# Human-Centered Reinforcement Learning

## Stephanie Milani

July, 2025

CMU-ML-25-108

Machine Learning Department
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA

### Thesis Committee

| | |
|---|---|
| Fei Fang | Carnegie Mellon University (Chair) |
| Geoffrey J. Gordon | Carnegie Mellon University |
| Hong Shen | Carnegie Mellon University |
| Katja Hofmann | Microsoft Research |
| Oriol Vinyals | Google DeepMind |

Submitted in partial fulfillment of the requirements for the degree of
*Doctor of Philosophy.*

# Abstract

Artificially intelligent (AI) agents are increasingly entrusted with making decisions in medium- to high-stakes domains such as healthcare, education, transportation, and cybersecurity. In these applications, agents make sequences of decisions that influence real-world outcomes. Reinforcement learning (RL) offers a natural and powerful framework for training such agents through experience. However, despite recent progress, there exist key barriers to deployment and adoption.

First, a capable RL agent can behave in ways that violate human expectations. In collaborative or safety-critical settings, unintuitive actions can confuse users or even create new risks. For example, an autonomous vehicle that suddenly swerves to avoid an accident may still be perceived as unsafe. This perception risks lack of adoption, even if the agent is a safer driver overall. Developing agents that exhibit *intuitive* behavior is therefore often a prerequisite for human-AI coordination and trust. Second, in safety-critical and regulated domains, the ability to explain and audit AI decisions is increasingly becoming a formal requirement. However, most RL agents make decisions using deep neural networks, which are challenging for people to understand. As a result, *interpretable* decision-making emerges as an important problem to address. Third, it is often challenging for designers to fully specify the range of desired behavior for an agent. Therefore, designers often leverage simpler proxy goals through fixed, simple reward functions. If this proxy goal is mis- or under-specified, agents can behave in ways that are misaligned with what people actually want. As a result, an important challenge is ensuring that agents are *aligned* with human intent, goals, and values.

These challenges all share a common theme: they arise because RL agents interact with or make decisions on behalf of huamns in human environments. As a result, a key question for the future of AI is how to develop agents that operate well with people. This dissertation advances a human-centered approach to RL to build and investigate AI agents that are interpretable, intuitive, and aligned. Here, we present technical advances in designing and evaluating AI agents, addressing key research questions that emerge from human involvement. Toward the goal of intuitive behavior, we design the first RL agent to pass a navigation Turing test and investigate why people perceive its behavior as human-like. Turing toward the goal of interpretability, we identify and algorithms toward two new dimensions of interpretability in RL: maintaining transparency in multi-agent decision-making and reducing the reliance on human annotations. We contribute a new alignment framing (decision-making) and introduce an algorithm that learns policies whose decision-making aligns with human preferences. Toward the goal of behavioral alignment, we contribute a benchmark and datasets for training and evaluating agents on fuzzy, underspecified tasks. We conclude this thesis by discussing how future work can leverage these ideas toward AI agents that support human flourishing.

# Acknowledgements

This thesis was the result of support of many amazing people.

I am deeply grateful for my advisor, Fei Fang. Fei, thank you for taking me on as an intern in the summer of 2019 and subsequently as your PhD student. You are perhaps the most hardworking and dedicated person I know. I am inspired not only by your passion for details and hunger to understand things well, but also by your unwavering commitment to your students.

I have come to think of Geoffrey Gordon as a pseudo-second-advisor. From Geoff, I have deepened my appreciation for the critical importance of rigor in research; our shared appreciation for clean and elegant ideas has influenced how I approach problems. Perhaps most importantly, you've shown me that one is never too established to write code or to deeply understand a new area.

Hong Shen, thank you for welcoming me into your lab and acting as yet another mentor to me. I feel incredibly fortunate to have had the opportunity to interact with you and your research group. From you, I have learned how to ask real, human-centered questions and think critically and carefully about ideas.

Katja Hofmann, first and foremost, thank you for Malmo! It has been a pleasure working with you throughout my Ph.D., from co-organizing the MineRL competitions to my transformative internship at MSR Cambridge. From you, I learned an underappreciated but extremely necessary skill in research: how to manage projects effectively and scope appropriately.

Oriol, it's been wonderful to interact with you through the years. I've deeply appreciated your support of MineRL, your generosity with your time at conferences, and your enthusiasm for ideas (especially those that connect learning and games!). Our conversations have been a consistent source of encouragement and inspiration.

During my PhD journey, I have had the unique privilege to work with many amazing mentors. At MSR, I was fortunate to be mentored by Ian Little, Harm van Seijen, Ida Momennejad, Sam Devlin, Guy Leroy, and Evelyn Zuniga. Mariya Toneva has been an incredible source of support: thanks to her, I spent a lovely week in Germany to visit MPI. Thank you to Dave Abel for always finding time to chat at conferences and offer thoughtful advice, and to Zhiyu Chen, for teaching me how to do NLP research. I'm profoundly grateful to Weiran Shen, for mentoring me through my first PhD project. During my job search, I was lucky to receive guidance from Max Simchowitz, Kianté Brantley, Berk Ustun, Ken Holstein, Nihar Shah, Eugene Vinitsky, Ben Eysenbach, Lily Xu, Ryan Shi, Zhiyu Chen, Lei Li, and more. I want to give special thanks to Kianté and Eugene for answering my Slack messages and texts at all hours, and to Kianté, Berk, and Max for challenging me to tell the right story.

CMU Machine Learning Department is one of the best places to do AI research, and this could not be possible without the fantastic support from staff like Diane Stidle, Dorothy Holland-Minkley, and Laura Winter.

# Contents

# List of Figures

# List of Tables

# 1 Introduction

We stand at a critical inflection point in artificial intelligence (AI). In recent years, rapid advances in AI have driven its widespread integration for prediction and decision support. Already, doctors receive classification-based risk scores to identify patients with seizure risk [327], volunteers receive push notifications through AI-powered recommendation systems [297], and teachers use predictive models to identify students at risk of falling behind [100]. But these systems do not act. In contrast, instead of just supporting decisions, newer AI systems are being designed to make them. These systems operate as agents: they perceive their environment and select actions in response [276]. This interactive capability means that AI agents are being designed to serve as *decision makers* in various sectors with medium to high risks, including education [314], transportation [65, 317], healthcare [56, 229], and cybersecurity [42, 210]. This proliferation of AI agents represents an important shift: we are moving from systems that simply predict to those that act in complex, real-world environments.

This shift renews focus on a long-standing problem—how we can train such agents. Reinforcement learning (RL) provides a natural solution [311]. RL formalizes decision-making under uncertainty by training agents to interact with an environment to optimize long-term rewards. This framework aligns well with the structure of many real-world tasks: in education, an AI agent selects a sequence of modules to support students' long-term understanding [81]; in transportation, AI agents coordinate traffic signals based on projected delays [345]; in healthcare, AI agents deliver behavioral health nudges to improve long-term patient outcomes [105], and in cybersecurity, AI-powered defenders may accept immediate losses to protect long-term valuable resources [181]. But once these agents are embedded in environments that involve people—as collaborators, observers, or end-users—the problem changes. Assumptions that hold in purely task-driven settings often fall apart, limiting the applicability of standard RL formulations.

First, RL agents are typically trained only to maximize task performance, without considering the human perceptions of the resulting behavior. But to better predict the behavior of others, people use past experience to develop mental models of how others in their environment will behave [32]. By not considering these perceptions, an otherwise performant RL agent can behave in unpredictable or unintuitive ways [299]. In collaborative or safety-critical settings, actions that are not intuitive can confuse users or even create new risks [73]. For example, a robot might quickly hand objects to a person to complete the task more efficiently, but the speed may feel too abrupt to the human partner, which may lead to a preference to complete the task alone at the cost of efficiency. Broadly, this unpredictable behavior can dissolve trust [70], risking lack of adoption even if the agent is a better partner overall. As a result, developing agents that exhibit *intuitive* behavior emerges as a prerequisite for human-AI coordination and trust.

Second, achieving strong performance in complex, high-dimensional domains typically requires representing policies with expressive function approximators, such as convolutional neural networks [226] and transformers [38, 50]. However, these architectures are notoriously opaque. Consequently, people face challenges in anticipating, trusting, or intervening in the agent's decision-making process [161, 241]. This lack of transparency is particularly important in safety-critical and regulated domains, in which the ability to explain and audit AI decisions is increasingly a formal requirement for deployment [93]. As a result, many applications now emphasize *interpretability* as a necessary aspect of agent design.

Third, because designers struggle to fully specify the range of desired behavior for an agent, they often leverage simpler proxy goals through fixed, simple reward functions. Although, in practice, many real-world objectives evade precise specification, this reward function is assumed to be precisely defined and consistent across all users [226, 283, 284]. If the proxy goal is mis- or under-specified with respect to the real objective, agents can behave in ways that are misaligned with what people actually want. As a result, there will often be some method for the agent to achieve high reward that is not what the designer intended [60, 162, 169, 179]. Furthermore, even if these agents behave the way we want, they may not make decisions in a way that aligns with our reasoning. Suppose, for example, a medical agent recommends a treatment due to spurious correlations (such as the patient's ZIP code rather than relevant clinical features). Then, the agent's decision is aligned, but its reasoning is not. We want agents that are *aligned* not only in their decisions but in how decisions are made.

Together, these challenges reveal a broader gap between traditional RL and the demands of real-world deployment: agents must be designed not only to optimize performance, but to work effectively with people. As a result, a key question for the future of AI is how to develop agents that operate well with people. Addressing this question requires rethinking how agents are trained and evaluated by treating people as integral to the process. **This dissertation advances a human-centered approach to RL to build AI agents that are intuitive, interpretable, and aligned.**

**Broader Work Toward Deployment.**  Although this dissertation focuses on human-centered RL, my research agenda broadly considers the current or future deployment of AI agents. This work spans three additional focal areas: improving sample efficiency, enabling coordination in multi-agent systems, and grounding learning algorithms in domains. Sample-efficient learning is an important barrier to deployment because real environment interactions may be costly or otherwise difficult to obtain. As a result, I have contributed to improved training schemes for multi-agent RL [333] and transformer-based RL [38], as well as a framework for transfer learning that captures when we can successfully transfer learned information across problems [255]. To spur further research interest and advances in this area, I have worked on designing [137] and releasing the influential MineRL benchmark and competitions [116, 115, 117, 155, 216]. Toward the goal of improved multi-agent learning, I contributed to work enabling agents to learn better representations for coordination tasks [333], as well as a benchmark and competition using Pokémon as a testbed for various multi-agent challenges, including adaptation to strategic opponents and long-context reasoning and planning [158]. Finally, toward the goal of building AI for concrete applications, I have contributed to methods for diverse domains, including cybersecurity [84, 215], education [211], and healthcare [342].

# 1.1 Overview of Contributions

We now provide an overview of the main contributions this dissertation makes to human-centered reinforcement learning, specifically to building agents that are intuitive, interpretable, and aligned.

