phenomena.

111 **Concept-based Markov Games** We extend Markov games in two important ways. First, we assume  
112 that, in addition to a state space  $\mathcal{S}$ , there exists an interpretable concept-state space  $\mathcal{C}$ , where each  
113 concept-state  $c \in \mathcal{C}$  is a vector of human-understandable concepts that describe key features of the  
114 environment (e.g., agent positions, object states). This assumption holds in many RL domains like  
115 Atari [29] and Melting Pot [21]. Second, we allow agents to generate estimates  $\hat{c}$  of the concept state  
116 as part of their internal representation. In sum, at each time-step  $t$ , the environment produces both a  
117 state  $s_t$  and concept-state  $c_t$ , and each agent  $i$  selects both actions  $a_{i,t}$  and concept estimates  $\hat{c}_{i,t}$ .

118 **Concept Bottleneck Architecture** There are many ways in which concepts can be modeled within  
119 the RL framework. Inspired by CBMs [19], we examine policy architectures in which an agent’s  
120 action is conditioned entirely on its own concept estimates. To this end, we factorize an agent  $i$ ’s  
121 policy  $\pi_i(a_i|o_i)$  into a function  $f_i: \mathcal{O}_i \rightarrow \mathcal{C}_i$  mapping observations  $o_i$  to concept estimates  $\hat{c}_i$ ; and  
122  $\pi_i^{\text{act}}: \mathcal{C}_i \rightarrow \mathcal{A}_i$  mapping concept estimates  $\hat{c}_i$  to actions  $a_i$ . Composing  $f_i$  and  $\pi_i^{\text{act}}$  yields a standard  
123 policy mapping observations to actions:  $\pi_i(a_i|o_i) = \pi_i^{\text{act}}(f_i(o_i)) = \pi_i^{\text{act}}(a_i|\hat{c}_i)$ . Given an observation  
124  $o_{i,t}$  for some agent  $i$  and time  $t$ , the bottleneck first produces concept estimates  $\hat{c}_{i,t} = f_i(o_{i,t})$ , then  
125 uses those estimates alone to select an action  $a_{i,t} \sim \pi_i^{\text{act}}$ . Conditioning actions on concepts creates an  
126 intrinsically interpretable policy, as the policy is forced to provide a human-understandable rendering

127 of the factors driving its decision making. Importantly,  $\hat{c}_i$  can be also used to interpret multi-agent  
 128 behavior. For example, if two agents  $i$  and  $j$  collide while moving in the environment, examining  $\hat{c}_i$   
 129 and  $\hat{c}_j$  may identify which of the agents incorrectly modeled the location of its teammate.

130 **Concept Bottleneck Learning** Under the CBP framing, learning a strong but interpretable policy re-  
 131 quires both (i) learning to predict concepts accurately, and (ii) learning to select actions from those con-  
 132 cepts effectively. To achieve the former, we introduce a concept loss:  $L_C(\mathbf{c}, \hat{\mathbf{c}}) = \sum_{j=1}^{|C|} L_{C_j}(c_j, \hat{c}_j)$ ,  
 133 where each component  $L_{C_j}$  measures the error between the  $j$ 'th predicted concept and its concept  
 134 label. Each  $L_{C_j}$  is a supervised loss, with a specific form dictated by the value of  $c_j$ —mean-squared  
 135 error if  $c_j$  is scalar, log loss if  $c_j$  is binary, etc. In practice, each pair of concepts  $\mathbf{c}$  and concept  
 136 estimates  $\hat{\mathbf{c}}$  are stored in an agent's replay buffer during training alongside standard  $(s, a, r, s')$  tuples.

137 To learn a policy from concept estimates, we attach  $L_C$  as an auxiliary loss to the reward-based loss  
 138 defined by any base RL algorithm (e.g., PPO [33], etc). In general, if  $L_{RL}$  is a generic reward-based  
 139 loss, we construct a joint concept bottleneck policy loss  $L_{CBP} = L_{RL} + \lambda L_C$ , where the concept  
 140 loss coefficient  $\lambda$  weights the relative importance of concept prediction.

141 **Behavioral Analysis via Concept Intervention** Intervention is a powerful tool for understanding  
 142 the statistical relationships that exist between random variables (e.g. in a graphical model) by  
 143 conducting experiments in which the value of a random variable is fixed in order to observe that  
 144 variable's effect on other variables[31]. A key feature of our method is its support of this type of  
 145 intervention analysis. For an agent  $i$ , let  $\text{Int}(\hat{c}_j, \bar{c}_j, i)$  be an intervention over the  $j$ 'th concept estimate  
 146  $\hat{c}_j$  from  $i$ 's bottleneck with the replacement value  $\bar{c}_j$  (in the simplest case,  $\bar{c}_j = 0$ ). We can then  
 147 observe the effect, if any, that  $\bar{c}_j$  has on the agent's behavior by comparing its impact on reward:

$$\mathbb{E} \left[ \sum_t r(s_t, \pi_i^{\text{act}}(\hat{c}_j)) \right] - \mathbb{E} \left[ \sum_t r(s_t, \pi_i^{\text{act}}(\bar{c}_j)) \right] \quad (1)$$

148 In mutual reward settings, this serves as a proxy for the intervention's impact on team behavior.<sup>1</sup>  
 149 Here we propose two examples of behavioral tests enabled by this approach:

150 **Detecting Coordination:** Let  $i$  and  $j$  be two agents. For  $i$  to coordinate with  $j$ ,  $i$  must condition its  
 151 policy on information about  $j$  ( $j$ 's position, orientation, etc). If the agents are coordinating and we  
 152 remove  $i$ 's concept estimates pertaining to  $j$ , we should expect a decrease in reward. Conversely, if  
 153 reward does not degrade, then  $i$  and  $j$  must not be explicitly coordinating (i.e., directly using signals  
 154 from each other). By intervening over concepts pertaining to an agent's teammates, therefore, we can  
 155 identify whether or not agents are coordinating.

156 **Exposing Lazy Agents:** Here we define a lazy agent as one that does not contribute to increasing team  
 157 reward through its own actions. For an agent to contribute, it must at least condition its policy on  
 158 information about its own interactions with the environment. A lazy agent, therefore, is one that has  
 159 learned a sub-optimal policy that does not encode such information, leading to unproductive behavior.  
 160 It follows that we can test for the *degree of laziness* of an agent by replacing its concept estimates  
 161 *about itself* and measuring the resulting degradation of team reward. If performance remains the  
 162 same with the agent incapacitated, we conclude that it is a lazy agent.

163 **Modeling Inter-Agent Social Dynamics** Here we show that it is possible to deepen our intervention-  
 164 based interpretability by learning a sparse graph over CBP intervention outcomes. Specifically, we  
 165 take each intervention  $\text{Int}(\hat{c}_j, \bar{c}_j, i)$  to be a random variable that is described by some agent-related  
 166 features, such as the average reward  $\mathbf{r}_{\hat{c}_j} = \{r_{\hat{c}_j}^1, \dots, r_{\hat{c}_j}^n\}$  each agent receives under the intervention<sup>2</sup>.  
 167 By performing  $\text{Int}(\cdot)$  iteratively over each concept  $c \in C$  and for each agent  $i \in \{1, \dots, N\}$ , we can  
 168 construct a matrix  $\mathbf{X}$  wherein each row is the outcome of an intervention:

$$\mathbf{X} = \begin{bmatrix} \text{Int}(\hat{c}_j, \bar{c}_j, 1) \\ \text{Int}(\hat{c}_k, \bar{c}_k, 2) \\ \vdots \\ \text{Int}(\hat{c}_l, \bar{c}_l, N) \end{bmatrix} = \begin{bmatrix} r_{\hat{c}_j}^1 & r_{\hat{c}_j}^2 & \dots & r_{\hat{c}_j}^n \\ r_{\hat{c}_k}^1 & r_{\hat{c}_k}^2 & \dots & r_{\hat{c}_k}^n \\ \vdots & & & \\ r_{\hat{c}_l}^1 & r_{\hat{c}_l}^2 & \dots & r_{\hat{c}_l}^n \end{bmatrix}$$

<sup>1</sup>We carefully examine the possibility of OOD examples and propose ways to test this in Appendix C.4.

<sup>2</sup>Outcomes can also be defined in terms of other features (e.g., resources collected, agent proximity, etc.)

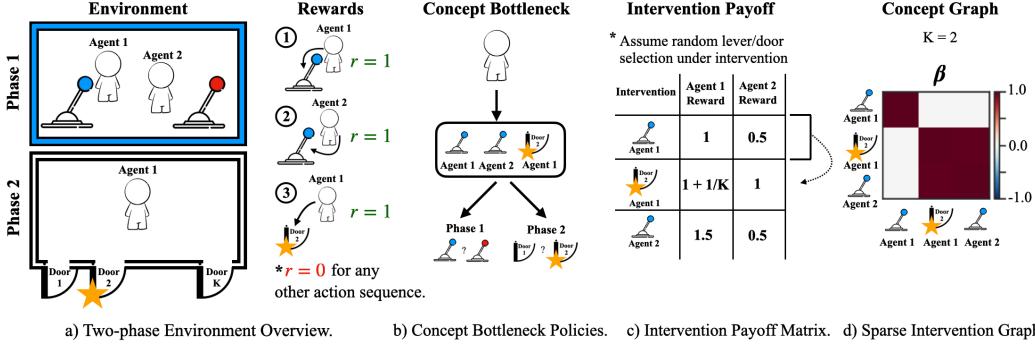


Figure 2: **a)** A two-phase game. In phase one, Agent 1 and 2 choose from two levers (red, blue) using a special observation (wall color) and receive individual rewards for selecting the correct lever (blue) in the correct sequence (Agent 1 moves first, then Agent 2). If done correctly, the agents move on to the second phase. In phase two, Agent 1 selects from  $K$  doors, also using a special observation (gold star), and receives an individual reward for selecting the correct door. **b)** Both agents use a CBP: Agent 1 has concepts for both lever color and the correct door, while Agent 2’s estimates lever color alone. **c)** It is possible to compute each agent’s expected payoff under concept intervention. **d)** The graph obtained by running Equation (2) over pairs of concept interventions from the payoff matrix. Each edge is a statistical relationship between two interventions (color = edge weight)

169 We can then learn a graph that encodes the pairwise relationships between interventions by performing  
170 Lasso neighborhood selection [28] over  $\mathbf{X}$ . Specifically, for the intervention  $\text{Int}(\hat{c}_j, i)$ , we solve:

$$\min_{\beta_{ij}} \|\mathbf{X}_{ij} - \mathbf{X}_{\setminus ij} \beta_{ij}\|_2^2 + \alpha \|\beta_{ij}\|_1 \quad (2)$$

171 where  $\mathbf{X}_{ij}$  is the outcome of  $\text{Int}(\hat{c}_j, \bar{c}_j, i)$ ,  $\mathbf{X}_{\setminus ij}$  represents the remaining  $N \times |C|-1$  interventions,  
172 and  $\alpha$  is an  $L_1$ -penalty, enforcing sparsity. Under this graphical interpretation, the coefficient  
173 vector  $\beta_{ij} \in \mathbb{R}^{n-1}$  is used to establish edges between intervention outcomes. Crucially, non-zero  
174 edge weights in  $\beta_{ij}$  mark the similarity of  $\text{Int}(\hat{c}_j, i)$ ’s outcome to other interventions. Outcome  
175 similarities may hint at the nature of agent behavior. We provide a comprehensive overview of Lasso  
176 neighborhood selection in Appendix D.