## 1.1.1 Part I: Intuitive

Autonomous agents are increasingly deployed in interactive settings, in which people observe and interpret the behavior in a specific context. The ability to anticipate and interpret an agent's actions is central to trust, oversight, and collaboration. However, AI agents are typically trained only to maximize task performance, not to be predictable. As we will show, high performance alone does not guarantee that behavior will be intuitive to a human observer, meaning that achieving agents that act intuitively means intentionally shaping behavior to align with human expectations.

Chapter 3 focuses on this problem in the context of goal-directed navigation, a domain where we expect people to have strong human priors due to the ubiquity of navigation. Inspired by Turing's original imitation game, we operationalize intuitive behavior as that which is human-like. Specifically, we ask whether human judges can distinguish between the movements of AI agents and real people in a video game. This chapter introduces the first AI agent to pass a Turing Test for embodied navigation [212], establishes new methodology for measuring human-likeness, and provides empirical insights into how humans judge behavior as human-like.

### Chapter 3: Building and Understanding Human-Like Agents

Here, we investigate how people evaluate the human-likeness of AI agent behavior in a familiar and embodied setting: goal-directed navigation. In this work, our objective is not simply to make agents move in a human-like manner, but also to understand which behavioral features shape human judgments and how those judgments can guide the design of future AI agents. To this end, we propose a novel RL agent. We collect hundreds of crowd-sourced assessments comparing the human-likeness of navigation behavior generated by our agent and baseline AI agents with human-generated behavior. Our RL agent is the first to pass the Navigation Turing Test.

---

**Further work**

---

Not included in this thesis, we collaborate with domain experts in mental health to contribute the first LLM-based patient simulator for training mental health professionals [342]. We demonstrate how using a human-centered design process enables us to build more realistic patient simulations. We evaluate the patient simulator with real mental health professionals and trainees, demonstrating the promise of the tool for use in classroom settings.

## 1.1.2 Part II: Interpretable

If interpretability is necessary—whether it is due to regulatory demands or stakeholder needs—achieving it is particularly challenging in sequential decision-making problems. Although interpretability has been studied extensively in supervised learning, most existing methods assume access to static datasets and a single decision-maker, making them ill-suited for RL. These models often rely on abundant human supervision to label human-meaningful structures, like concepts, to facilitate interpretable predictions. In RL, agents collect data through interaction, posing practical challenges for this human-in-the-loop training. Furthermore, because these techniques are developed for predictive settings, they do not account for the possibility of multiple agents acting in an environment simultaneously. To overcome these challenges, this thesis develops methods for learning interpretable and performant policies while minimizing human annotation effort (Chapter 6) [353] and for the multi-agent setting (Chapter 5) [220].

### Chapter 4: Reducing Human Annotation Burden for Concept-Based Policies

Concept-based interpretabilityis a promising technique from supervised learning, where predictions are made using human-understandable attributes instead of opaque raw inputs. By training the model both to have high task accuracy and to accurately match experts' concept labels, these models learn a high-level concept-based representation that is simultaneously meaningful to humans and useful for machine learning tasks. In RL, we uncover a critical limiting assumption: prior work assumes concept labels are readily available during training. This requirement poses a significant limitation: it necessitates continuous real-time concept annotation, which either places an impractical burden on human annotators or incurs substantial costs in API queries and inference time when employing automated labeling methods.

To overcome this limitation, we introduce a novel training scheme, LICORICE, that enables RL agents to efficiently learn a concept-based policy by only querying annotators to label a small set of data [353]. We show how LICORICE reduces human labeling efforts to 500 or fewer concept labels in three environments, and 5000 or fewer in two more complex environments, all at no cost to performance. We also present the first exploration of VLMs as automated concept annotators for interpretable RL, finding them effective in some cases but imperfect in others. Our work significantly reduces the annotation burden for interpretable RL, making it more practical for real-world applications. More broadly, we are the first to uncover and address this fundamental assumption in concept-based explanations.

---

**Further work**

Not included in this thesis, we demonstrate how we can transform growing a decision-tree policy as an RL problem. It can equivalently be thought of as a concept selection process: each "meta"-timestep involves choosing a concept (feature) and value to test. This work enables RL agents to simultaneously learn good behavior and concise interpretable representations [321].

---

**Chapter 5: Extracting Interpretable Policies for Multi-Agent Coordination**

In many real-world tasks, a human operator will be responsible for a *team* of cooperating decision-making agents. In some tasks, like cybersecurity, there may also be an adversary or set of adversaries that at odds with the goals of the team. Previous work does not handle multi-agent transparent decision making, in which all agents must successfully act independently using interpretable policies while remaining robust to a variety of attackers, if applicable.

We propose the first algorithms, IVIPER and MAVIPER, for training interpretable policies in the multi-agent setting, combining ideas from model compression and imitation learning [217, 220]. In particular, we show that we can successfully train decision-tree policies that achieve high performance in environments that require coordination while maintaining interpretability through inherently-interpretable decision-tree policies. This approach is the first to demonstrate that coordinated behavior is still possible under interpretability constraints. Furthermore, we demonstrate that MAVIPER-trained agents remain robust to a variety of adversaries, which is an essential step toward real-world deployment.

> **Further work**
>
> Not included in this thesis, we extend this work to scale to more agents with an eye toward traffic signal control applications [51]. Through a novel algorithm, HYDRAVIPER, we demonstrate how we can achieve comparable performance to MAVIPER in a manner that is more scalable in the number of samples and in wall-clock time.

## 1.1.3  Part IV: Alignment

Many real-world applications involve fuzzy, context-dependent objectives that are difficult to encode mathematically. Consider the challenge of training an AI assistant to clean a house according to one's individual standards or navigate complex social situations. The "correct" behavior in these contexts depends heavily on individual preferences and situational factors. However, traditional RL relies on well-defined reward functions, which can struggle to capture the behavior we want from agents. This alignment challenge is further complicated by the need to handle trade-offs between multiple, potentially conflicting objectives and ways in which agents may arrive at decisions. An autonomous vehicle must balance safety, efficiency, passenger comfort, and adherence to traffic law, all while making decisions based on reasonable attributes. This thesis develops techniques to align agent decision-making with human preferences (Chapter 6) [221], as well as a benchmark and datasets for enabling agents to learn aligned behavior from human feedback for tasks that evade precise specification (Chapter 7) [214, 213].

## Chapter 6: Aligning Reinforcement Learning Policies with Human Feedback

Alignment techniques only focus on behavioral alignment of completions or trajectories, rather than the policies themselves. But in many cases, there is rich preference information about *how* agents should make decisions, not just *what* decisions they should make. For example, in loan approval systems, stakeholders may prefer policies that use information from more diverse feature categories rather than concentrating on a few. Interpretable RL provides an ideal foundation for policy-space alignment because it has established a rich vocabulary of measurable policy attributes—from decision tree depth to rule set size to feature importance—that naturally capture the dimensions along which humans have preferences about decision-making approaches. However, an unaddressed challenge in interpretable RL is enabling AI agents to integrate preference feedback into the policy generation process. Existing methods only collect feedback only after training is complete, neglecting opportunities to *inform* the learning process. To address this gap, we propose a novel framework, PASTEL, to train preference-aligned interpretable policies [221]. To our knowledge, we are also the first to formalize the notion of policy-space alignment.

PASTEL interleaves preference learning with an evolutionary algorithm, using updated preference estimates to guide the generation of better-aligned policies, and using newly-generated policies to query users to refine the preference model. Evolutionary algorithms enable the exploration of the full space of policies; however, it is intractable to maintain separate preference estimates—like win rates—for each individual policy in this infinite space. To handle this challenge, we propose to represent policies as feature vectors consisting of a finite set of meaningful attributes. To maximize the value of each user query, we employ a novel filtering technique to avoid presenting policies that are dominated in all dimensions, as repeated selections of clearly superior policies provides little information. We validate our method with experiments on synthetic preference data on two RL environments. We show that it produces RL policies that are better-aligned with user preferences and more efficient in the number of user queries.

### Further work

Not included in this thesis, we investigate the decision-making process of a large agentic model for Minecraft using a swathe of interpretability techniques [150]. We diagnose a safety failure case and provide a simple fix to better align the agent.

## Chapter 7: Learning from Human Feedback for Fuzzy Tasks

As AI agents become more capable, ensuring they pursue objectives aligned with human values becomes increasingly critical. Learning from human feedback (LfHF) represents a promising approach for training agents on complex tasks where specifying reward functions fails to capture human intent. However, the field lacks standardized benchmarks with high-quality human demonstrations and evaluation protocols needed to advance alignment research. We address this gap by introducing the BASALT Evaluation and Demonstrations Dataset (BEDD), the largest open dataset for learning from human feedback on visually complex, open-ended tasks that mirror real-world alignment challenges. This chapter combines two works [213, 214] to provide a more extensive overview of the BASALT benchmark, approaches, and resulting datasets.

BEDD consists of a collection of 26 million image-action pairs from nearly 14,000 videos of human players completing the BASALT tasks in Minecraft. It also includes over 3,000 dense pairwise human evaluations of human and AI agents. By dense, we mean that each evaluation includes a comparison of the relative task-completion performance of the agents and at least four additional questions—such as which agent was more human-like—resulting in 27,905 comparison points. Each dense comparison additionally includes natural language justifications. We establish concrete methodological standards for the BASALT benchmark and provide detailed analyses revealing that no AI agent has yet matched human performance. Even top-performing agents struggle with subjective aspects of task completion, demonstrating there exist challenges in learning what humans actually want versus what they specify. Our analysis highlights critical failure modes where agents technically complete tasks while violating implicit human expectations.

**Further work**

---

Not included in this thesis, we establish the BASALT challenge and host two years of NeurIPS competitions based on the challenge [154, 290, 291]. Most notably, BASALT is the first benchmark for evaluating fine-tuning from human feedback techniques for embodied agents.

# 2  Background and Preliminaries

In this chapter, we introduce important terms and background for reinforcement learning (RL), including Markov decision processes (MDPs), deep RL, and learning from human feedback. We then establish important concepts in interpretable RL. We include important notation in Table 2.1.

## 2.1  Reinforcement Learning Basics

RL refers to the study of agents that learn to make sequential decisions through trial and error, guided by reward signals. It encompasses both the formal problem setting and the suite of algorithms that aim to solve it. RL agents learn from interaction, without explicit supervision, by receiving evaluative feedback in the form of rewards. The learner and decision-maker is called the agent. The thing it interacts with, comprising everything outside the agent, is called the environment. These interact continually, the agent selecting actions and the environment responding to those actions and presenting new situations to the agent. A complete specification of an environment defines a task. The boundary between the agent and the environment is conceptually important: it delineates the extent of the agent's influence. The agent selects actions, but the resulting state and reward are determined by the environment, possibly stochastically.

### 2.1.1  Markov Decision Processes

To formalize this interaction, we use the mathematical framework of Markov decision processes (MDPs). MDPs [258] are a classical formalization of sequential decision making, and, as such, are used widely to model single-agent RL problems.