177 *Illustrative Example:* To build intuition, we present the following mathematical example. Consider  
178 the two-phase sequential game in Figure 2a. In the Phase 1 of this game,  $N = 2$  agents spawn in a  
179 room that contains two levers (red and blue). One lever is set as a reward-giving lever (blue here)  
180 and is identified as such by an observable feature (wall color). Agent 1 acts first and must choose a lever  
181 to pull, receiving an individual reward ( $r = 1.0$ ) for selecting the correct lever. Agent 2 then acts,  
182 also receiving  $r = 1.0$  for selecting the correct lever. If both agents select the correct lever, the game  
183 proceeds to Phase 2; otherwise it is terminated. In Phase 2, Agent 1 spawns in a new room containing  
184  $K$  unique doors (Agent 2 does not participate). Agent 1 must select a door to exit the room. One  
185 reward-giving door produces an individual reward ( $r = 1.0$ ), while the other  $K-1$  doors produce  
186 zero reward. As before, an observable feature (gold star) indicates which door is the correct one.

187 Next, assume both agents select actions with a CBP (see Figure 2b)—Agent 1’s CBP conditions  
188 actions on two concepts: lever color and the door indicator; and Agent 2’s CBP estimates lever color  
189 alone. If we assume that intervening over any concept estimate reduces decision-making to a random  
190 policy, we can compute the expected rewards for each agent under intervention, as shown in Figure 2c.  
191 Computing Equation (2) over these expected rewards with  $K = 2$  returns a graph with a strong  
192 bi-directional edge between Agent 1’s door color concept and Agent 2’s lever color concept (see  
193 Figure 2d). This edge reveals an important inter-agent relationship: Agent 1’s door selection in Phase  
194 2 relies heavily on Agent 2’s lever selection in Phase 1. In contrast, Agent 1’s lever selection in Phase  
195 1 does not depend on Agent 2 (evident by the lone self-loop for Agent 1’s lever color intervention).

## 196 5 Experiments

197 In this section, we evaluate the extent to which our method addresses the following questions:

198 **Identifying Emergent Coordination** In cooperative settings, can concept intervention identify emergent  
199 coordination from policies that act independently, or how much coordination the environment

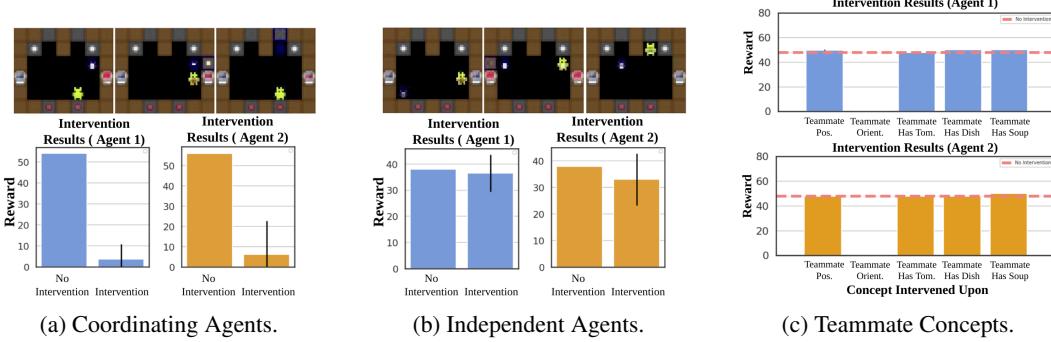


Figure 3: Performance degradation under concept intervention uncovers (a) policies that coordinate and (b) policies that act independently. (c) For coordinating agents, intervening over teammate-related concepts exposes the factors driving the team’s coordination strategy. Interestingly, teammate orientation is relied upon heavily by both agents for coordination.

demands? Can it expose lazy agents and, in general, measure an agent’s contribution to the multi-agent system? Moreover, if agents *are* coordinating, what specific features underlie that coordination?

**Emergent Strategic Behavior** In competitive environments, do CBPs reveal strategic behaviors as they emerge during multi-agent self-play training? For example, can intervention help us better understand both inter-team emergent roles and intra-team counter-strategies that develop over time? **Inter-Agent Social Dynamics** In mixed-incentive environments, how well does concept intervention identify inter-agent social dynamics? Can it expose the *functional connectivity* of the multi-agent system; or reveal chains of dependencies that might represent brittle inter-agent overfitting?

We use three environments from Melting Pot [21] for our experiments. Melting Pot environments vary along a number of multi-agent axes—*incentive alignment* (cooperative, mixed motive, adversarial), *coordination requirements*, etc—and therefore provide a strong benchmark for studying emergent behavior. Specific environments are described in the subsections below, with further details (including concept states) in Appendix B. For training, we use PPO [33], which has been shown to be state-of-the-art for multi-agent settings [9, 49], and train in a fully-decentralized fashion. We augment the PPO objective with our concept loss (called ConceptPPO moving forward). Additional training details and hyperparameter sweeps are provided in Appendix C.

## 5.1 Cooperative: Identifying Emergent Coordination

We first evaluate CBPs as a tool for understanding emergent coordination in fully-cooperative settings. We use the game Collaborative Cooking, where a group of agents inhabit a kitchen-like environment and must collaborate to find ingredients (tomatoes), complete recipes (bringing tomatoes to and from cooking pots), and deliver food as quickly as possible. We train 10 ConceptPPO policies with a concept cost weight of  $\lambda = 0.1$  in this environment and use them in our evaluation below.

**Do agents learn to coordinate?** We aim to distinguish policies that learn coordination from those that solve the task independently. We evaluate each of the trained ConceptPPO policies over 100 test-time trajectories (and 5 seeds each) while intervening over all of the concepts related to each agent’s teammate. The average cumulative reward under intervention is shown in Figure 3.

Figure 3 investigates a cooking environment with duplicate sets of ingredients (tomatoes) and tools (bowls, pots) on both sides. Therefore, both independent strategies and highly coordinated strategies can complete the task. In the trajectory in Figure 3a, we see an emergent strategy in which two agents coordinate. The orange agent works the bottom of the environment picking up tomatoes and bringing them to the cooking pot, while the blue agent stays in the top running dishes to and from the pot to deliver soup. Intervention over teammate-related concepts leads to a catastrophic drop in reward, as the agent’s coordination clearly hinges on an accurate modeling of its teammate. In Figure 3b, a different strategy emerges, where agents complete the task independently on opposite sides of the kitchen. Here, intervention does not hurt reward, as neither agent needs to closely monitor its teammate to complete the task. This result therefore indicates that coordination detection through intervention is accurate and reliable with CBPs. In Appendix C.3, we study several cooking

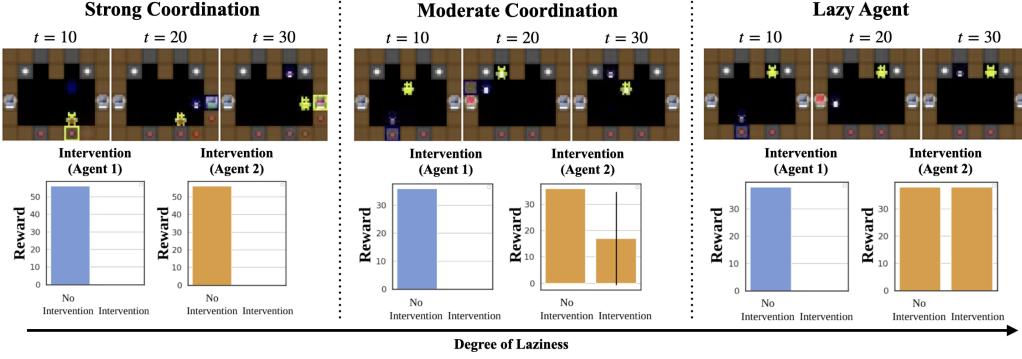


Figure 4: Intervening on an agent’s concepts pertaining to itself exposes its *degree of laziness*.

237 environment layouts and average the reduction in reward observed when intervening on teammate  
238 concepts across all policies. We can thus measure the coordination demands of each environment.

239 **How do they coordinate?** Next, we examine the policies that *have* learned to coordinate and pinpoint  
240 the specific inter-agent features that drive their coordination (Figure 3c). Rather than intervening over  
241 all teammate-related concepts at once, we intervene one at a time. If agents are coordinating using a  
242 particular signal, the corresponding concept intervention should result in a sharp drop in performance.

243 Interestingly, performance only drops when intervening over the orientation concept. This suggests  
244 that agents are primarily using the orientation of other agents as a coordination signal, which is  
245 curious, as we might expect coordination to involve multiple sources of inter-agent information. For  
246 completeness, we conduct two supporting analyses. First, we rule out the possibility that intervening  
247 with a fixed orientation creates an OOD (or adversarial) input. We do this using an empirical sample  
248 of agent orientations to manufacture interventions that are both in- and out-of-distribution; and show  
249 that our results are consistent in both cases (see Appendix C.4.1). Second, we re-run this intervention  
250 over additional ConceptPPO policies trained without orientation as a concept. These results show  
251 that agent coordination latches onto yet another concept (teammate has\_soup) as its primary signal  
252 (see Appendix C.4.2). The fact that agents learn a policy in which coordination hinges only on one  
253 particular signal reveals the brittleness of deep multi-agent reinforcement learning policies and their  
254 tendency to rely on highly specialized conventions. This brittleness is one of the reasons that MARL  
255 agents often fail to generalize to novel partners, or have poor zero-shot coordination performance [16].

256 **Identifying Lazy Agents** Our method also allows us to test for *coordination failures*. Here we test  
257 for lazy agents by replacing each agent’s concept estimates about itself, including its own position,  
258 orientation, etc. If an agent is acting productively in the environment, removing this information will  
259 greatly hinder its performance; and team performance as a whole. If performance does not degrade,  
260 the agent likely is not contributing to the task. Figure 4 shows the results of this test for three policies,  
261 each differing in the strength of their contribution to the team.

262 The magnitude of performance degradation is a direct function of the productivity of each agent.  
263 When agents are both contributing (left), concept intervention has a strong negative impact on the  
264 performance of the agents. In the extreme lazy agent case where one agent does not move over the  
265 course of a trajectory (right), reward is not impacted by intervention. Interestingly, this intervention  
266 also captures intermediate effects. For example, when only one agent is productive, but both agents  
267 share a workspace, there is a small but noticeable drop in performance (middle). Thus, these results  
268 not only demonstrate that CBPs give us a unique way of diagnosing lazy agents, but also shows that  
269 it can quantify the *degree of laziness* of each agent as a function of performance degradation.