More concretely, the agent and environment interact at each of a sequence of discrete time steps, $t = 0, 1, 2, 3, \ldots$. At each time step $t$, the environment provides some representation of the state $S_t \in \mathcal{S}$, where $\mathcal{S}$ is the set of possible states. On that basis, the agent selects an action $A_t \in \mathcal{A}_t(s)$, where $\mathcal{A}_t(s)$ is the set of actions available in state $s$. In part as a consequence of the agent's action, the environment responds by returning a numerical reward $R_{t+1} \in \mathcal{R}$ and transitioning to a new state $S_{t+1}$. This interaction between the MDP and agent produces a sequence or trajectory $\tau = [S_0, A_0, R_1, S_1, A_1, R_2, \ldots]$. We use the term *rollout* to refer to generating such a trajectory by executing a policy within the environment, typically for a fixed number of time steps or until termination.

| Symbol | Description |
| --- | --- |
| $\mathcal{M}$ | Markov decision process |
| $s, s'$ | States |
| $a$ | Action |
| $r$ | Reward |
| $t$ | Discrete timestep |
| $\mathcal{S}$ | Set of all states |
| $\mathcal{A}$ | Set of all actions |
| $\mathcal{A}(s)$ | Set of all actions available in state $s$ |
| $P$ | Transition probability function |
| $R$ | Reward function |
| $\gamma$ | Discount factor $\gamma \in [0, 1]$, determining the relative importance of future rewards compared to immediate rewards |
| $\pi_\theta$ | Policy with parameters $\theta$, decision-making rule |
| $\mu$ | Behavior policy (off-policy) |
| $\tau$ | Trajectory (episode) |
| $G$ | Return (cumulative reward) |
| $V_\pi$ | State-value function under policy $\pi$ |
| $Q_\pi((s, a))$ | Action-value function under policy $\pi$ |
| $u_\theta(\tau)$ | Utility of trajectory $\tau$ |
| $d^\pi(s)$ | State distribution induced by policy $\pi$ |

Table 2.1: **Summary of notation for RL used throughout this dissertation.** When appropriate, we will occasionally use subscript $t$ to denote the timestep.

In a finite MDP, the sets of states, actions, and rewards $\mathcal{S}, \mathcal{A}, \mathcal{R}$ all have a finite number of elements. This finite structure allows for exact representation and computation of policies and value functions, though as we will see, real-world problems often require function approximation due to large or continuous spaces. In this setting, $R_t$ and $S_t$ are well-defined and depend only on the preceding state and action. That is, we can write the probability of particular instantiations of these random variables, $s' \in \mathcal{S}$ and $r \in \mathcal{R}$, occurring at time $t$ given particular values of the preceding state and action as:

$$p(s', r \mid s, a) \doteq \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}, \tag{2.1}$$

for all $s', s \in \mathcal{S}, r \in \mathcal{R}$, and $a \in \mathcal{A}(s)$. This function $p$ that governs the next state and reward defines the *dynamics* of the MDP. It is sometimes convenient to disentangle the transition dynamics $P$ from the reward function $r$, as we will see for e.g., learning from human feedback. For example, we may want to compute the state-transition probabilities $P : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \to [0, 1]$ or the expected rewards over state-action-next-state triplets $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$. These quantities are useful for *model-based* RL, where the transition and reward functions may be estimated from data and used for learning or planning with the learned dynamics model.

In an MDP, the environment's dynamics are fully specified by the transition function $p(s', r \mid s, a)$, which gives the probability of transitioning to state $s'$ and receiving reward $r$ after taking action $a$ in state $s$. This function encodes the *Markov property*, which asserts that the future is conditionally independent of the past given the present state and action. Formally, for all $s, s' \in \mathcal{S}, a \in \mathcal{A}(s)$, and $r \in \mathcal{R}$,

$$\Pr(S_{t+1} = s', R_{t+1} = r \mid S_0, A_0, R_1, \ldots, S_t = s, A_t = a)$$
$$= \Pr(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a).$$

This assumption is not a restriction on the dynamics themselves, but on how the state is represented. A state $s$ is said to be *Markov* if it captures all relevant aspects of the agent–environment history necessary to predict future outcomes. That is, the state must summarize all information from the past that can influence what happens next.

A core feature of RL is the use of a reward signal to define the agent's objective. The agent's goal is to maximize the total amount of reward it receives over time, making the reward function the central mechanism for specifying goals. Rather than being given explicit instructions, the agent infers desirable behavior by seeking to maximize a scalar quantity called the return, denoted $G_t$, which is some specific function of the reward sequence (often a summation). In this thesis, we primarily consider the *episodic* setting, which means that the agent–environment interaction decomposes naturally into finite-length episodes. Each episode begins in an initial state, which is typically sampled from some distribution of starting states $\rho$ and ends in a special state called the terminal state. To avoid overly myopic behavior, we typically define the return using a discount factor $\gamma \in [0, 1]$, which determines the weight placed on future rewards. Then we can write the expected discounted return as:

$$G_t \doteq \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}.$$

When $\gamma$ is close to 0, the agent will prioritize immediate rewards; when $\gamma$ is close to 1, the agent will prioritize long-term outcomes. Although this thesis primarily considers episodic MDPs, we note that the MDP framework is very flexible and admits a number of different formulations, such as the undiscounted and continuing cases.

## 2.1.2 Policies and Value Functions

Equipped with our MDP definition, we need to specify how an agent makes decisions. RL methods specify how the agent changes its policy as a result of its experience. Formally, a policy is a mapping from states to probabilities of selecting each possible action available in those states. The policy is denoted $\pi$, where $\pi(a|s)$ is the probability that $a_t = a$ if $s_t = s$. Learning then involves adjusting the policy based on experience, with the goal of improving the agent's behavior over time. Different RL algorithms achieve this in different ways.

Improving a policy requires some evaluation of the quality of states and/or actions. Value functions provide a means for this evaluation, specifying how good it is for an agent to be in a particular state or to take a specific action from a particular state. Almost all RL algorithms involve estimating these value functions. The state-value function $V_\pi(s)$ under policy $\pi$ is defined as:

$$V_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s] = \left[ \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1} \mid S_t = s \right], \quad \forall s \in \mathcal{S}. \tag{2.2}$$

Similarly, the action-value function $Q_\pi(s, a)$ represents the expected return when taking action $a$ in state $s$ and following $\pi$ thereafter:

$$Q_\pi(s, a) \doteq \mathbb{E}[G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]. \tag{2.3}$$

These functions provide a formal basis for evaluating and comparing policies. In particular, $Q_\pi(s, a)$ is useful for selecting actions, while $V_\pi(s)$ summarizes the overall desirability of states under policy $\pi$.

The optimal value functions describe the best possible performance achievable in the environment. The *optimal state-value function* $V^*(s)$ gives the maximum expected return achievable from state $s$, and is defined as:

$$V^*(s) \doteq \max_\pi V_\pi(s).$$

Optimal policies also share the same optimal action-value function $Q^*(s, a)$, defined by the Bellman optimality equation:

$$Q^*(s, a) \doteq \max_\pi Q_\pi(s, a) = \mathbb{E} \left[ R_{t+1} + \gamma \sum_{s'} P(s' \mid s, a) \max_{a'} Q^*(s', a') \right],$$

$\forall s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$. Here, the expectation is over the possible next states $s'$, and $R_{t+1}$ denotes the expected immediate reward for taking action $a$ in state $s$. Once one has $Q^*$, it is relatively easy to determine an optimal policy. In fact, for finite MDPs, we can precisely define an optimal policy using value functions. Because value functions define a partial ordering over policies, a policy $\pi$ is at least as good as another policy $\pi'$ if its expected return is greater than or equal to that of $\pi'$ for all states.

In practice, agents must estimate these value functions from experience. RL algorithms differ in how that experience is collected. In *on-policy* methods, the agent improves its policy using data gathered by that same policy. In contrast, *off-policy* methods allow learning from data collected by a different policy, known as the *behavior policy*, denoted $\mu$. Off-policy learning is useful not only when interaction with the environment is limited, but also when reusing past experience is desirable for improving sample efficiency. Sample efficiency becomes particularly critical in large or high-dimensional environments in which tabular representations of the policy and/or value function(s) are no longer feasible.

### 2.1.3  Deep Reinforcement Learning

In simple, finite MDPs, it is often feasible to represent policies and value functions as explicit tables, with one entry for each state or state-action pair. This setting is called the *tabular* case. For an excellent overview of tabular methods, we refer an interested reader here [311]. However, many real-world settings have continuous state or action spaces, or otherwise have too large of a state space to enumerate. As a result, function approximation becomes necessary to generalize across states and support learning from limited experience.

Modern RL often uses deep neural networks as function approximators, a setting known as deep reinforcement learning (deep RL) [226]. In this case, the policy is often parameterized by a neural network with parameters $\theta$ and actions are sampled according to the probability distribution output by this network:

$$a_t \sim \pi(s_t; \theta).$$

To learn effectively, an RL agent interacts with the environment to collect trajectories and periodically updates the network parameters $\theta$ to maximize cumulative reward.

To support learning, the agent may also approximate the value function $V^\pi(s)$ or the action-value function $Q^\pi(s, a)$ using a separate neural network with parameters $\psi$. These estimates are used to guide the policy update, for example by reducing variance or bootstrapping future returns. In actor–critic methods, the policy (actor) and value function (critic) are trained jointly, with the critic estimating the long-term value of states or actions under the current policy (e.g., as in Equations (2.2) and (2.3), respectively). Modern deep RL algorithms such as PPO [284] and SAC [119] use these learned approximations to improve sample efficiency and training stability while optimizing policies over high-dimensional observation and action spaces.

### 2.1.4  Learning from Human Feedback

So far, we have assumed that $R$ is given by the environment. But, importantly, in many cases $R$ is *designed* by a person for a particular task. In these cases, it may be challenging or impossible for a human to write down a function that not only covers all possible edge cases but also does not lead to specification gaming or degenerate behavior. As a result, *learning from human feedback* has emerged as a paradigm for tasks in which a reward function is difficult to specify or incomplete.

More broadly, different modalities of human feedback provides additional learning signals, enabling the agent to learn key components of the environment rather than assuming they are fully specified. But, most commonly, this feedback is used to estimate $R$, which will be our primary focus for this section. In this section, we introduce two key forms of feedback—demonstrations and preferences—that we build upon in later chapters.

**Learning from Demonstrations**

A demonstration typically consists of a trajectory $\tau$ of state-action pairs collected by some expert policy $\pi^*$. There exist multiple ways to leverage such data, including (but not limited to) inverse RL, in which we aim to recover the reward function of the expert, and imitation learning, in which we aim to imitate the expert's policy. Within imitation learning, techniques differ in the assumptions they make about the learner's ability to interact with the environment or query the expert. Some methods operate purely offline, learning from a fixed dataset of demonstrations, while others require online access to either the environment or additional expert supervision. However, here we focus on two forms of imitation learning, as we will build on this class of techniques in Chapter 5.