## 270 5.2 Emergence of Strategic Behavior

271 There, we intervene over orientation (i.e. R1’s estimate of every other agent’s orientation individually)  
272 and flag position (i.e. R1’s estimate of the red flag and blue flag positions). These interventions give  
273 us insight into how important agent vs. environmental features are to the red agents’ policies. The  
274 key takeaway is that the pattern of intervention outcomes reveals the roles that both agents have learned. R1 is  
275 negatively impacted by interventions over its own orientation, the blue team’s attacking agent (B1),  
276 and its home flag. This pattern is consistent with the behavior of a base defender. R2 is negatively  
277

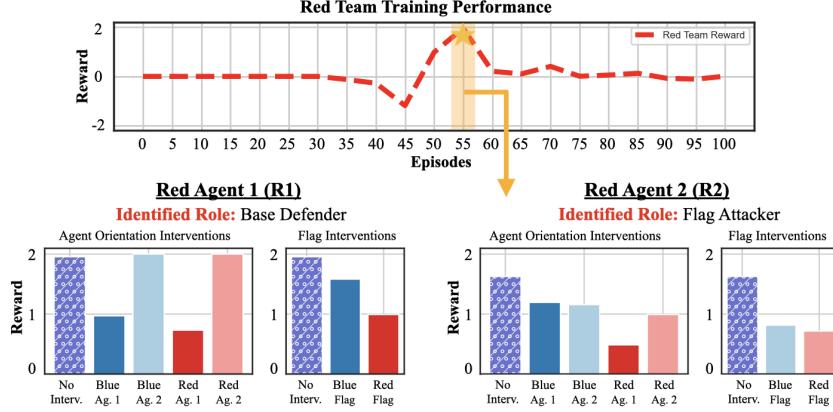


Figure 5: **Emergent strategy in Capture the Flag.** The red team discovers a dominant strategy at episode 55. For this episode, we intervene over orientation (e.g., R1’s estimate of every other agent’s orientation individually) and flag position (e.g., R1’s estimate of both flag positions). These interventions give us insight into how important agent vs. environmental features are to the red agents. Specifically, they reveal two emergent roles: (i) Red agent R1 is identified as a flag defender agent—it models B1 (the blue attacking agent) and its home flag and is only negatively impacted by interventions over those concepts; (ii) R2 interacts with both blue agents, both flags, and is strongly impacted by intervention. R2 is thus identified as a flag attacking agent.

278 impacted by interventions over both flags and each of the other agents. This means that R2 is a flag  
279 attacker, navigating between the opponent base to take flags and its home base to return them.

280 Mixed cooperative-competitive environments induce an automatic skill learning curriculum for  
281 MARL agents similar to that of self-play [36], which can lead to highly-sophisticated behaviors such  
282 as the emergent tool use observed in team hide-and-seek by Baker et al. [3]. Here we investigate the  
283 extent to which CBPs can expose these strategic behaviors during training. We train independent  
284 CBP agents in Melting Pot’s Capture the Flag environment (outlined further in Appendix B.4) and  
285 intermittently perform interventions over agent concept estimates during training. As before, we  
286 measure the magnitude of reward degradation.

287 Results from this analysis are shown in Figure 5. The reward curve represents one oscillation of self-  
288 play strategy for the red team—the red team receives negative reward in episode 45, but eventually  
289 counters with a dominant strategy in episode 55. Oscillations such as this are indicative of a newly  
290 learned strategic behavior by one or both agents [3]. But what strategic behavior was learned?

291 Conducting intervention analysis at episode 55 reveals more precisely how the agents behave beyond  
292 red agents capturing the flag. Analysis of the first red agent (R1) shows that it is negatively impacted  
293 by interventions over itself, the blue team’s attacking agent (B1), and its home flag. This pattern is  
294 consistent with the behavior of a *base defender*. Interventions over the second red agent (R2) expose  
295 a different pattern—the agent is negatively impacted by interventions over both flags and each of  
296 the other agents. This pattern suggests that R2 is a *flag attacker*, navigating between the opponent  
297 base to take flags and its home base to return them. A qualitative analysis of the agents’ behavior is  
298 performed in Appendix C.9 and confirms these results. Thus, our intervention technique can be used  
299 *during training* to augment reward-based analysis and investigate strategic behaviors *as they emerge*.

### 300 5.3 Social Dilemmas: Inter-agent Social Dynamics

301 Social dilemmas further extend the scope of emergent behavioral complexity to include exploitation,  
302 free-riding, and public resource sharing. We evaluate our method’s ability to uncover such interactions  
303 in Melting Pot’s Clean Up, an environment in which agents must balance selfish behavior (harvesting  
304 fruit) with public service (cleaning a river is necessary for fruit to grow).

305 As we’ve seen, a natural way to investigate inter-agent dynamics with CPBs is by looking at outcomes  
306 with and without intervention. In Figure 6, we do the same for each agent in Clean Up and find that  
307 the rewards of Agent 1 and 2 are correlated under intervention. We also see an increase in Agent 1  
308 and 2’s idleness and a decrease in their inter-agent distance. A reasonable conclusion, therefore, is

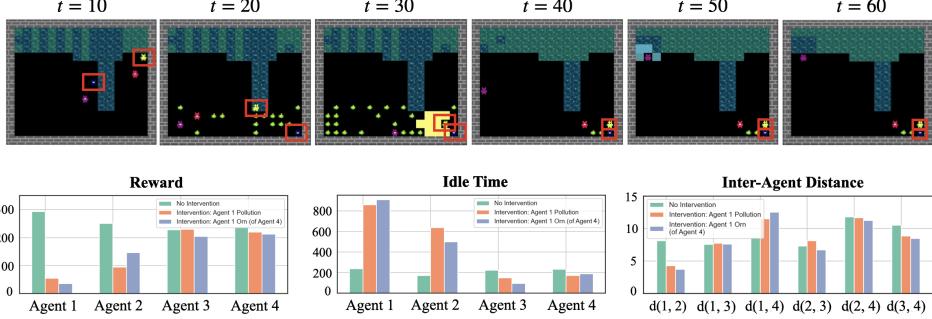


Figure 6: Learned social dynamics in Clean Up. Interventions with respect to reward and key behavioral features (agent idle time, inter-agent distance) an interesting relationship between Agents 1 (Blue) and 2 (Yellow). Both agents’ interventions cause their rewards to decrease, idle time to increase, and distance to decrease, suggesting that Agent 1 and 2 collide as a result of these interventions (the qualitative result at the top confirms this).

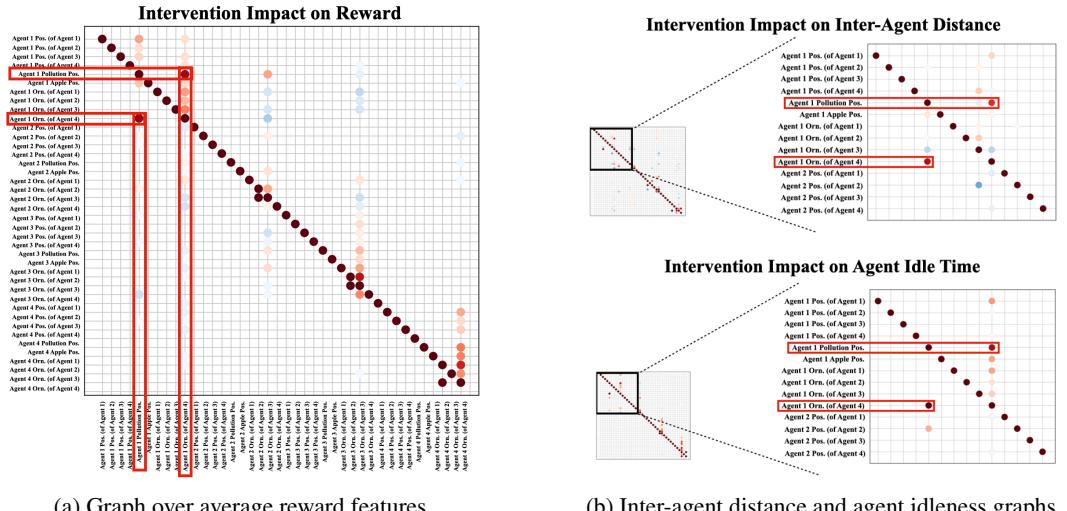


Figure 7: Social interactions in Clean Up. Each graph shows the result of running Equation (2) over all agent-concept pairs. (a) Intervening with respect to agent reward shows prominent bi-directional edges between Agent 1’s estimate of the closest pollution and Agent 1’s estimate of Agent 4’s orientation (red highlights). (b) Intervening with respect to behavioral features results in the same prominent edges. The correlation between Agent 1 and 2’s rewards, idle time, and distance (their collision in Figure 6) is therefore not the result of interrupting their coordination. Rather, it is the result of a chain of events that begin with a perturbation of agent 1’s perception. Thus, our graph learning strategy has uncovered a chain of events that links the emergent behaviors of each agent).

309 that Agent 1 and Agent 2 are coordinating spatially and require accurate estimations of each other  
310 to complete the task. However, as we’ll show, our method can uncover a more descriptive *chain of*  
311 *social interactions* that lead to these intervention outcomes. The trajectory in Figure 6 confirms this.

312 To investigate, we learn a series of graphs by computing Equation (2) over all agent-concept interventions  
313 using average reward, pairwise inter-agent distance, and the number of idle (non-moving) steps  
314 as features, respectively. Figure 7a shows the graph computed over rewards (in matrix form), which  
315 exposes a number of intriguing concept relationships across agents. First, there are *no* identifiable  
316 relationships between Agent 1 and 2 (no edges link interventions over Agent 1’s concepts pertaining  
317 to Agent 2 or vice versa). This means that Agent 1 and 2 are in fact not relying directly on each other’s  
318 concepts. Next, we find a strong bi-directional relationship between Agent 1’s closest pollution  
319 concept and Agent 1’s estimate of Agent 4’s orientation. Importantly, this same bi-directional edge  
320 also exists in the graphs computed over inter-agent distance and idleness (see Figure 7b).

321 The graphs therefore reveal the following: First, the lack of a link between the Agent 1 and 2 in the  
322 graph suggests that their relationship is not one of coordination, but rather the result of an incidental

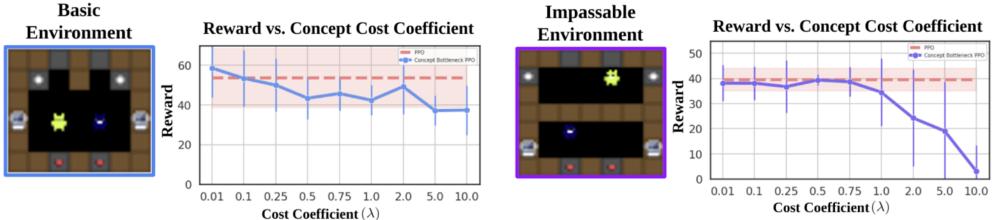


Figure 8: Asymptotic reward of ConceptPPO vs. PPO. For  $\lambda \leq 0.5$ , ConceptPPO matches the performance of non-concept-based PPO.

interaction in the environment (they are not explicitly modeling each other). Next, the bi-directional edge reveals that Agent 1’s performance is largely dependent on its exploitation of the public service of Agent 4. Specifically, Agent 4 runs to and from the river to clean, at which point apples appear in the area of the patch that Agent 4 left. Agent 1 uses both Agent 4’s orientation and its estimate of pollution (which Agent 4 changes) to decide when it should run to the area of the patch left by Agent 4 to consume apples. Because Agent 1’s policy is so reliant on these two concepts, intervening over either of them causes Agent 1’s behavior to collapse; and this happens in such a way that it interrupts Agent 2. Concretely, intervening over both Agent 1’s pollution concept and Agent 1’s estimate of Agent 4 lead to the same outcome: a collision between Agent 1 and Agent 2 (see Figure 7b).

In sum, the chain of dependencies is not  $\text{Agent 1} \leftrightarrow \text{Agent 2}$ , but rather  $\text{Agent 1} \rightarrow \text{Agent 4} / \text{Pollution} \rightarrow \text{Agent 2}$  and rather than coordinating, we have discovered that Agent 1 and Agent 2 have overfit to brittle behavioral patterns that are easily disrupted. This is in turn why there is a correlation between their rewards under intervention. Our graph learning technique over CBP interventions thus reveals complex inter-agent dynamics that could take humans many hours to detect.