**Imitation learning.** In imitation learning [140, 236], the goal is to learn a policy $\hat{\pi}$ that reproduces similar behavior to the expert $\pi^*$. A standard approach is *behavioral cloning*, which treats the expert data as i.i.d. samples for supervised learning. Here, the learner does not have environment or expert access after data collection. Given a dataset of expert-labeled state-action pairs $\mathcal{D} = \{(s_i, a_i)\}_{i=1}^N$, the agent minimizes a loss function such as cross-entropy:

$$\hat{\pi} \doteq \arg\min_{\pi \in \Pi} \sum_{(s,a) \in \mathcal{D}} \ell(\pi(a \mid s), a).$$

While imitation learning is frequently applied to human demonstrations, it also generalizes to settings where the expert is a neural network policy, which may be trained with RL. In this case, imitation learning can serve as a form of distillation [97, 133], which can enable the transfer of behavior to a learner with other desirable properties (such as interpretability or fewer parameters to reduce inference time).

**Interactive imitation learning.** A standard issue in behavioral cloning is distribution shift. In sequential decision-making tasks, this assumption often fails: small errors in the learned policy can lead to states not covered by the training data, resulting in compounding errors that degrade performance. In other words, $d^\pi(s) \neq d^{\pi^*}(s)$, where $d$ is the distribution induced by following a particular policy. Interactive imitation learning reduces this distributional mismatch by collecting expert labels on states encountered by the learner during training. For example, DAgger [270] frames imitation learning as an iterative algorithm that alternates between policy execution and expert supervision. Roughly speaking, at each iteration $m$, the current policy $\hat{\pi}_m$ collects trajectories, and the expert $\pi^*$ provides action labels for the encountered states. The algorithm aggregates these state-action pairs into a growing dataset,

$$\mathcal{D}_{m+1} \doteq \mathcal{D}_m \cup \{(s, \pi^*(s))\},$$

and trains the next policy to mimic $\mathcal{D}_{m+1}$. Over time, the aggregated dataset reflects the distribution of states that the policy is likely to encounter under its own execution, thereby reducing distributional shift.

## Learning from Preferences

Preferences serve as another useful signal when reward evades precise specification. Rather than assigning absolute scores to individual trajectories, humans may simply indicate which of two behaviors they prefer. These comparisons can be used to recover a latent utility function that guides policy learning. Preference-based learning enables alignment with implicit human objectives even when those objectives cannot be explicitly codified. We will build on this class of techniques in Chapters 6 and 7.

**Preference modeling.** Typically, we need to make some assumption on the structure that generates the underlying preferences. A standard assumption is that the underlying user preferences follow the Bradley-Terry model [28], which posits that each trajectory $\tau$ is associated with a latent utility $u_\theta(\tau)$ with parameters $\theta$. Given two trajectories $\tau_1, \tau_2$, the probability that a person prefers $\tau_1$ over $\tau_2$ is:

$$\Pr(\tau_1 \succ \tau_2 \mid \theta) \doteq \frac{\exp(u_\theta(\tau_1))}{\exp(u_\theta(\tau_1)) + \exp(u_\theta(\tau_2))}. \tag{2.4}$$

This formulation is analogous to the Elo rating system used in competitive games [88], where preferences induce a ranking over alternatives based on their underlying utilities.

**Learning and optimization.** Given a dataset $\mathcal{P} \doteq \{(\tau_i^1, \tau_i^2, y_i)\}_{i=1}^N$ of preference-labeled pairs, we can then estimate $\theta$ by maximizing the likelihood:

$$\mathcal{L}(\theta) \doteq \sum_{i=1}^N \left[ y_i \log \Pr(\tau_i^1 \succ \tau_i^2 \mid \theta) + (1 - y_i) \log \Pr(\tau_i^2 \succ \tau_i^1 \mid \theta) \right].$$

Once we learn the utility function $u_\theta(\tau)$, it can be used to train policies via standard RL [57]. In many settings, we define a reward function $\hat{R}_\theta(s, a)$ such that the cumulative reward over a trajectory approximates $u_\theta(\tau)$. The agent then solves:

$$\hat{\pi} \doteq \arg\max_\pi \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t \hat{R}_\theta(s_t, a_t) \right]. \tag{2.5}$$

It has become especially influential in large-scale settings, such as language model fine-tuning [246, 262], where aligning models with human values or judgments is more tractable through comparison data than explicit reward design.

Figure 2.1: **Interpretable RL taxonomy and its relationship to the RL problem.** Explanations may explain individual actions $a$ in terms of the input state $s$ (Feature Importance), influential experiences from training (Learning Process and MDP), and long-term behavior (Policy-Level).

## 2.2 Interpretable Reinforcement Learning

There has been increasing interest [34, 319] in *explaining* AI agents, including those trained using RL, to gain insight into the agent's decision-making process. This increased interest is in part due to large initiatives such as the DARPA Explainable Artificial Intelligence project [112], which evolved from a narrower focus on supervised learning to include broader tasks in AI. Here, we explain important concepts in interpretable RL using our novel taxonomy for different threads of work in the space [219]. This taxonomy focuses on the RL problem, with clear ties to key parts of the RL framework.

### 2.2.1 An RL-First Taxonomy

The most common taxonomy for organizing explainable machine learning methods groups them into two categories: intrinsic and post-hoc. *Intrinsic* interpretability refers to the direct construction of understandable models, such as small decision trees [187, 272]. The use of intrinsically interpretable models produces a system that is already compliant with interpretability requirements [274]. *Post-hoc* interpretability refers to creating an additional model to explain an existing one. Post-hoc and intrinsic interpretability are also useful properties for understanding and organizing XRL techniques. Indeed, we may want to learn an intrinsically interpretable policy or explain it using a number of post-hoc techniques, such as constructing a surrogate model that imitates the non-interpretable policy. This type of taxonomy fails to emphasize the unique challenges and opportunities that emerge when applied to RL. In particular, it does not capture important components of the RL problem, such as action selection and long-term behavior.

We now explain our RL-first taxonomy. Figure 2.1 depicts how the categorizations correspond to parts of the RL framework. Since a practitioner starts with a goal for using explanations, we prioritize partitioning based on goal: which part of the agent or learning process does one wish to explain? *Feature importance* methods identify the features that affect an agent's action choice $a_t$ for the input state $s_t$. *Learning process and MDP* methods show the past experiences or the components of the MDP that led to the current behavior. *Policy-level* methods illustrate the long-term behavior of the agent.

This taxonomy is more tightly-coupled to the RL problem to enable practitioners to more easily select explanations that better suit their needs. For example, if a practitioner wants to investigate the importance of specific features in the agent's decision-making process at a given state, they can use feature importance explanations. Doing so may enable them to diagnose whether agents are utilizing appropriate features for action selection and adjust their algorithms accordingly. In contrast, if a practitioner wants to understand the agent's interaction with the environment and its learning process over time, they can focus on learning process and MDP explanations. If a practitioner wants to summarize the long-term behavior of the agent, they can use policy-level explanations. Because this decomposition first prioritizes *what* the practitioner wants to understand, it provides more of a use-case-centered perspective on explanations for RL.

### 2.2.2 Feature Importance

Feature importance explanations provide *local* insights into at an agent's decision-making process by identifying which components of the input state $s$ were most influential in selecting a particular action $a$. Formally, given $s$, these methods analyze $a$ via the policy $\pi$ or (action-)value function $Q_\pi$. There exist many possible ways to explain this local decision-making.

Some techniques are faithful, in that they reflect the true features and decision boundaries used by the agent. In contrast, other methods generate post-hoc *rationalizations* that approximate how a human might explain the same behavior. These rationalizations may improve human understanding but do not necessarily reflect the agent's internal computations. In this thesis, we primarily focus on methods that produce intrinsically interpretable policies, rather than generating explanations for an existing opaque model. As a result, we now turn to a discussion of different policy representations that support interpretability.

### Policy Representations

Although there are many types of representations that could be considered interpretable, such as natural language, this section focuses on two broad classes of policy representations: decision trees and programs and logical structures. In this thesis, we primarily focus on decision trees as a policy representation (Chapters 5 and 6), though we note that many algorithms designed for extracting decision-tree policies can be easily augmented to admit alternative structures.

**Decision trees.**   Decision trees are hierarchical models that recursively partition the input space along a specific feature using a cutoff value. These models produce axis-parallel partitions: internal nodes are the intermediate partitions, and leaf nodes are the final partitions. When used to represent policies, the internal nodes represent the features and values of the input state that the agent uses to choose its action, and the leaf nodes correspond to chosen actions given some input state. Decision trees are often considered interpretable because the reasoning process behind each decision is explicitly encoded in the tree's structure [187]. A user can examine the path from root to leaf to understand which features influenced a particular action. There exist many algorithms for training or extracting decision-tree policies [17, 24, 66, 361] and Q-functions [148]. In addition to classic decision trees, there exist a number of other extensions to this structure, including so-called soft decision trees [113, 97], which are decision trees with sigmoid activations rather than Boolean operators, oblique decision trees [351, 68], and other more expressive tree forms [79, 129, 177].

**Programs and logical structures.**   Beyond decision trees, a variety of symbolic representations have been proposed to capture policies in more expressive yet interpretable forms. One class of approaches represents policies as programs [334], written in a domain-specific or general-purpose language. Other approaches leverage more general logical structures: Tsetlin machines [330] use finite automata to encode logical rules; conceptual embedding methods [67] produce interpretable intermediate representations; and fuzzy controllers [128], based on fuzzy logic [357], enable approximate reasoning via continuous truth values. However, despite their ability to resolve the issues with classic decision trees, most of the aforementioned approaches have not been evaluated with user studies to understand their actionability (do the models enable users to take meaningful decisions?), comprehensibility (does the target audience understand the model?), or preferability (do end users prefer the model over alternative models, including uninterpretable ones?).

### 2.2.3  Learning Process or MDP

These explanations provide additional information about the effects of the training process or the structure of the underlying MDP. These methods reflect how learning dynamics or reward design influence the final behavior. However, these explanations may be global or local, depending on what is being explained (e.g., behavior or individual actions).

Most commonly, the agent reveals its learned transition dynamics $\hat{P}$ through e.g., structural causal models [200] or generative models of the environment [49, 275], but other techniques decompose the reward function to produce explanations in terms of the agent's objectives. Most work in this area focuses on presenting interpretable views of $Q_\pi(s, a)$ by decomposing $R$ into a set of additive terms with semantic meaning, called reward decomposition [152, 9]. Reward-based explanations may help laypeople understand how reward-based objectives influence the decision-making process of a RL policy [9]. Still other techniques determine training points that influenced the learned policy [106]. Influence commonly refers to a change in some function value (e.g., $Q_\pi(s, a)$). These approaches can provide insight into data efficiency, failure cases, or unexpected dependencies in learning.

## 2.2.4 Policy-Level

A less-studied but important category of explanations is those which describe an agent's longer-term behavior. These policy-level explanations typically summarize long-term behavior through abstraction or representative examples. They answer questions like, "how will the agent behave over time?". By definition, these explanations are global.