#### 337 5.4 Bottleneck Performance

To evaluate our method more generally, we compare the asymptotic performance of CBPs relative to traditional, non-concept-based PPO. We emphasize that PPO has been shown to be a state-of-the-art algorithm for decentralized multi-agent learning [49]. We measure performance as the average cumulative reward obtained over 100 test-time trajectories (and five random seeds each). Figure 8 provides an overview of these results for two Collaborative Cooking environments. There we find evidence that CBPs can match the performance of PPO across each of our environments for small values of the concept cost coefficient ( $\lambda \leq 0.5$ ). This is an important result from the perspective of interpretability. It demonstrates that, if  $\lambda$  is tuned appropriately, it is possible to train intrinsically-interpretable policy networks—where, notably, decisions are expressed in human-understandable concepts—without sacrificing in task performance. Importantly, it also demonstrates the sufficiency of the concept set. Figure 8, and specifically the results for the Impassable Environment, show that over-valuing concept prediction loss causes performance to degrade. Performance falls for  $\lambda > 0.75$  and, in all but the basic environment, collapses to zero for larger values ( $\lambda \geq 5.0$ ). This breakdown occurs because, for high  $\lambda$ , the total gradient from  $L_C(\nabla L_{RL} + \lambda \nabla L_C)$  is dominated by the gradient from the concept-based loss  $L_C$ . In this case, gradient descent is not able to move the policy network’s parameters in the direction of  $\nabla L_{RL}$ , preventing policy optimization. We present our complete results across all environments in Appendix C.7.

## 355 6 Conclusion

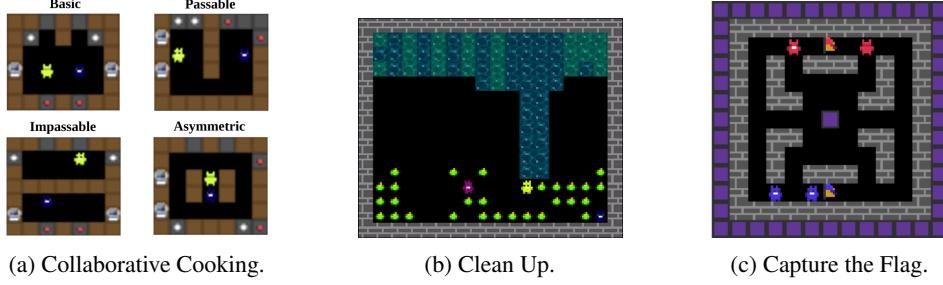
We introduced Concept Bottleneck Policies as an interpretable, concept-based policy learning method for MARL and demonstrated that they are effective for understanding emergent multi-agent behavior. In particular, CBPs support concept intervention, which can be used to identify when a multi-agent team has learned to coordinate, what inter-agent features drive that coordination, and to what extent coordination is required in an environment. Moreover, concept intervention helps expose coordination failures like lazy agents and complex inter-agent dynamics. We discuss broader impacts, limitations, and future work in Appendix A.

363 **References**

- 364 [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David,  
365 Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. Do as i can, not  
366 as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- 367 [2] Dan Amir and Ofra Amir. Highlights: Summarizing agent behavior to people. In *Proceedings*  
368 *of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages  
369 1168–1176, 2018.
- 370 [3] Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew,  
371 and Igor Mordatch. Emergent tool use from multi-agent autocurricula. *arXiv preprint*  
372 *arXiv:1909.07528*, 2019.
- 373 [4] Osbert Bastani, Carolyn Kim, and Hamsa Bastani. Interpretability via model extraction. *arXiv*  
374 *preprint arXiv:1706.09773*, 2017.
- 375 [5] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Agata Lapedriza, Bolei Zhou, and Antonio Torralba.  
376 Understanding the role of individual units in a deep neural network. *Proceedings of the National*  
377 *Academy of Sciences*, 117(48):30071–30078, 2020.
- 378 [6] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy  
379 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large  
380 scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- 381 [7] Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca  
382 Dragan. On the utility of learning about humans for human-ai coordination. *Advances in neural*  
383 *information processing systems*, 32, 2019.
- 384 [8] Zhi Chen, Yijie Bei, and Cynthia Rudin. Concept whitening for interpretable image recognition.  
385 *Nature Machine Intelligence*, 2(12):772–782, 2020.
- 386 [9] Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS  
387 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft  
388 multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- 389 [10] Xiaowen Dong, Dorina Thanou, Michael Rabbat, and Pascal Frossard. Learning graphs from  
390 data: A signal representation perspective. *IEEE Signal Processing Magazine*, 36(3):44–63,  
391 2019.
- 392 [11] Jessica Zosa Forde, Charles Lovering, George Konidaris, Ellie Pavlick, and Michael L Littman.  
393 Where, when & which concepts does alphazero learn? lessons from the game of hex. In *AAAI*  
394 *Workshop on Reinforcement Learning in Games*, volume 2, 2022.
- 395 [12] Ghost Town Games. Overcooked. <https://store.steampowered.com/app/448510/Overcooked/>,  
396 2016.
- 397 [13] Asma Ghandeharioun, Been Kim, Chun-Liang Li, Brendan Jou, Brian Eoff, and Rosalind W  
398 Picard. Dissect: Disentangled simultaneous explanations via concept traversals. *arXiv preprint*  
399 *arXiv:2105.15164*, 2021.
- 400 [14] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-  
401 based explanations. *Advances in Neural Information Processing Systems*, 32, 2019.
- 402 [15] Bradley Hayes and Julie A Shah. Improving robot controller transparency through autonomous  
403 policy explanation. In *2017 12th ACM/IEEE International Conference on Human-Robot*  
404 *Interaction (HRI)*, pages 303–312. IEEE, 2017.
- 405 [16] Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster. “other-play” for zero-shot  
406 coordination. In *International Conference on Machine Learning*, pages 4399–4410. PMLR,  
407 2020.

- 408 [17] Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia  
 409 Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al.  
 410 Human-level performance in 3d multiplayer games with population-based reinforcement learn-  
 411 ing. *Science*, 364(6443):859–865, 2019.
- 412 [18] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al.  
 413 Interpretability beyond feature attribution: Quantitative testing with concept activation vectors  
 414 (tcav). In *International conference on machine learning*, pages 2668–2677. PMLR, 2018.
- 415 [19] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been  
 416 Kim, and Percy Liang. Concept bottleneck models. In *International Conference on Machine  
 417 Learning*, pages 5338–5348. PMLR, 2020.
- 418 [20] Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. Multi-  
 419 agent reinforcement learning in sequential social dilemmas. *arXiv preprint arXiv:1702.03037*,  
 420 2017.
- 421 [21] Joel Z Leibo, Edgar A Dueñez-Guzman, Alexander Vezhnevets, John P Agapiou, Peter Sunehag,  
 422 Raphael Koster, Jayd Matyas, Charlie Beattie, Igor Mordatch, and Thore Graepel. Scalable  
 423 evaluation of multi-agent reinforcement learning with melting pot. In *International Conference  
 424 on Machine Learning*, pages 6187–6199. PMLR, 2021.
- 425 [22] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In  
 426 *Machine learning proceedings 1994*, pages 157–163. Elsevier, 1994.
- 427 [23] Siqi Liu, Guy Lever, Zhe Wang, Josh Merel, S. M. Ali Eslami, Daniel Hennes, Wojciech M.  
 428 Czarnecki, Yuval Tassa, Shayegan Omidshafiei, Abbas Abdolmaleki, Noah Y. Siegel, Leonard  
 429 Hasenclever, Luke Marris, Saran Tunyasuvunakool, H. Francis Song, Markus Wulfmeier, Paul  
 430 Muller, Tuomas Haarnoja, Brendan Tracey, Karl Tuyls, Thore Graepel, and Nicolas Heess.  
 431 From motor control to team play in simulated humanoid football. *Science Robotics*, 7(69):  
 432 eab0235, 2022. doi: 10.1126/scirobotics.ab0235.
- 433 [24] Siqi Liu, Guy Lever, Zhe Wang, Josh Merel, SM Ali Eslami, Daniel Hennes, Wojciech M  
 434 Czarnecki, Yuval Tassa, Shayegan Omidshafiei, Abbas Abdolmaleki, et al. From motor control  
 435 to team play in simulated humanoid football. *Science Robotics*, 7(69):eab0235, 2022.
- 436 [25] Prashan Madumal, Tim Miller, Liz Sonenberg, and Frank Vetere. Explainable reinforcement  
 437 learning through a causal lens. In *Proceedings of the AAAI conference on artificial intelligence*,  
 438 volume 34, pages 2493–2500, 2020.
- 439 [26] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent  
 440 variational exploration. *Advances in Neural Information Processing Systems*, 32, 2019.
- 441 [27] Thomas McGrath, Andrei Kapishnikov, Nenad Tomašev, Adam Pearce, Demis Hassabis, Been  
 442 Kim, Ulrich Paquet, and Vladimir Kramnik. Acquisition of chess knowledge in alphazero.  
 443 *arXiv preprint arXiv:2111.09259*, 2021.
- 444 [28] Nicolai Meinshausen and Peter Bühlmann. High-dimensional graphs and variable selection  
 445 with the lasso. *The annals of statistics*, 34(3):1436–1462, 2006.
- 446 [29] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
 447 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint  
 448 arXiv:1312.5602*, 2013.
- 449 [30] Shayegan Omidshafiei, Andrei Kapishnikov, Yannick Assogba, Lucas Dixon, and Been Kim.  
 450 Beyond rewards: a hierarchical perspective on offline multiagent behavioral analysis. *arXiv  
 451 preprint arXiv:2206.09046*, 2022.
- 452 [31] Judea Pearl. Causal inference in statistics: An overview. 2009.
- 453 [32] Julien Perolat, Bart de Vylder, Daniel Hennes, Eugene Tarassov, Florian Strub, Vincent de Boer,  
 454 Paul Muller, Jerome T Connor, Neil Burch, Thomas Anthony, et al. Mastering the game of  
 455 stratego with model-free multiagent reinforcement learning. *arXiv preprint arXiv:2206.15378*,  
 456 2022.