At a high level, techniques fall into three sub-categories: summarizing behavior via experiences encountered during training, converting internal model components into human-interpretable representations, and extracting clusters or abstract states. Within the first sub-category of techniques, we aim to describe the agent's long-term behavior based on influential trajectories (not individual transition tuples), as in Section 2.2.3. These techniques primary differ with respect to the selection criteria for choosing the trajectories and how to combine them for the summary [8, 80, 139, 173]. The second sub-category of techniques primarily focuses on extracting RNN policies [168]—typically, into finite state representations [41, 242]. This approach differs from Section 2.2.2 because RNN policies depend not only on the current state but also the trajectory history encoded in the hidden state, making per-step, state-action explanations insufficient. The third and final sub-category of techniques explain the overall behavior of the agent in terms of how it will act when it encounters similar states, which can then be organized as graphs over important landmarks [306] or expected future transitions [322]. They exclude information in the state aggregation techniques by definition of abstraction.

# Part I

# Intuitive

# 3 Building and Understanding Human-Like Agents

## 3.1 Introduction

One aspect of human-centered RL involves creating agents that exhibit behavior that is *natural* and *intuitive* to those who interact with it. Consider a common driving scenario: a car attempts to merge into another lane during traffic. A nearby human driver observes the merging vehicle's gradual deceleration, anticipates the intended merge, and adjusts accordingly. Coordination emerges from shared expectations and familiar behavioral patterns. Now imagine that the merging vehicle is controlled by an AI agent trained using RL. The agent, rewarded for minimizing travel time and avoiding collisions, may learn to execute a maneuver that is optimal under that reward function: it accelerates aggressively to close the gap, then brakes sharply to merge at the last moment. Although technically successful (there were no collisions and the maneuver was efficient), this behavior can appear erratic or threatening to nearby drivers. It violates the informal social contract of driving, where behaviors signal intent and facilitate coordination.

More generally, when an RL agent behaves in a way that is familiar to humans, its behavior becomes more predictable and easier to interpret. This predictability is critical for building appropriate trust and understanding between humans and AI systems as we move into more integrated environments. As a result, in this chapter, we define intuitive behavior as behavior that resembles human behavior, since people naturally interpret and respond to actions through the lens of their own experiences. Indeed, prior work has shown that agents perceived as human-like are more enjoyable to interact with [304] and are better suited for collaborative tasks that require shared understanding [39, 118, 281, 339]. In safety-critical domains like autonomous driving, AI-powered vehicles must behave sufficiently human-like for human drivers to interpret, anticipate, and act in their presence [127].

However, human-likeness is not a guaranteed byproduct of high task performance. This intuition is illustrated in Figure 3.1. In some cases, two policies may achieve the same task performance, but one will do so in a way that also satisfies human likeness. In other cases, there may be an explicit trade-off: actions that satisfy the objective of expressing human likeness may be in conflict with actions that achieve high task performance. In either case, this disconnect—optimizing for $x$ (e.g., navigation efficiency) while hoping for $y$ (e.g., intuitive social behavior)—is a well-documented failure mode in AI systems [162]. As a result, if we want $y$, producing agents that exhibit human-like behavior, we must explicitly design and optimize for it.

Figure 3.1: **High task performance is not sufficient for human-like behavior.** In this example, a hypothetical agent is tasked with moving from the initial state to the goal state. The actions that are driven by optimizing for task performance can differ from those that would be driven by expressing human-like behavior. We suspect that these actions can result in inconsequential differences in the overall task performance while having substantial differences in subjective evaluations of human likeness.

What constitutes human-like behavior is not universal. Instead, it is context-dependent. Although we can hand-craft attributes that we believe correspond to human-like behavior, it remains unclear whether these attributes correspond to the criteria people actually use when making comparative judgments about agent behavior in context. To evaluate and refine these assumptions, human evaluations enable studying both the choices people make when identifying human-like behavior and the justifications they provide. Understanding the behaviors that contribute to people's perceptions of human likeness is a foundational first step towards achieving general human-like behavior of artificial agents.

To study these questions, we desire a behavior that is ubiquitous in human environments: navigation. Specifically, we study agents that must move through a 3D embodied environment from one point to another. We focus on 3D navigation tasks because they provide a rich yet tractable setting to isolate human-likeness in agent movement. Unlike open-ended tasks with many confounding variables, navigation offers a clean way to study how humans interpret movement patterns. Furthermore, this form of navigation is pervasive in many video games, making it a key area of interest for game developers [7, 86]. More generally, it is considered fundamental to embodied biological intelligence [125, 238], making it of interest to cognitive scientists [233, 268] and researchers broadly interested in intelligent behavior [348]. It has also been a key area of interest in HCI [329] due to how people (or robots) navigate in real, augmented, or entirely virtual spaces.

To study human-like navigation in video games, we leverage a Turing test. The Turing test is a well-established methodology for studying whether AI is indistinguishable from humans. This navigation-specific Turing test, called the Human Navigation Turing Test (HNTT) [77], asks human judges to indicate which of two videos demonstrates more human-like behavior. The judges then justify their decision and indicate their certainty about their choice. However, in the HNTT work, the authors did not instantiate a statistical test to definitively conclude whether an agent passed the HNTT. But both studied AI agents exhibited non-human-like behavior. As a result, producing an agent that passes the HNTT still remains an open challenge.

To this end, we design a novel agent to pass the HNTT. To assist with our design of a human-like agent, we inspect the resulting behavior of the two baseline agents from prior work [77]. With these insights, we design our novel agent—the *reward-shaping* agent—using simple and intuitive techniques. We then conduct a behavioral study on Amazon Mechanical Turk (MTurk) of the HNTT to investigate the behavior of our agent and the baselines. To determine whether agents pass the HNTT, we propose a firm criterion: a statistical test that determines whether human judges distinguish between human and agent behavior at a level that is *significantly different* from chance. We then validate the conclusion of previous work: the two baseline agents are not sufficiently human-like because they do not pass the HNTT. In contrast, human judges cannot reliably distinguish between the behavior of our *reward-shaping* agent from one controlled by a person. To our knowledge, this agent is the first to pass the HNTT.

To understand these assessments, we analyze the free-form responses to determine which characteristics people believe are representative of human and AI navigation behavior. We annotate the responses with codes that summarize the provided rationale. These annotations show that there are key differences between how people characterize human-like and non-human-like behavior. Specifically, people utilize the same high-level characteristics when describing human-like and non-human-like behavior, but the presence or absence of these characteristics strongly informs their judgments. Based on the findings of our analysis, we summarize considerations when developing and evaluating the human likeness of AI agents.

**Contributions.** In summary, we make the following contributions:

- We contribute a novel RL agent that exhibits more human-like navigation behavior.

- We validate that this agent passes the HNTT by conducting a large-scale behavioral study with 192 participants (1152 human-agent comparisons) and using our newly-proposed criterion. This is the first agent that passes the HNTT, meaning its behavior is indistinguishable from a human's.

- To more deeply understand the factors that contribute to assessments of human likeness, we conduct an extensive qualitative analysis of 395 free-form responses. Using these annotations, we find that there are key differences between how people characterize human-like and non-human-like behavior. We believe these findings are useful for developing and evaluating the human likeness of AI agents in embodied navigation, video games, and beyond.

## 3.2 Related Work

In this section, we review prior work on building human-like agents, including a detailed section on *reward shaping*, as we will use this technique later. We then review prior work on evaluating human-like AI agents, identifying key gaps that motivate our approach.

**Building human-like agents.**     Researchers have taken various approaches to address the challenge of developing human-like AI agents in games, including learning from demonstrations [145, 157, 208], RL [12, 98, 208, 363], and more [224, 318]. RL [311] provides a generally-applicable set of algorithms for learning to control agents in settings including (but not limited to) modern game environments [7, 116, 176, 313, 336]. It also offers significant benefits as an approach for generating navigation behavior [7]. In particular, the use of RL may enable more complex navigation abilities (such as grappling or teleportation) while being relatively hands-off for game designers to use.

**Reward-shaping.**     Although RL agents learn effective navigation by maximizing reward, they make no consideration for the *style* with which they act [7]. If these approaches are to be adopted in commercial game development, practitioners have firmly asserted that controlling style is essential [146]. As an extreme example, RL approaches that have recently defeated world champion human players at modern games demonstrated unusual behaviors [135] that made collaborative play between human and AI in mixed teams far less successful [22]. Simply maximizing the task-specific reward signal is unlikely to produce human-like agents. Reward shaping [236, 347] is a simple yet powerful technique that allows practitioners to clearly specify the desired agent behavior. This approach involves crafting a reward signal that provides dense feedback to the agent. We use reward shaping to help generate more human-like behavior because it is a flexible, intuitive way for practitioners to control the agent's behavior by specifying objectives instead of optimizing unintuitive hyperparameters. Additionally, reward shaping can be used with any RL algorithm, making it possible to swap in and out the underlying algorithm as needed.

**Evaluation of human-like AI.**     There is no standard set of metrics for evaluating human-like AI. One paradigm involves measuring human similarity with proxy metrics for human judgments. Some of this work measures the task performance of the AI [318, 363], but, as we will show, this metric is an insufficient proxy for human similarity. Other work assesses how well the AI agent can predict the subsequent human action [145] or align its behavior with people [224, 316], but these metrics do not include evaluations with real people, which is vital for assessing human likeness in games [104] and beyond [358]. In contrast, studies that do involve human evaluations tend to be small-scale surveys [98, 208] that often offer only a preliminary investigation into the specific characteristics that inform beliefs and do not include a Turing test [104, 163]. Work that uses a Turing test often does not investigate the behaviors or provide concrete metrics [232], or it focuses on assessing the full spectrum of game behaviors [12, 98, 208]. Due to the complexity, providing concrete recommendations to game designers is challenging. In contrast, we isolate a specific but widely-used behavior: point-to-point navigation. To perform our assessment, we utilize the setup of the recently-proposed Human Navigation Turing Test [77]; however, we propose and perform a deeper evaluation of human assessments of AI and human behavior.

Figure 3.2: **Navigation task as observed by study participants (screenshot, left), and detail of the mini map of the game level (right).** Agents spawn on the island outside of the main map, which is shown in the bottom portion of the mini map on the right. They must jump to the main area and navigate to the goal location. The light blue containers in the left screenshot represent the goal location.

## 3.3 Problem Setup

We use the same navigation task and game from previous work [77]. Underlying the game is the crucial mechanic of goal-directed navigation: players must move from one location to another to collect key items, go to drop-off areas, and engage in combat with other players. As a result, navigation between points represents an abstraction of the most common task in the game. To concentrate on navigation-specific characteristics, we use a simplified game version that excludes other mechanics and objectives.

**Navigation task.**   We instantiate the navigation task in the same way as prior work: a single agent must navigate to a target location. The left screenshot of Figure 3.2 shows this location, indicated by the three blue containers. Before the player moves, the navigation target spawns uniformly at random in one of 16 possible locations, denoted by the green crosses in the right-hand image of Figure 3.2.

**HNTT.**   The HNTT asks human judges to identify which of two navigation behaviors more closely resemble how people navigate in *reality*. This phrasing aims to capture how *convincing* an agent is [190], in contrast to another Turing test interpretation: whether a human or AI agent *controls* an entity. We choose this phrasing because we want to create *convincing* NPCs that contribute to an immersive game experience; we do not wish to deceive the player into thinking that an agent is controlled by a person. We discuss these considerations in Section 3.8. This objective also aligns with our goals in bringing AI agents into real-world settings: we want people to collaborate with robots in their workplaces and homes; we do not want to convince them that a human is controlling the factory robot or their robot vacuum.