- 457 [33] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
 458 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 459 [34] Pedro Sequeira and Melinda Gervasio. Interestingness elements for explainable reinforcement  
 460 learning: Understanding agents' capabilities and limitations. *Artificial Intelligence*, 288:103367,  
 461 2020.
- 462 [35] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driess-  
 463 che, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mas-  
 464 tering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489,  
 465 2016.
- 466 [36] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur  
 467 Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. Mastering  
 468 chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint*  
 469 *arXiv:1712.01815*, 2017.
- 470 [37] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:  
 471 Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*,  
 472 2013.
- 473 [38] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smooth-  
 474 grad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- 475 [39] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving  
 476 for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- 477 [40] Wolfgang Stammer, Patrick Schramowski, and Kristian Kersting. Right for the right concept:  
 478 Revising neuro-symbolic concepts by interacting with their explanations. In *Proceedings of the*  
 479 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3619–3629, 2021.
- 480 [41] DJ Strouse, Kevin McKee, Matt Botvinick, Edward Hughes, and Richard Everett. Collaborating  
 481 with humans without human data. *Advances in Neural Information Processing Systems*, 34:  
 482 14502–14515, 2021.
- 483 [42] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
 484 *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
- 485 [43] Eugene Vinitsky, Nathan Lichtenlé, Kanaad Parvate, and Alexandre Bayen. Optimizing mixed  
 486 autonomy traffic flow with decentralized autonomous vehicles and multi-agent reinforcement  
 487 learning. *ACM Transactions on Cyber-Physical Systems*, 7(2):1–22, 2023.
- 488 [44] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Jun-  
 489 young Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster  
 490 level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 491 [45] Cathy Wu, Abdul Rahman Kreidieh, Kanaad Parvate, Eugene Vinitsky, and Alexandre M Bayen.  
 492 Flow: A modular learning framework for mixed autonomy traffic. *IEEE Transactions on*  
 493 *Robotics*, 38(2):1270–1286, 2021.
- 494 [46] Sarah A Wu, Rose E Wang, James A Evans, Joshua B Tenenbaum, David C Parkes, and  
 495 Max Kleiman-Weiner. Too many cooks: Bayesian inference for coordinating multi-agent  
 496 collaboration. *Topics in Cognitive Science*, 13(2):414–432, 2021.
- 497 [47] Zhao Yang, Song Bai, Li Zhang, and Philip HS Torr. Learn to interpret atari agents. *arXiv*  
 498 *preprint arXiv:1812.11276*, 2018.
- 499 [48] Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar.  
 500 On completeness-aware concept-based explanations in deep neural networks. *Advances in*  
 501 *Neural Information Processing Systems*, 33:20554–20565, 2020.
- 502 [49] Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising  
 503 effectiveness of ppo in cooperative, multi-agent games. *arXiv preprint arXiv:2103.01955*, 2021.



(b) Clean Up.

(c) Capture the Flag.

(a) Collaborative Cooking.

(b) Clean Up.

(c) Capture the Flag.

(a) Collaborative Cooking.

(b) Clean Up.

(c) Capture the Flag.

## 504 A Discussion

505 **Limitations** One of the limitations of the proposed method is that it requires manual specification of  
 506 human-interpretable concepts to be used in the policy bottlenecks at training time. As illustrated in the  
 507 experiments section of the paper, these are often intuitive characteristics of the specific environment  
 508 targeted (e.g., positions and orientations of agents, states of various objects that can be interacted  
 509 with, etc.), which can be then automatically extracted from the environment simulator. In future work,  
 510 it would be interesting to learn such key relevant concepts automatically. Another potential limitation  
 511 of the method is that the combination of traditional reward-based RL loss and concept-loss can make  
 512 for a slightly complex loss landscape. However, as illustrated by our experiments, one mechanism for  
 513 controlling this is via the concept loss  $\lambda$  parameter. To make the selection of this parameter easier, a  
 514 baseline of the agent policies ran *without* the concept bottlenecks can be used to establish reasonable  
 515 bounds on their expected performance in the task. Subsequently, the parameter can be swept over  
 516 (similar to our experiments) to identify a  $\lambda$  that leads to accurate concept learning with minimal  
 517 impact on final performance. Finally, we note that we do not provide theoretical guarantees that CBPs  
 518 converge to a specific solution concept (i.e. refinement of the Nash equilibrium). In general, such  
 519 guarantees are difficult to provide when combining MARL with neural networks and the addition of  
 520 a concept-based loss further complicates this picture; though we do point out that the lack of such  
 521 guarantees in many prior MARL works highlights the importance of the interpretability problem.

522 **Broader Impacts** In terms of potential societal implications, modeling of interactive decision-making  
 523 and concept estimations of agents in such systems might be used by adversarial actors to intervene  
 524 on concepts and intentionally cause performance degradation or agents misbehaving. On the other  
 525 hand, concept bottlenecks themselves could be a means of identifying such adversarial attacks, by  
 526 detecting shifts in concept estimations that were otherwise unexpected.

527 **Future Work** There are a number of interesting avenues for future work in the area of concept  
 528 bottlenecks for MARL. First, taking inspiration from the literature on POMDPs and belief-space  
 529 planning, we can explore extensions to our architecture where agents can leverage both histories  
 530 of observations and histories of concepts in training and concept estimation. This may also be  
 531 combined with architecture improvements, such as giving agents memory (e.g., through recurrent  
 532 policies). Further, there is room to explore the use of concept prediction accuracy as an intrinsic  
 533 reward that incentivizes agents to explore subsets of the state space in which they do not predict  
 534 concepts accurately (thereby improving their concept estimates). An important future work is to  
 535 extend our simple graph-based analysis technique to more complex graph learning paradigms [10].

## 536 B Additional Environment Details

### 537 B.1 Setup

538 The following details are common across all of our environments:

- 539 • *States*: Each environment is a grid world. Grid cells can be filled by an agent or an  
 540 environment-specific object.

- 541     • *Observations*: Agents receive partial multi-modal observations consisting of their own  
 542       position and orientation in the grid, as well as a partial RGB rendering of a 5-cell × 5-cell  
 543       window centered at the agent.  
 544     • *Actions*: Agents can execute one of 8 actions: no-op, move {up, down, left, right}, turn  
 545       {left, right}, and interact.

546     **B.2 Collaborative Cooking**

547     A visualization of the four cooking environments from our experiments are shown in Figure 9a.  
 548     Based on the game Overcooked [12] and subsequent work that has developed Overcooked-like  
 549       environments for studying multi-agent coordination [7, 46], Melting Pot’s Collaborative Cooking is  
 550       a game in which a group of agents inhabit a kitchen-like environment and must collaborate to find  
 551       ingredients, complete recipes, and deliver finished dishes as quickly as possible. Solving the cooking  
 552       task requires sophisticated coordination, involving both task partitioning—splitting a recipe into  
 553       parts—and role assignment—distributing sub-tasks among agents. For these reasons, Collaborative  
 554       Cooking is investigated in a number of prior works [46, 7, 41] and is emerging as a strong benchmark  
 555       for multi-agent learning. The Collaborative Cooking game for  $N$  agents is defined as follows:

- 556     • *Recipe*: (i) Bring a tomato to a cooking pot (3 times), (ii) Wait for soup to cook in the pot  
 557       (20 time-steps), (iii) Bring a dish to the cooking pot, (iv) Pour soup from the pot into the  
 558       dish, (v) Deliver soup to the delivery location. In practice, solving this task from scratch  
 559       in its entirety is extremely difficult, as each of the aforementioned steps requires agents to  
 560       execute a series of movement and interaction actions in sequence.  
 561     • *States*: Grid cells can be filled by an agent or any of the following items: Floor, Counter,  
 562       Cooking Pot, Dish, Tomato.  
 563     • *Reward*: By default, agents share a positive reward for completing the entire recipe outlined  
 564       above. In practice, however, solving the cooking task with this sparse reward alone is  
 565       infeasible (completing each of the recipe steps through random exploration is prohibitively  
 566       challenging) and successful approaches in prior works either pair learning agents with  
 567       helpful bot agents or introduce “densified” pseudorewards to augment the agents’ learning  
 568       signal [20]. We implement the latter, giving agents a small positive reward for completing  
 569       steps (i) and (iii) of the recipe. More concretely, we define the following three-part reward:

$$r_t = \begin{cases} 20, & \text{if soup cooked and delivered.} \\ 1, & \text{if tomato placed in cooking pot.} \\ 1, & \text{if soup poured into dish.} \\ 0, & \text{otherwise.} \end{cases}$$

- 570     • *Concepts*: We assume that each environment supports the following concepts (and concept  
 571       types): (i) agent position (scalar); (ii) agent orientation (scalar); (iii) whether or not an agent  
 572       has a tomato, dish or soup (binary); (iv) cooking pot position (scalar) (v) the progress of the  
 573       cooking pot (scalar); (vi) the number of tomatoes in the cooking pot (categorical); and (vii)  
 574       the position of each tomato and dish (scalar).

575     **B.3 Clean Up**

576     The Clean Up game is shown in Figure 9b. Clean Up is a classic social dilemma in which agents  
 577       are forced to balance selfish behaviors with public goods. Each agent is rewarded proportionally to  
 578       the number of apples it harvests. However, the rate at which apples respawn after they are eaten is  
 579       directly related to the amount of pollution that exists in the river. Over time, pollution builds up in the  
 580       river if it remains uncleared until eventually, apples are prevented from growing altogether. Agents  
 581       must learn to periodically clean the river to ensure that apples continuously grow, even though they  
 582       are not directly incentivized to do so by default. Therefore, it is possible for agents to develop both  
 583       emergent “free-loading” behaviors—harvesting apples without cleaning the river—and emergent  
 584       public service behaviors—cleaning the river while trying to harvest as many of the remaining apples  
 585       after cleaning. This makes Clean Up a fitting environment with which to study the extent to which  
 586       inter-agent social dynamics can be revealed. Some environment-specific details are outlined below:

- 587 • *States*: Grid cells can be filled by an agent or any of the following items: Floor, Wall, Apple,  
 588 Clean River, Polluted River.
- 589 • *Reward*: Agents receive an individual positive reward for collecting apples. By default,  
 590 agents are not rewarded for cleaning the river. We found that this greatly increased the  
 591 amount of training time required for agents to learn to clean the river (because they have no  
 592 incentive to do so). Because we are motivated by understanding social behaviors once they  
 593 have emerged, we augmented each agent’s reward with a small positive reward for cleaning  
 594 the river. Specifically, we define the following reward for the cleaning task:

$$r_t = \begin{cases} 1, & \text{for eating an apple.} \\ 0.01, & \text{for cleaning pollution.} \\ 0, & \text{otherwise.} \end{cases}$$

- 595 • *Concepts*: We assume that each environment supports the following concepts (and concept  
 596 types): (i) agent position (scalar); (ii) agent orientation (scalar); (iii) closest pollution  
 597 position (scalar); (iv) closest apple position (scalar).

#### 598 B.4 Capture the Flag

599 The Capture the Flag game is shown in Figure 9c. Capture the Flag is a well-known competitive  
 600 environment in the multi-agent literature [17]. The game pits two multi-agent teams (red, blue)  
 601 head-to-head in an arena-style environment. The red team is spawned at the top of the environment,  
 602 and the blue team is spawned at the bottom. Both team have at their “home base” a flag of their team’s  
 603 color. The goal of each team is to capture the opposing team’s flag as many times as possible within  
 604 the game’s time horizon; while also preventing the opposing team from capturing the team’s home  
 605 flag. The primary defense mechanism for an agent is a “zap” action that reduces an opponent agent’s  
 606 health if it is in the zap radius, and forces the agent to drop the flag (if it is holding one). In addition  
 607 to zapping other agents, agents can choose a “paint” action that colors the floor tiles with the painting  
 608 agent’s home color. Agents of the opposite color cannot move over painted tiles until they paint it  
 609 their own color, so the painting mechanism can be used strategically to slow down the movement of  
 610 opponents through certain corridors. There is also a set of purple tiles in the environment (one in  
 611 the center, many surrounding the arena) that serves as an indicator of whether or not a flag has been  
 612 taken from its home base. This tile serves as an observation for each agent, informing that agent of  
 613 the current state of both team’s flags (purple = no flags taken, blue = blue flag taken, red = red flag  
 614 taken, gray = both flags taken).