Figure 3.3: **All agents successfully learn to navigate**. This plot shows the average amount of time needed to solve the task as a function of training time (shaded area is the standard deviation). For reward-shaping, N=3; for hybrid and symbolic, N=4. Curves are smoothed with a rolling window of 200. Importantly, all agents converge to solve the task in around 60 steps (around 12 seconds of in-game time). In contrast, agents start out needing around 140 steps (around 28 seconds of in-game time) to solve the task. This shows that performance differences are not responsible for differences in human likeness, and standard metrics of task performance are insufficient to assess human likeness.

**Baseline agents.** Previous work [77] introduce two types of RL agents: a *symbolic* and a *hybrid* agent. However, participants could accurately detect human players above chance when evaluating these agents, meaning that people did not perceive their behavior as sufficiently human-like. To progress toward the goal location, the agents take actions from a prespecified set (called an action space). This action space consists of 8 possible actions: do nothing, move forward, and move left and right (30, 45, and 90 degrees on each side). The agents receive a dense reward signal $R_{\text{base}}$ that only encourages successful navigation to the goal as quickly as possible. It consists of the following terms: a -0.01 per-step penalty to encourage efficiency, a -1 one-time penalty for dying because the agent may fall off the map, an incremental reward for approaching the goal, and a +1 reward for reaching the goal. The main difference between these two agents is the observations that they take as input. The *symbolic* agent receives only a semantic, low-dimensional representation as input; the *hybrid* agent also receives an image input.

## 3.4 Building a Human-Like Reinforcement Learning Agent

To design our *reward-shaping* agent, we analyze the *hybrid* and *symbolic* agents to find influential characteristics on the previous judgments of human likeness. Based on this analysis, we introduce a novel agent for the HNTT: the *reward-shaping* agent. This agent extends the *hybrid* agent with two critical changes to promote learning human-like behavior: additional terms to the reward signal and an expanded available action space. To test whether our contributions result in differences in perceptions of human likeness, we fix all other components of our *reward-shaping* agent to be the same as the *hybrid* agent.

**Identifying problematic behavior.**   Because the *symbolic* and *hybrid* agents previously exhibited non-human-like behavior, we inspect examples of their generated navigation and isolate three classes of problematic behavior. Agents would wildly swing camera angles or make sudden turns (**P1**), frequently collide with walls (**P2**), and move more slowly than expected (**P3**).

**Fixing problematic behavior with reward shaping.**   To correct these behaviors, we use reward shaping [269] by including terms to penalize undesired behavior. We start with the same reward function used by the *hybrid* and *symbolic* agents to encourage navigation, $R_{\text{base}}$. To discourage abrupt camera movements (**P1**), we include a penalty $\alpha$ for swift camera angle changes over a set 0.15 difference threshold value between consecutive timesteps. To reduce wall collisions (**P2**), we evaluate whether a collision is detected at time $t$ (denoted as collision$_t$), then apply a penalty $\beta$ of -0.05 if a collision is indeed detected. To discourage stagnation of movement (**P3**), we provide a penalty $\eta$ of -0.01 if the distance traveled between timesteps, defined as dist$_t$, is lower than an environment-specific threshold value of 220 map units. Putting it all together, the new shaped reward function $R_{\text{rs}}$ becomes:

$$R_{\text{rs}}(s_t, a_t) \doteq R_{\text{base}}(s_t, a_t) + R_{\text{shape}}(s_t, a_t), \text{ where}$$
$$R_{\text{shape}}(s_t, a_t) \doteq -\alpha \cdot \mathbf{1}[\Delta\theta_t > 0.15] - \beta \cdot \mathbf{1}[\text{collision}_t] - \eta \cdot \mathbf{1}[\text{dist}_t < 220].$$

**Fixing problematic behavior with action-space shaping.**   To encourage smoother control and avoid abrupt turns, we utilize an approach similar to action-space shaping [156]. We introduce finer-grained controls by adding more available actions to the agent. Specifically, we extend the action space to 14 actions from the previous 8 by increasing the degrees of turning left and right from 3 to 6. The updated list of turning degrees for this agent is: 18, 36, 45, 54, 72, and 90 on each side.

**Producing high-quality navigation.**   We train all agents to a similar level of performance on the navigation task to ensure that task skill is not responsible for the perceived differences in human likeness. We measure task proficiency using the number of steps needed to reach the goal, where each step corresponds to around 5 seconds of real-time play. Figure 3.3 confirms that all agents learn to reliably reach the goal, indicated by the performance near the end of training. A skilled agent now takes approximately 60 steps to complete the task (about 12 seconds of real-time play). This result ensures that differences in the human-likeness of assessments are not due to differences in the ability of the agents to solve the task.

## 3.5  Experimental Design

To understand what characteristics people believe are indicative of human likeness, we conduct a behavioral study with human participants. Our setup closely follows prior work [77], but, for completeness, we detail the full study design here.

| Study | Number of Participants | Number of Trials |
|---|---|---|
| Human vs. *hybrid* | 50 | 6 |
| Human vs. *symbolic* | 50 | 6 |
| Human vs. *reward-shaping* | 92 | 6 |

Table 3.1: **Conditions tested in each study and the number of trials per condition.** The *hybrid* and *symbolic* studies replicate prior work [77] to validate switching to a crowd-sourcing platform.

**Experimental task.**   We ask each participant to act as a judge by completing a survey of 6 HNTT trials. In each trial, the judge views two side-by-side video stimuli of people or agents completing the navigation task. Each video corresponds to a trajectory $\tau$. After watching these videos, the judge indicates which video they believe navigated more like a human would in the real world, provide a justification of their response, and report their certainty. More specifically, the judges answer the following questions:

1. **Which agent navigates more like a human would in the real world?** This decision is a forced binary choice.

2. **Why do you think this is the case? Please provide details specific to the video.** The judge answers this question as a free-form response.

3. **How certain are you of your choice?** The judge answers this question on a 5-point Likert scale, with choices ranging from extremely certain to extremely uncertain.

To mitigate subject learning effects from sequentially viewing multiple videos, we do not reveal to the judges which of the videos featured an reinforcement learning agent. In other words, participants complete each task and, in the end, do not know the ground-truth labels for whether the agent featured in each trajectory was controlled by a human or not.

**Experimental procedure.**   We complete 3 studies; each study pits a human-controlled agent against a different RL agent (Table 3.1).[1] Within each study, all judges view the same 6 trials but in a randomized order per judge. Within each trial, the ordering of the two videos is randomized, such that the answer cannot be inferred by presentation order. Each participant first reads through an introduction page with the required task instructions. They then complete a consent form and read through a background page with brief details about the game. They answer a series of questions to assess their comprehension of the task and familiarity with video games. Finally, participants engage in the 6 HNTT trials. Figure 3.4 shows the comprehension and familiarity questions and an example HNTT trial.

---

[1]Data use under MSR-LA license. License details are in the original authors' GitHub https://github.com/microsoft/NTT. This link includes *all* data used in this study, including from our *reward-shaping* agent.

I understand this task takes approximately 30 minutes, and that I won't be paid extra if I take longer or won't be paid if I've completed this task before.
○ I agree
○ I disagree

I understand this HIT has a 1-hour duration, and I can return the HIT at any time within this 1-hour, but I won't be paid for returned tasks or partial completions.
○ I agree
○ I disagree

I understand that I need to complete all the questions.
○ I agree
○ I disagree

How familiar are you with Third Person Action* video games?
*game where the camera during gameplay is primarily in a third-person perspective
○ Never heard of them
○ I am aware but have never played them
○ I play only sometimes
○ I play on a regular basis
○ Other

How familiar are you with the video game [title]?
○ Never heard of it
○ I am aware but have never played it
○ I play only sometimes
○ I play on a regular basis
○ Other

(a)

Please watch the videos below. Then, answer the questions below. One video is an AI agent, the other could be an AI agent OR a human. The objective is to identify **which video navigates more like a human would in the real world.** Assume the human is a competent player and knows the map.

Which video navigates more like a human would in the real world?

**Video A navigates more like a human**
○

**Video B navigates more like a human**
○

Why do you think this is the case? Please provide details specific to the videos on this page.

How certain are you of your choice?
○ Extremely certain
○ Somewhat certain
○ Neither certain nor uncertain
○ Somewhat uncertain
○ Extremely uncertain

(b)

Figure 3.4: **HNTT survey questions.** The screenshot in (a) shows the comprehension and familiarity questions (asked once per participant). We gauge the participant's familiarity with the time the task will take, understanding of task completion, familiarity with third-person action video games, and familiarity with the video game used in the survey. The screenshot in (b) depicts one HNTT trial. We ask participants to choose their response to the human likeness question, justify it, and indicate their certainty.

**RL details.** We train all three agents with PPO [284], a popular deep RL algorithm that is empirically robust and effective in a wide range of tasks [89]. The *hybrid* and *reward-shaping* agents receive an additional input of a 32x32 cropped depth buffer visual input which the *symbolic* agent do not. The visuals present a third-person view of the agent in the environment. To process this additional visual channel, the *hybrid* and *reward-shaping* agents are equipped with a convolutional neural network which learns to extract high level visual features that are then concatenated with a representation of the symbolic inputs. We train each of the three agent architectures for 15 hours, the equivalent of 10 million training timesteps, on at least 3 different random seeds. We train all agents using Tensorflow 2.3 [2] and the OpenAI Baselines PPO implementation [78] with a distributed sampler. For a full list of training hyperparameters, please refer to Table 1. We choose them because we found this set to perform best on preliminary experiments. To effectively train agents in a complex video game setting, we use a distributed approach leveraging an in-house sample collection framework and Azure cloud resources. Training samples are collected from a scaleset of 20 low priority GPU virtual machines (Azure NV6), each running 3 video game instances. The samples are then sent to one training head node, a CPU-only Azure E32s memory-optimized virtual machine.

**Navigation video production.**   A key part of the study is the videos that we show to the human judges. For the data and videos of human-controlled agent navigation, we use the publicly-available sample published by previous work [77]. We sample human videos from the 40 videos under their "study 1" protocol. To generate the AI agent navigation data, we select each agent's most recently saved checkpoint. Then, we instantiate a new session and deploy the agent in the game 100 times, producing 100 total navigation videos per agent. To produce the video stimuli that we show to the human judges in the study, we sample the recordings uniformly at random.

**Study standardization.**   We implement several measures to standardize the videos and ensure that any measurement noise applied to all conditions. First, we check that any changes in light applied similarly across conditions. Second, we design the timing of the stimuli to ensure that participants had sufficient time to engage in and provide meaningful responses in all trials. Third, because the goal locations may differ depending on the game-controlled initialization, we match the goal locations of the human videos with the AI agent videos. Consequently, we use different human videos for different studies. Fourth, we apply the post-processing steps from prior work [77]), including masking identifying information, adding a "For Research Purposes Only" watermark, and cutting out the last few seconds of the human videos (this change corrects an effect of the data collection process, where the human players manually ended their recording, adding a few seconds at the end of the videos).

**Other experimental control.**   We use the MTurk crowd-sourcing platform [249], which is widely used for data collection and research due to its scalability as long as researchers implement appropriate steps for quality control [142]. We set the following MTurk requirements for survey participation: location is United States, age is 18 or older, and language is English. We do not collect demographic information or any other personally identifiable information. To target more experienced MTurk Workers, we set the following Human Intelligence Task (HIT) qualifications: HIT Approval Rate greater than 98%, Number of HITs Approved greater than 500, and a qualification to prevent repeat responses. To incentivize quality, we include a bonus payment for each high-quality response. We review the free-form answers to find low-quality or suspected bot responses; for example, we exclude from analysis responses with high instances of typos, copy/pasted answers, or nonsensical wording. We pay all participants who completed the task for the HIT, even if their response was identified as low-quality. The low-quality responses do not receive the bonus payment. We pay on average 15 USD per hour. We obtain approval for our studies from our Institutional Review Board (IRB) and informed consent from each participant. We include details of the study and a description of any potential participant risks in the consent form.

## 3.6  Assessing Human-Likeness

We first aim to identify which agents pass the HNTT. Because existing work demonstrates differences in assessment ability depending on expertise, we seek to identify whether this phenomenon holds in our setting. We finally seek to investigate the relationship between self-reported uncertainty and accuracy when assessing the agents. We instantiate the following research questions:

| Choice | Free-Form Response | More Human-Like | Less Human-Like |
|--------|--------------------|------------------|-----------------|
| B | The character in Video B runs in <mark>straight lines</mark> and <mark>goes to where he needs to be going</mark>. The character in Video A is <mark>running in circles, into objects</mark>, etc. | Smoothness of movement +; Goal directed + | Collision avoidance −; Goal directed − |

Table 3.2: **Example coded response to the question, "Which video navigates more like a human would in the real world?"**. This judge believes Video B is more human-like (leftmost column). The highlighted text illustrates the in-text annotation. The judge notes that the more human-like agent runs in straight lines (more human-like code: smoothness of movement +) and navigates to the goal (more human-like code: goal directed +), while the agent that they believe is less human-like runs in circles (less human-like code: goal directed −) and into objects (less human-like code: collision avoidance −).

**RQ 1.** Which agents are judged as being human-like?

**RQ 2.** Do judges exhibit greater accuracy in assessing human likeness as a function of their experience with games?

**RQ 3.** What is the relationship between the accuracy of human judges and their self-reported uncertainty?

## 3.6.1 Analysis

To answer **RQ 1**, we propose a firm criterion for deciding whether an agent is sufficiently human-like, formalizing the question: *are human judges unable to distinguish between agent and human behavior*? We implement this criterion as a statistical test that determines whether human judges distinguish between human and agent behavior at a level significantly different from chance. We instantiate this test by computing the 95% confidence interval (CI) for the median of the human-agent comparisons with bootstrap sampling. If the confidence interval includes $0.5$ (chance-level agreement), then the agent passes the HNTT. In other words, if the human judge essentially flips a fair coin to choose which of the two is more likely to be human, then they are unable to distinguish between the two. To calculate the 95% CI, we run each bootstrap over 10000 iterations.

For both **RQ 2** and **RQ 3**, we compare our variables of interest with *accuracy*. We define accuracy to mean that the participant identified that the human navigating was more human-like than the AI navigating. Specifically, we report the *median* accuracy. To answer **RQ 2**, we compare whether the accuracy of the participants of the study is related to the self-reported familiarity with action games in general and the specific game in the study. To answer **RQ 3**, we examine the self-reported uncertainty of the judges and its relationship to accuracy.

| Agent | Median Accuracy (IQR) [95% CI] | Median Uncertainty (IQR) |
|---|---|---|
| *symbolic* | $0.83 \, (0.67 - 1.00) \, [0.67, 1.00]$ | $2.17 \, (1.67 - 2.42)$ |
| *hybrid* | $0.83 \, (0.67 - 1.00) \, [0.83, 1.00]$ | $1.92 \, (1.33 - 2.25)$ |
| *reward-shaping* | $0.50 \, (0.33 - 0.67) \, [0.44, 0.63]$ | $2.17 \, (1.75 - 2.67)$ |

Table 3.3: **Full summary statistics of accuracy and uncertainty.** We show the median accuracy (IQR=Q1-Q3) for each agent, reported as non-parametric measures of central tendency and spread; we report the 95% confidence interval and median uncertainty (IQR=Q1-Q3) of the human-agent comparisons for each agent. Only the *reward-shaping* agent passes the HNTT.

## 3.6.2 Results

Using the aforementioned analysis, we now answer **RQ 1-3**. Table 3.3 shows the full summary statistics, which we compute over the full dataset from our survey. We find that, of the three agents, only our reward-shaping agent passes the HNTT, making it the first agent to do so (**RQ 1**). We also find that game familiarity does not influence ability to assess human likeness (**RQ 2**) and that self-reported uncertainty is lower for the reward-shaping and *symbolic* agents than the *hybrid* agent.

### RQ 1: Human-Likeness

**The *symbolic* and *hybrid* agents do not pass the HNTT.** We first confirm the findings of the previous work: the two baseline agents do not pass the HNTT. The judges exhibit median accuracies of $0.83$ for both agents, indicating that they indeed distinguish the agents from humans higher than chance level. However, recall that our test is based on the CI including chance-level agreement ($0.5$). Importantly, neither CI includes chance-level agreement (*symbolic* agent, 95% CI=[0.67, 1.0]; *hybrid* agent, 95% CI=[0.83, 1.0]), meaning that neither agent passes the HNTT.

**Only the *reward-shaping* agent passes the HNTT.** In contrast, our *reward shaping* agent passes the HNTT: the 95% confidence interval includes 0.5 (chance-level agreement). This result suggests that the judges cannot consistently differentiate between the *reward shaping* agent and the human player (*reward shaping* agent, median accuracy=0.50, 95% CI=[0.44, 0.63]). The computed CI is also small and largely centered around 0.5. [2] We, therefore, answer our **RQ 1**: only the *reward-shaping* agent is judged as human-like. This finding means that we are the first to build an agent that passes a Turing test for navigation.

---

[2]Because the sample sizes of the trials differ (50 samples for the human vs. *hybrid* and human vs. *symbolic* conditions; 92 samples for the human vs. *reward-shaping* condition), we subsample the data for the *reward-shaping* agent to 50 samples, then run the bootstrap sampling procedure 100 times. We find that the computed CI always contains 0.5 in each run of the bootstrap. The average median accuracy is $0.50$, with a variance of $0.00$; the averaged CI is [0.44, 0.63], with a variance of 0.01 for the lower bound and 0.00 for the upper bound.

## RQ 2: Game Familiarity and Ability to Assess Human-Likeness

Because we want to understand which factors influence assessments of human likeness, we assess whether the ability to assess human likeness is significantly influenced by self-reported game familiarity. We ask participants whether they are familiar with the specific game used in the study, as well as video games in general.

**There is no relationship between game familiarity and ability to accurately assess the human likeness of the RL agents.** We perform a multiple linear regression analysis to test whether either type of game familiarity significantly predicts accuracy in assessing human-likeness. Through this analysis, we find there is no relationship between either of the self-reported familiarities and accuracy for all agents (*symbolic*: $R^2 = 0.01, F(2,47) = 0.21, p = 0.814$; *hybrid*: $R^2 = 0.06, F(2,47) = 0.26, p = 0.261$; *reward-shaping*: $R^2 = 0.02, F(2,89) = 0.94, p = 0.393$). Decomposing the results further, neither specific game familiarity (*symbolic*: $\beta = -0.01, p = 0.878$; *hybrid*: $\beta = -0.03, p = 0.377$; *reward-shaping*: $\beta = -0.02, p = 0.522$) nor general game familiarity (*symbolic*: $\beta = 0.03, p = 0.525$; *hybrid*: $\beta = 0.06, p = 0.109$; *reward-shaping*: $\beta = 0.04, p = 0.191$) predicted accuracy. These findings suggest that game familiarity is generally *not* predictive of accuracy for this specific task, answering **RQ 2**. In contrast, previous findings have demonstrated a relationship between the ability to assess human likeness and familiarity with the domain of study. This result may differ because we study a relatively simple setting, in which most people have strong priors about what constitutes human likeness. Navigating by walking or running is an activity that most people perform or observe daily, so we likely have a strong internal sense of human-like movement—even if we are not familiar with games that require navigation.

## RQ 3: Confidence and Ability to Assess Human-Likeness

We assess the median uncertainties of participants (Table 3.3); perhaps counterintuitively, lower values indicate more certainty (less confidence) and higher values indicate less certainty (higher confidence).

**Human judges report lower certainty about the *symbolic* and *reward-shaping* agents.** Participants report higher certainty when assessing the *hybrid* agent (median=1.92, IQR=(1.33-2.25)) compared with the *reward-shaping* (median=2.17, IQR=(1.75-2.67)) and *symbolic* agents (median=2.17, IQR=(1.67-2.42)). We believe that the judges may have been less certain of their assessments of the *symbolic* agent due to differences in the lengths of the videos: on average, the videos of the symbolic agents were 8.3 seconds long, whereas the hybrid agent videos were 15.3 seconds long. They may have not had enough time with the agent to confidently assess it. Although the judges felt less certain about their assessments of the *reward-shaping* agent compared to the *hybrid* agent, they more accurately detected human behavior in the presence of the *hybrid* agent. This result is interesting because it aligns well with accuracy: in general, human judges seem to be well-calibrated. When they are more accurate (as when assessing the *hybrid* agent), they are also more confident; when they are less accurate (as when assessing the reward-shaping agent), they are also less confident. This result answers **RQ 3**.

Figure 3.5: **Analysis pipeline for understanding characterizations of human-like and non-human-like behavior.** Using a smaller subset of the free-form responses, two annotators iteratively convene to generate a list of codes and protocol for coding the larger subset of data. We then analyze the resulting data to understand which characteristics the human judges believe are indicative of human-like or non-human-like behavior.

## 3.7  Assessing Human-Like Characteristics

We ask the following research questions to analyze the *characteristics* that underlie the assessments of human likeness:

**RQ 4.**  Are there key differences between how people characterize human-like and non-human-like behavior? Does this differ when the agent does or does not pass the HNTT?

**RQ 5.**  What is the relationship between the characteristics that people use to assess human likeness and their ability to accurately assess it?