615 Sophisticated Capture the Flag agents can learn a number of complex strategies, such as attacking an  
 616 opponent’s flag, battling opponent agents for territory, and base camping to protect the home flag.  
 617 Agents are forced to learn these competitive behaviors from scratch and refine those behaviors during  
 618 training. From the perspective of Concept Bottleneck Policies, therefore, Capture the Flag is a strong  
 619 case study of the emergence of strategic behaviors over time. Some environment-specific details are  
 620 outlined below:

- 621 • *States*: Grid cells can be filled by an agent or any of the following items: Wall, Red Flag,  
 622 Blue Flag, Red Paint, Blue Paint, Flag Indicator Tile.
- 623 • *Reward*: Agents receive a shared team reward for capturing the opposing team’s flag (by  
 624 running it back to base). Agents are also rewarded individually for picking up a flag from the  
 625 opposing team’s base, returning a previously stolen flag to base, and zapping an opponent  
 626 flag carrier. Agents are therefore most strongly motivated to act offensively, but must learn  
 627 to balance attacks with defensive strategy. In sum, we define the following reward for the  
 628 task:

$$r_t = \begin{cases} 1, & \text{for capturing the opponent’s flag..} \\ 0.01, & \text{for picking up the opponent’s flag.} \\ 0.01, & \text{for returning a flag to base.} \\ 0.01, & \text{for zapping an opponent flag carrier.} \\ 0, & \text{otherwise.} \end{cases}$$

- 629 • *Concepts*: We assume that each environment supports the following concepts (and concept  
 630 types): (i) agent position (scalar); (ii) agent orientation (scalar); (iii) flag position (scalar);

Table 1: Hyperparameter sweeps for training ConceptPPO and PPO. Swept values are shown in braces and highest-performing values are bolded.

Hyperparameters	
Name	Value
Training Steps	25e6
Batch Size	{64, 128, 256, <b>512</b> , 1024}
Learning Rate	{1e-3, <b>1e-4</b> , 1e-5}
Gradient Norm	{0.1, <b>0.5</b> , 1.0, 5.0, 10.}
PPO Unroll Length	{4, 8, <b>16</b> , 32}
PPO Clipping $\epsilon$	{0.01, <b>0.05</b> , 0.1, 0.2, 0.3}
PPO Entropy Cost	{0.001, <b>0.01</b> , 0.05}
PPO Value Cost	{0.75, 0.9, <b>1.0</b> }

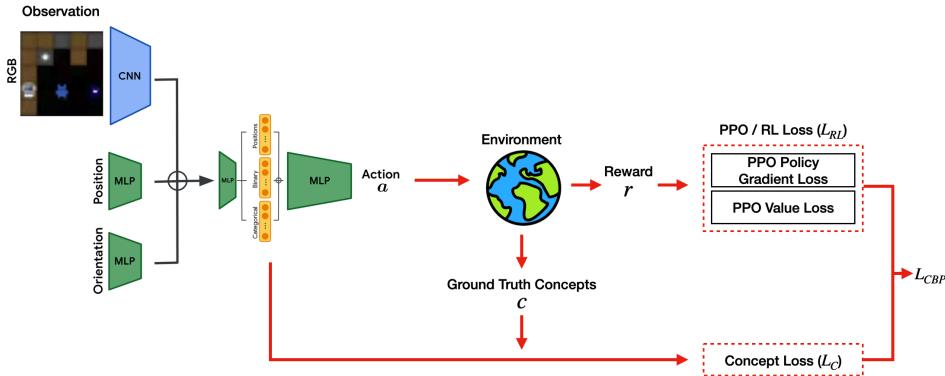


Figure 10: An architecture diagram showing how the concept bottleneck loss  $L_{CBP}$  is computed in  $L_C$ . Crucially, PPO is applied only over  $L_{RL}$ , which captures environmental rewards (e.g., delivering in collaborative cooking). Separately, concept loss  $L_C$  is computed in a supervised fashion using concept labels that are extracted from the environment. The concept loss is **not** incorporated into an agent’s reward function and therefore does not incentivize its behavior directly.

631 (iv) whether or not an agent has the opponent’s flag (binary) (scalar); (v) floor paint color  
 632 (categorical); (vi) flag indicator tile color (categorical).

## 633 C Training Details

### 634 C.1 Architecture and Hyperparameters

635 For both PPO and ConceptPPO, each agent’s policy network consists of CNN and MLP encoders  
 636 (for image and position/orientation inputs, respectively), followed by a two-layer MLP and a linear  
 637 mapping that compresses the encoded inputs into concept predictions. Concept estimates are fed  
 638 through a two-layer MLP, which produces the final action. ReLU activation is used throughout  
 639 (except in the bottleneck layer itself). As a baseline, we use vanilla PPO (no concept loss) with  
 640 the same architecture. We train 10 individual policies across each of the following values of  $\lambda$ :  
 641  $\{0.01, 0.1, 0.25, 0.5, 0.75, 1.0, 2.0, 5.0, 10.0\}$ . We conducted a wide hyperparameter sweep to train  
 642 both ConceptPPO and PPO, which is summarized in Table 1.

### 643 C.2 Loss Breakdown

644 In this section, we present a more detailed discussion of the concept bottleneck loss introduced  
 645 in Section 4 and, specifically,  $L_C$ . The objective  $L_{CBP}$  is a function of two separately computed  
 646 objectives: (i)  $L_{RL}$ , which is a reward-based loss that operates only over environmental rewards (e.g.,  
 647 delivering in collaborative cooking); and (ii)  $L_C$ , which is a supervised loss computed over concept  
 648 labels that are extracted from the environment. We are leveraging PPO to optimize each agent’s

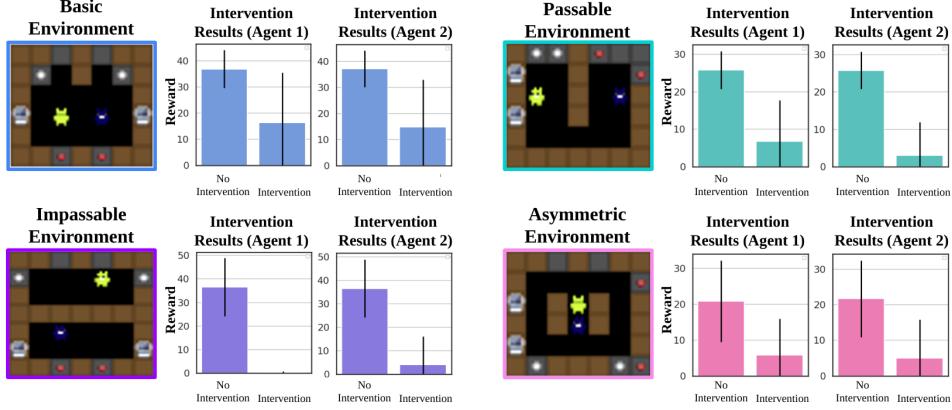


Figure 11: **Does the environment require coordination?** Averaging the impact of concept intervention over all policies trained in an environment reveals the extent to which coordination is required by that environment. In the impassable environment, agents cannot solve the task without coordinating, leading to consistent performance drops under intervention. In the basic environment, policies that coordinate (and therefore fail under intervention) are averaged with policies that act independently (and are uninterrupted by intervention), so the overall impact of intervention is less severe.

649 policy, which includes terms for both value error and generalized advantage estimation. A crucial  
 650 detail of our architecture is that PPO’s value and advantage estimation operates over environmental  
 651 rewards only. The accuracy of concept predictions is not incorporated into an agent’s reward in any  
 652 way—i.e. we are not adding an auxiliary/intrinsic rewards associated with the concept predictions  
 653 to the advantage function estimate. This is done intentionally so as not to bias agent behavior away  
 654 from states that result in high concept error prediction.

### 655 C.3 Environmental Demands of Coordination

656 For each Concept PPO and PPO policy trained in our cooking environments, we mask out, for each  
 657 agent, all of the concepts related to that agent’s teammate. We measure the average cumulative  
 658 reward attained over 100 trajectories each. The results of this intervention test are shown in Figure 11.  
 659 The level of coordination required by the environment is apparent from the severity of performance  
 660 degradation that results from intervening on each agent’s estimates of its teammate. Most notable  
 661 is the contrast between the basic environment (blue) and the impassable environment (purple). The  
 662 impassable environment requires strict coordination—agents can only access a subset of ingredients  
 663 and must pass items to each other across the center divider. In this case, intervention performance  
 664 drops to near-zero across all policies. In the basic environment, on the other hand, there are no  
 665 obstacles and agents have access to their own supply of ingredients; and so both policies in which  
 666 agents coordinate and policies in which agents act independently are successful. Consequently, the  
 667 impact of intervention is much less severe, as the performance of policies that coordinate (and fail  
 668 under intervention) is averaged in with independent policies that are unaffected by intervention.  
 669 Altogether, these results indicate that our method accurately distinguishes environments that require  
 670 coordination from those that do not.

671 We also present results for this analysis in two additional cooking environments: passable and  
 672 asymmetric. As in the basic and impassable environments, we find that the coordination demands  
 673 of the asymmetric and passable environments can be revealed through concept intervention. Both  
 674 environments involve moderate coordination—the most efficient strategy for bringing items to and  
 675 from the cooking pot involves passing them over the counter from one agent to another—but this  
 676 coordination is not required (like in the impassable environment). For this reason, we still see a  
 677 significant drop in the performance of our agents under intervention, but not a full decrease to zero.

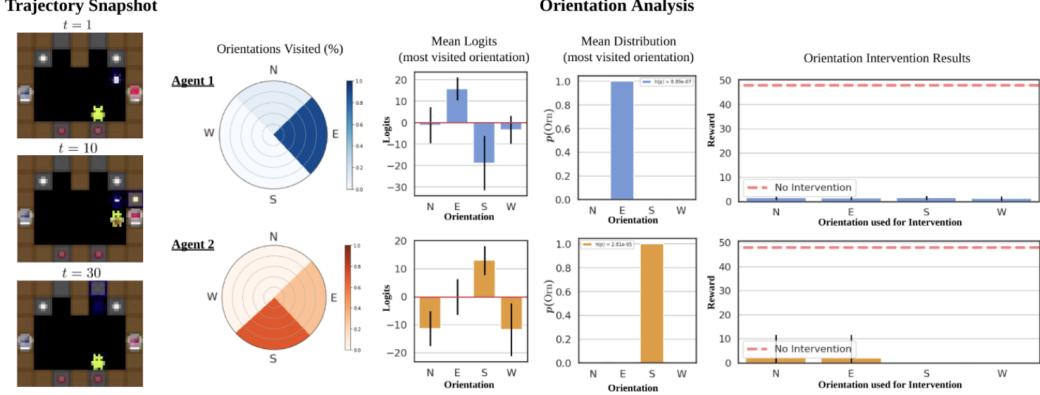


Figure 12: Overview of our method for constructing in-distribution concept masks for intervention. Using test-time trajectories, we compute an empirical distribution of the orientations experienced by each agent. Next, we compute the mean logits (and corresponding distribution) produced by each agent for the most frequently visited orientation of its teammate. Finally, we report the results of intervening with this mean logits vector in place of each of the four cardinal orientations. The results of this intervention are consistent, regardless of whether the mask used was in- or out-of-distribution.

#### 678 C.4 Factors of Coordination Analysis (cont'd)

##### 679 C.4.1 In vs. Out-of-Distribution Orientation Analysis

680 Here we further examine the impact of intervening on orientation. First, we compute an empirical  
 681 distribution of orientations that each agent visits over 100 test-time trajectories (across five random  
 682 seeds each). From those trajectories, we compute both the mean logits vector produced by each  
 683 agent’s concept bottleneck for the *most frequently visited orientation*, and the distribution represented  
 684 by those logits. Crucially, the mean logits vector for the most frequently visited orientation can be  
 685 used as an in-distribution value for concept interventions—it is an orientation estimate that each  
 686 agent has likely seen before. Similarly, we can create an out-of-distribution concept intervention  
 687 mask by permuting the mean logits vector such that the probability mass shifts a different cardinal  
 688 orientation. Using these manufactured orientations, we perform the same intervention technique as  
 689 before, iteratively replacing each agent’s orientation concept with the mean logits vector in place of  
 690 each cardinal direction, and measure performance of the multi-agent team.

691 To provide intuition for this technique, we ground it in the cooking task used for our experiments.  
 692 Consider, for example, the tomato-picking agent in the trajectory snapshot of Figure 12, who primarily  
 693 faces south and east—the directions needed to pick tomatoes from the bottom counter and place them  
 694 in the cooking pot on the right-hand side. The tomato-picking agent’s teammate (the waiter agent)  
 695 must learn to accurately model these orientations to satisfy its concept prediction objective, and so  
 696 frequently passes low-entropy distributions for south and east to its policy network. Intervening on  
 697 the waiter agent’s orientation estimate with a low-entropy distribution for south and east, therefore,  
 698 creates an in-distribution mask, whereas intervening with a low-entropy distribution for north or west  
 699 creates an out-of-distribution mask.

700 The results of this intervention test are shown in Figure 12, alongside the orientation distributions,  
 701 mean logits, and the distribution represented by those logits. Interestingly, the previously observed  
 702 degradation of performance as a result of intervening on teammate orientation is upheld, regardless  
 703 of whether the mask value is in- or out-of-distribution.

704 This provides further evidence that the agent’s reliance on orientation is a legitimate artifact of their  
 705 emergent strategy and not an adversarial or OOD example.

##### 706 C.4.2 Intervention Without Orientation

707 To further investigate the surprising use of orientation as the primary signal driving the emergent  
 708 behavior of our learning agents, we train a set of ConceptPPO policies in our basic environment  
 709 without orientation as a concept (across each  $\lambda$  value as before). We then perform the same iterative

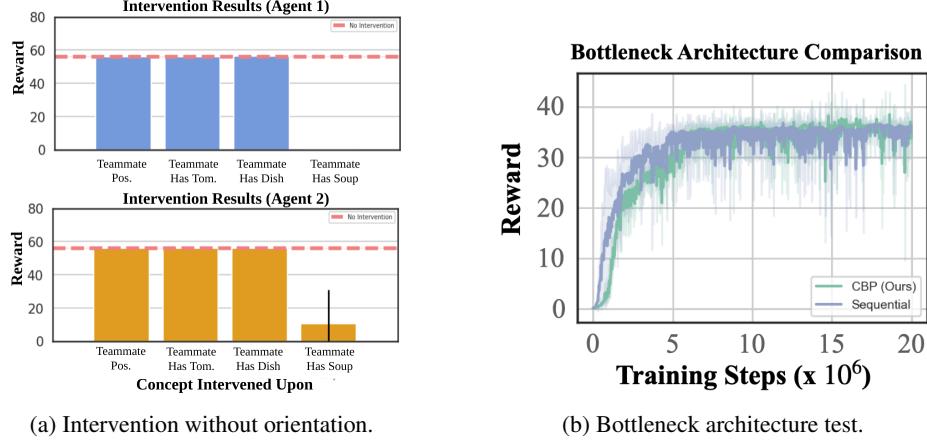


Figure 13: Additional bottleneck analyses. (a) Results of iterative concept intervention on Concept PPO agents trained without orientation. (b) Asymptotic performance of joint CBP architecture (ours) vs. sequential CBP architecture.

710 intervention analysis over the concepts pertaining to each agent’s teammate (as outlined in Section 5.1  
 711 (excluding orientation, of course). The results of this analysis are shown in Figure 13a. After  
 712 removing each agent’s orientation concept, we find that the agents latch onto a new concept—“has  
 713 soup”—as the primary concept that drives their coordination.

### 714 C.5 Bottleneck Architecture Comparison

715 Consider a concept bottleneck composed of two functions: (i)  $g : x \rightarrow c$  mapping inputs to concept  
 716 estimates; and (ii)  $f : c \rightarrow y$  mapping concept estimates to outputs. The work of Koh et al. [19]  
 717 compares three architectures for concept bottlenecks in supervised learning settings that lean these  
 718 functions in different ways:

- 719 • **Independent Bottleneck:** This architecture learns functions  $\hat{f}$  and  $\hat{g}$  independently, where  
 720  $\hat{g}$  is trained separately to predict concept estimates  $\hat{c}$  from inputs  $x$  using ground truth  
 721 concepts  $c$  as supervision.  $\hat{f}$  is trained to predict outputs  $\hat{y}$  from ground truth concepts  $c$ ,  
 722 using ground truth labels  $y$  as supervision.
- 723 • **Sequential Bottleneck:** This architecture learns  $\hat{g}$  first, then uses concept predictions from  
 724  $\hat{g}$  to train  $\hat{f}$ . Specifically,  $\hat{g}$  is trained in the same manner as in the independent bottleneck.  
 725  $\hat{f}$  is then trained to predict outputs  $\hat{y}$  from concept estimates  $\hat{c} = \hat{g}(x)$ , using ground truth  
 726 labels  $y$  as supervision.
- 727 • **Joint Bottleneck:** This architecture trains  $\hat{f}$  and  $\hat{g}$  together. The models are trained jointly  
 728 to predict  $\hat{y} = \hat{f}(\hat{g}(x))$  using an objective that is a weighted sum of loss terms for concept  
 729 prediction (using ground truth concepts  $c$  as supervision) and target label prediction (using  
 730 ground truth labels  $y$  as supervision).

731 Our proposed CBP architecture is most like the joint bottleneck. Here we compare the performance of  
 732 our proposed architecture to a MARL equivalent of the sequential bottleneck. A simple way to obtain  
 733 an approximation of the sequential bottleneck using the same architecture (shown in Figure 10) and  
 734 objective (defined by  $L_C$ ) as CBPs is by the stopping gradient flow between the concept estimator  
 735 network and policy head during backpropagation.

736 We train both CBPs (our proposed joint bottleneck architecture) and sequential bottlenecks in the  
 737 basic cooking environment shown in Figure 9a. Comparing the asymptotic reward achieved by both  
 738 approaches shows that joint CBPs (our architecture) and the sequential bottleneck (implemented via  
 739 stop gradient) perform comparably. This gives us confidence that the joint architecture is not unfairly  
 740 surpassing the performance of other bottleneck architectures by “hacking” concept estimates. We  
 741 investigate concept leakage more concretely in the following subsection.

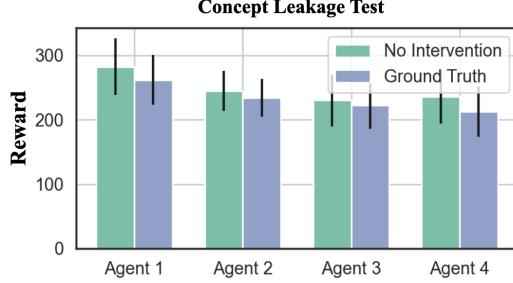


Figure 14: Concept leakage test. We perform an intervention over each agent’s concept estimates with the ground truth concept labels from the environment. There is a very slight decrease in performance during intervention, which indicates that some concept leakage does exist. However, the small amount of performance degradation indicates that agents are faithfully learning to estimate concepts.

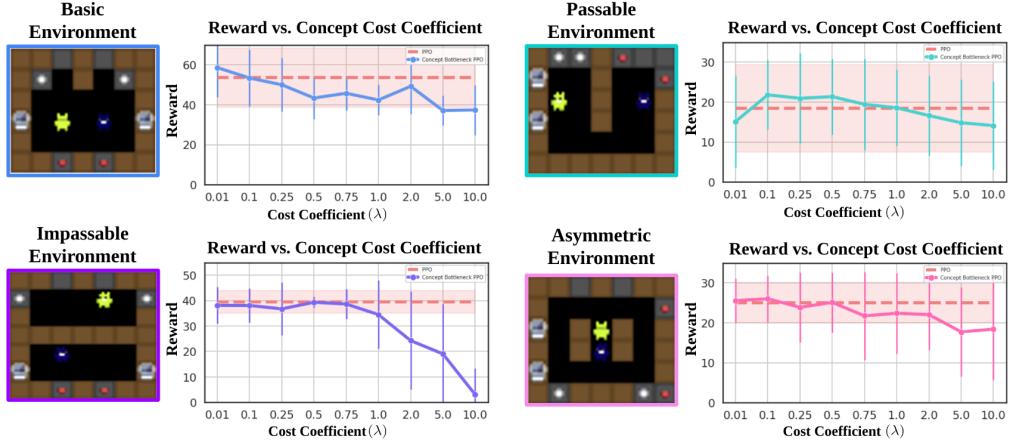


Figure 15: Asymptotic reward of ConceptPPO vs. PPO. For  $\lambda \leq 0.5$ , ConceptPPO matches the performance of non-concept-based PPO.

## 742 C.6 Concept Leakage Analysis

743 To examine the extent to which concept leakage occurs as a result of the joint objective outlined in  
 744  $L_C$ , we perform an evaluation of CBPs while intervening over each agent’s concept estimates with  
 745 ground truth concept values (for all concepts). Specifically, at each time-step, we replace each agent’s  
 746 concept estimates with the ground truth concept values extracted from the environment. Intuitively, if  
 747 concept leakage is a rampant issue in our proposed architecture and agents are learning to "hack"  
 748 their concept estimates to encode additional side channel information, we should see performance  
 749 degrade significantly as a result of this intervention. We perform this analysis over 100 test-time  
 750 trajectories (across 5 random seeds each) using the trained CBP policies from the Clean Up analysis  
 751 in Section 5.3.

752 The results of this analysis are shown in Figure 14, which compares the average reward under the  
 753 aforementioned ground truth intervention analysis to the average reward obtained by agents with no  
 754 interventions (using concept estimates as usual). Interestingly, there is a slight decrease in performance  
 755 that occurs as a result of the ground truth intervention, indicating that some concept leakage may be  
 756 present in the agents’ CBPs. However, the small magnitude of performance degradation indicates  
 757 that agents are reliably encoding the true concept values in their bottleneck predictions.

## 758 C.7 Bottleneck Performance

759 To evaluate our method more generally, we compare the asymptotic performance of Concept PPO  
 760 and PPO. We measure performance as the average cumulative reward obtained over 100 test-time  
 761 trajectories (and five random seeds each). Figure 15 provides an overview of these results. We find

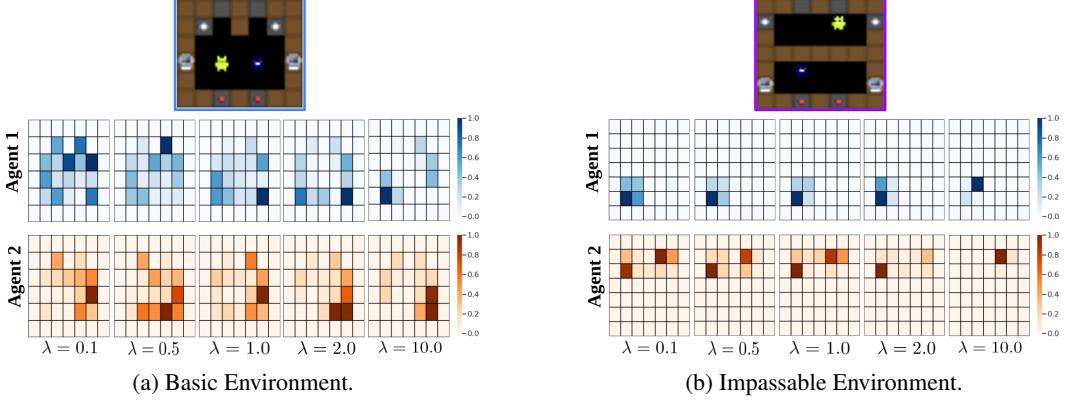


Figure 16: State visitation results. As concept cost coefficient  $\lambda$  increases, behavior collapses.

evidence that ConceptPPO can match the performance of PPO across each of our environments for small values of the concept cost coefficient ( $\lambda \leq 0.5$ ). This is an important result from the perspective of interpretability. It demonstrates that, if  $\lambda$  is tuned appropriately, it is possible to train intrinsically-interpretable policy networks—where, notably, decisions are expressed in human-understandable concepts—without sacrificing in task performance. Importantly, it also demonstrates the sufficiency of the concept set.

Figure 15 also shows that over-valuing concept prediction loss causes performance to degrade. Performance falls for  $\lambda > 0.75$  and, in all but the basic environment, collapses to zero for larger values ( $\lambda \geq 5.0$ ). This breakdown occurs because, for high  $\lambda$ , the total gradient from  $L_C$  ( $\nabla L_{RL} + \lambda \nabla L_C$ ) is dominated by the gradient from the concept-based loss  $L_C$ . In this case, gradient descent is not able to move the policy network’s parameters in the direction of  $\nabla L_{RL}$ , preventing policy optimization.

### 773 C.8 State Visitation Analysis

To illustrate the behavioral changes induced by  $\lambda$ , we measure the distribution of states visited by each agent across a subset of the  $\lambda$  values used during training. The results for each agent are plotted as a heatmap in Figure 16. As  $\lambda$  increases, multi-agent behavior begins to break down, as the gradient signal from environmental reward gets over-shadowed by that of the concept prediction objective. In the most extreme cases (e.g.,  $\lambda = 10.0$  in the impassable environment) agents fail to make any progress in completing the task.

### 780 C.9 Role Assignment Results (Cont’d)

**781 Full Temporal Analysis** Here we supplement our analysis of role assignment in capture the flag  
 782 from Section 5.2 by presenting the complete set of reward and intervention curves over time (Figure 5). As before, we hone in on interventions over agent and flag-related concept estimates from  
 783 Agent 3 and Agent 4 (from the red team). Baseline performance with no intervention is shown as the  
 784 red-dashed line. As shown in Figure 5, test-time reward is zero both with and without intervention  
 785 in the early stages of training, as agents are first interacting with the environment. A first swing of  
 786 reward occurs around checkpoint 45, where it appears that the blue team has learned to capture the  
 787 red team’s flag (reward for the red agents is negative). This is followed by a quick counter-swing in  
 788 which the red team begins collecting positive reward. We will hone in on the intervention analysis at  
 789 checkpoint 55, as it is the first time we see the red team learn productive behavior. At checkpoint 55,  
 790 we see an interesting pattern in the intervention results for Agent 3 vs. Agent 4. Te team receives  
 791 positive reward, so we know that the red team is capturing the blue team’s flag, but the intervention  
 792 analysis reveals more precisely how that is done. First, in the top row of Figure 17, intervening over  
 793 Agent 3’s orientation concepts shows that the red team is negatively impacted by interventions over  
 794 Agent 3’s orientation estimate for Agent 1 (a blue agent) and Agent 3 (itself). Intervening over Agent  
 795 4 results in a slightly different pattern—every intervention over orientation concepts degrades the  
 796 red teams performance! As we have seen previously, agents tend to model the other agents that  
 797 they interact with the most, suggesting that Agent 4 is interacting with all of the other agents in the  
 798

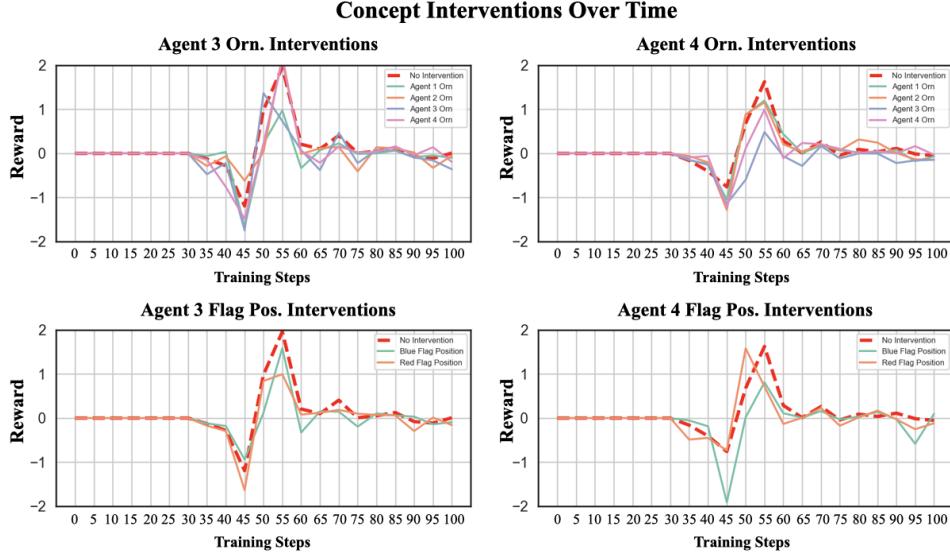


Figure 17: Concept interventions are performed intermittently over time for each agent. Results for the agents on the red team (Agent 3, Agent 4) are shown here. In the top row, we show results of an intervention over agent-related concepts (orientation), and in the bottom row we show results of an intervention over flag-related concepts (flag position). We find an interesting pattern in the intervention results at checkpoint 55. First, we find that Agent 4 is negatively impacted by all of the orientation interventions, whereas Agent 3 is only negatively impacted when intervening over the orientation of Agent 1 and Agent 3 (itself). Next, we find that Agent 4 is similarly negatively impacted by both flag interventions, whereas Agent 3 is only negatively impacted by an intervention over the red flag’s position. The intervention reveals, therefore, that Agent 4 is likely an attacking agent and Agent 3 is likely a defending agent.

799 environment, whereas Agent 3 is interacting with one blue agent. We see a similar pattern in the  
800 flag-related interventions as well (bottom row of Figure 5). There we find that, at checkpoint 56,  
801 Agent 3 is negatively impacted by an intervention over the red flag’s position, but not the blue flag’s  
802 position. Agent 4, however, is negatively impacted by intervention over both flag positions. This  
803 further suggests that Agent 3 interacts primarily

804 Altogether, the clues that we get from this intervention analysis suggest the following: (i) First, Agent  
805 4 interacts with each agent in the environment (spanning both teams) and also interacts with both  
806 flags. The intervention analysis reveals, therefore, that Agent 4 has learned an attacking strategy  
807 and is achieving positive reward by capturing the blue teams flag; (ii) Next, we see that Agent 3  
808 interacts primarily with its own flag (the red flag) and also interacts with one of the blue agents. The  
809 intervention analysis reveals, therefore, that Agent 3 has learned a defensive strategy, and is likely  
810 maintaining the red team’s positive reward by fending off a blue attacker.

811 **Qualitative Results** Figure 18 provides a snapshot of capture the flag behavior exhibited by the  
812 CBP agents’ at episode 55. As predicted by our intervention analysis in Section 5.2, one agent (R2)  
813 has learned an attacking role—traversing to the blue team’s base, capturing the flag, and returning it  
814 to the red base (and passing both blue agents along the way)—and the other agent (R1) takes on a  
815 defending role—camping at the red team’s base and defending it from a blue attacker. Altogether,  
816 these results confirm that intervention can be used *during training* to augment reward-based analysis  
817 and investigate strategic behaviors *as they emerge*.

## 818 D Lasso Neighborhood Selection

819 Lasso neighborhood selection is a simple method that poses graph learning as a Lasso regression  
820 problem [10]. Let  $X$  be an observation matrix that is composed of some set of random variables

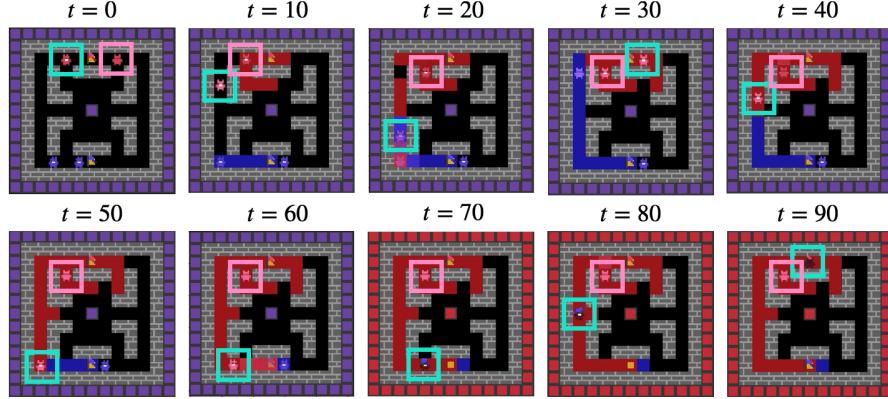


Figure 18: A Capture the Flag trajectory from trained CBP policies. The two red agents have learned emergent roles: attacker (teal square) and defender (pink square).

821 (rows), each of which is described by a set of features (columns). Lasso neighborhood assumes that  
 822 each variable can be approximated as a sparse linear combination of the observations (i.e. features)  
 823 of other variables. For a random variable  $x_i$ , this approximation is computed as:

$$\min_{\beta_i} \|\mathbf{X}_i - \mathbf{X}_{\setminus i}\beta_i\|_2^2 + \alpha\|\beta_i\|_1$$

824 where  $\mathbf{X}_i$  represents the features describing the variable  $x_i$  (i.e., transpose of the  $i$ 'th row of  $\mathbf{X}$ ),  
 825  $\mathbf{X}_{\setminus i}$  represents the features from rest of the variables (the remaining rows in  $\mathbf{X}$ ),  $\beta_i$  is a vector of  
 826 coefficients, and  $\alpha$  weights the L1-regularization term (enforcing sparsity).

827 Importantly, coefficients  $\beta_i$  determine which edges, if any, are connected to the node that represents  
 828  $x_i$ . In particular, for some additional variable  $x_j$ , an edge is established between  $x_i$  and  $x_j$  if either  $\beta_{ij}$   
 829 or  $\beta_{ji}$  is non-zero (or both). Intuitively, because a graph is a representation of pairwise relationships,  
 830 Lasso neighborhood selection posits that learning a graph is equivalent to learning a neighborhood  
 831 for each vertex—i.e., the other vertices to which it is connected—assumes that the observation at a  
 832 particular vertex may be represented by observations at the neighboring vertices.

833 In this work, we construct an observation matrix from the outcomes of interventions and use Lasso  
 834 neighborhood selection to model the relationships (i.e. similarities and dissimilarities) between  
 835 interventions.