### 3.7.1  Analysis

Figure 3.5 shows the overall analysis pipeline.  To enable comparison between an agent that does not pass the HNTT with one that does, we sub-sample responses from the *hybrid* agent and the *reward-shaping* agent studies. We first randomly sample a set of 53 responses to compute the initial agreement, called the *agreement sample*. We then construct the sample for analysis by randomly sub-sampling 3 free-form responses per judge for each study. To minimize bias, we shuffle responses before sampling. This procedure results in a dataset of 395 responses for our analysis of human-like characteristics.

| Code | Shorthand | Definition | Key Phrases | Example Snippet |
|---|---|---|---|---|
| Smoothness of movement | smooth | Quality of the agent's navigation or camera movement | Smooth, jerky, swerve, steady, straight, fluid | Movements are way more smooth |
| Goal directed | goal | How intentional the agent's behavior seems | Intention, focus, knew where to go | Deliberate camera movements |
| Collision avoidance | avoidance | Whether the agent avoids collisions | Collide, avoid, crash, runs into obstacle | Runs into a box |
| Environment receptivity | receptivity | Whether the agent understands and/or properly interacts with the environment | Explore, stay on path, collect power-ups | Ignores all the health/mana/etc |
| Intuition | intuition | The judge cannot pinpoint behaviors | Natural, feeling, seems to be | Just a feeling |
| Self-reference | self-reference | Relationship to the judge's own movement or play | Like I play | [Like] how I navigate with that ... view |

Table 3.4: **Annotation code definitions.** The codes used to label the free-form responses are in the leftmost column, accompanied by the corresponding shorthand for the codes (middle-left column) and the definition (middle column). The middle-right column lists the keywords and phrases that the annotators used to determine if a response could be labeled as containing a particular code. An example snippet of a response that would be labeled with that code is provided in the rightmost column.

**Coding process.** We follow a pair-coding approach to annotate the data. One annotator proposes an initial list of codes derived from previous work [368]. Following established notation [10], we have a set of $I$ items (or responses), labeled as at least one of the $K$ categories by $C = 2$ coders. We decompose each label as more or less human-like $H = (\texttt{more}, \texttt{less})$ and quantify its direction $D = (\texttt{more}, \texttt{less})$, when applicable. The direction refers to the item itself, not how human-like it is. For example, if we label item $i$ as *smoothness of movement* $(K)$, we note whether the judge considered the behavior human-like $(H)$ and whether they noted it as being more $+$ or less $-$ smooth $(D)$. The two annotators then convene to discuss the meaning of the codes and jointly code a set of 5 responses. Table 3.2 illustrates an example of a coded response. After that, the two annotators independently code the agreement sample with the initial set of codes, with the option to label responses as *other* and provide specific examples to enable revisions of the codes if other themes emerged. The two annotators iteratively reconvene to discuss disagreements and refine the codes. After multiple rounds of discussion, independent coding, and disagreement resolution, the annotators fix the set of codes (Table 3.4) and their inclusion criteria to label the full sample.

**Coding protocol.**    Because we aim to design human-like AI agents, we prioritize codes that could be utilized by AI designers. For that reason, when deciding on codes, we prioritize codes that refer to specific behaviors over more general ones. For example, a collision avoidance behavior could be coded as goal-directed; however, we code it only as collision avoidance. This protocol promotes the independence of categories while prioritizing specific, lower-level behaviors to use in designing agents. When coding, the annotators first consider whether the response could be categorized as a lower-level code, then move to more general codes if needed.

**Inter-annotator agreement.**    To ensure sufficient alignment between the two annotators and show that the aggregation of the independently-coded data is reasonable, we compute the inter-annotator agreement for the agreement sample. Specifically, we calculate agreement using binary Cohen's kappa $\kappa$ [62] over $K$, $D$, and $H$. The annotators achieve an overall average inter-annotator agreement of $0.84$, indicating that there is relatively high agreement for this data sub-sample. After fixing the list of codes, the annotators receive the sub-sample of data to code such that they overlap on 25% of the data (99 items). We also calculate the overall average $\kappa$ for this sample. This time, $\kappa = 0.78$, only $0.06$ lower than the agreement sample, indicating that our coding process is fairly general.

**Annotation codes.**    We now discuss the annotation codes and inclusion criteria in more detail. Table 3.4 includes these definitions and phrases that helped us identify the presence of each code. For each code, we provide a supporting example to give the reader a sense of what common responses may look like. *Smoothness of movement* refers to the quality of the agent's navigation or camera movement. This code considers both immediate jerky actions and temporally-extended zig-zagging behavior. *Goal directed* refers to how intentional the agent's behavior appears. We include descriptions of behavior that pertain to a perceived goal, even if that goal is not the primary one. We include the code *collision avoidance* because it is a long-standing area of research in the robotics community [285]. This code refers to intentional behavior to redirect from a potential crash. *Environment receptivity* aims to capture the agent's relationship with the game environment, its contextual understanding, and adherence to norms. In a real-world setting, this might look like a person walking on a path instead of the grass or crossing the street when permitted by a pedestrian signal. Any responses that refer to non-specific feeling that a behavior was more human-like are categorized as *intuition*. We include this code to capture instances where participants can identify what they believe is more human-like behavior but struggle to express it. Finally, we include *self-reference* as a code to capture when judges relate the agent's behavior to their own play.

**Disagreement resolution.**    During the iterative coding process, the annotators assess the likely causes of disagreements. After resolving easy-to-resolve issues, the remaining disagreements arise from individual differences in interpreting ambiguous natural-language responses. This cause means that neither annotator can be treated as more correct for disagreement resolution. The annotators, therefore, use the following disagreement resolution scheme. When a disagreement arises in at least one label for an item annotated by both annotators, we randomly choose an annotator to treat as correct and use their labels.

Figure 3.6: **Proportion of codes describing human-like and non-human-like behavior.** The judges more frequently characterize human-like behavior as being more smooth, receptive and responsive to the environment, and goal-directed. In contrast, participants more frequently describe non-human-like behavior as being less smooth, receptive and responsive to the environment, and goal-directed.

## 3.7.2 Results

In all plots, we use the shorthand version of the codes, noted in Table 3.4, along with the $+$ and $-$ notation. The $+$ and $-$ notation indicates the direction of the code. For example, smooth $+$ indicates that the participant referenced more smooth movement; smooth $-$ indicates a reference to less smooth movement.

### RQ 4: Characteristics of Human-Like and Non-Human-Like Navigation

Because the total number of codes the judges use may differ, we instead investigate the relative number of times a code was used compared to all codes used to describe either human-like or non-human-like behavior (human-like and non-human-like code proportions should sum to 1) (Figure 3.6).

**Human judges rely on similar high-level characteristics when assessing human-like behavior.**
We find that judges rely on similar high-level characteristics when characterizing human-like behavior. The judges most frequently refer to smoothness of movement, goal directedness, and environment receptivity when assessing the agents. The only difference is the ordering. Judges more frequently characterize human-like behavior as more smooth, receptive and responsive to the environment, and goal-directed.

Figure 3.7: **Codes describing human-like and non-human-like behavior, further decomposed by high- and low-accuracy judges**. We compare the proportions of codes the human judges use to describe human-like behavior (left) and non-human-like behavior (right). We decompose these codes by high- and low-accuracy judges to determine whether more accurate individuals rely on different features to rationalize their decisions. Interestingly, we see that the judges rely on similar characteristics but to different degrees.

**Agent type does not influence the proportion of characteristics used.** We investigate these responses based on agent type but do not find a difference between the resulting proportions. This finding is surprising because we might expect the characteristics to change substantially depending on the type of agent being assessed. However, this finding supports the idea that people may have relatively stable beliefs regarding what constitutes human-like behavior. Therefore, the rationale is only sometimes useful: in other words, looking for the jerkier agent only makes sense if the AI has not been designed to be less jerky than the person. We, therefore, conclude that, although people rely on different specific characteristics to determine human likeness, the general characteristics are relatively stable across different agents, answering **RQ 4**.

## RQ 5: Characteristics and Ability to Assess Human-Likeness

To investigate how the use of codes differs depending on the accuracy of the judge, we divide the participants into two groups. The first group consists of high-accuracy judges ($> 80\%$ of responses correctly choosing the human videos as being more human-like); the second group consists of low accuracy judges ($\leq 80\%$ of responses indicating the more human-like agents aligned with the human-generated navigation). After performing this decomposition, we then analyze the group-level codes that the judges more frequently use to describe human-likeness (Figure 3.7).

**High and low accuracy judges refer to categories at a similar rate when assessing human-likeness.** We first assess the category-level usage (so, over $K$ and $H$ but not $D$). We find that high- and low-accuracy judges generally rely on similar high-level categories. For example, both types of judges refer to the high-level code of smoothness of movement in $40\%$ of their codes when describing human-like behavior, and they refer to the high-level code of smoothness of movement in around $49\%$ of their codes when describing non-human-like behavior. These results indicate that there is no little in their tendency to rely on this characteristic to explain behavior.

**More accurate human judges exhibit different beliefs than less accurate ones.** However, when we factor in $D$, we see that high-accuracy judges more commonly describe smooth motion when describing human-like behavior; in contrast, low-accuracy judges more often mention smoothness as a characteristic of non-human-like behavior, further supporting the idea that people's beliefs about AI capabilities may inform their assessments. In our case, the low-accuracy participants seem to share a belief that an AI agent is more capable than a person (by producing more smooth or "perfect" navigation). Interestingly, both high- and low-accuracy judges utilize intuition and self-reference to a similar level of frequency when assessing human-like behavior. Taken together, these findings answer **RQ 5**.

## 3.8  Discussion and Future Directions

This chapter establishes that we can indeed design human-like agents using relatively simple and generalizable techniques. Our study reveals that only the agent *designed* to display more human-like behavior passed our test of human likeness, highlighting the importance of explicitly incorporating these objectives when building agents. Although this result does not preclude the emergence of a particular capability, it emphasizes that incorporating our goals into our optimization target is important for achieving it.

However, determining what exactly constitutes human likeness requires careful consideration from designers. We also highlight how conducting mixed-methods evaluations of feedback from human judgments can assist in more deeply understanding the relationship between the ability to assess human-likeness, beliefs about human-likeness, and other factors (such as uncertainty). Although conducted in a limited scope of video game navigation, our findings and methodology should assist with future work on designing and evaluating human-like agents. In this section, we discuss the trade-offs of our design choices, then discuss future work toward building more general human-like agents.

### 3.8.1  Design Choices

Our study requires a number of design choices. These choices include the scope of the behavior studied, the phrasing of the primary question asked, and the type of evaluation conducted (e.g., observational or interactive). We discuss some of the trade-offs of these choices, including the scope of the study and potential misunderstanding of the question.

## Scope of Human-Like Behavior Studied

Our study specifically evaluates the human-likeness of third-person perspective point-to-point navigation behavior of agents in a video game. We choose to conduct a deeper investigation of a more narrowly-defined behavior. This type of focused, behavior-specific study offers several key advantages over broader evaluations of AI human-likeness. By concentrating on third-person navigation behavior, we can gain deeper insights into the specific processes that distinguish human movement from algorithmic approaches. Broader investigations may conflate different behaviors, making it challenging to isolate the factors in the specific context.

**Other forms of navigation behavior that warrant investigation.** Although this type of navigation is present in many settings, like pedestrian navigation in driving simulator [338] there are many other forms of navigation that exist in both real-world and virtual environments. Each of these types presents unique challenges and requires different strategies for designing human-like behavior. Although our study does not address all types of navigation, it provides a valuable starting point for evaluating the human-likeness of agents in one specific type of navigation. The codes that we identify are general enough to provide a starting point for researchers to analyze different forms of navigation. For instance, collision avoidance is a general characteristic that is persistent in many domains featuring diverse types of navigation, like driving and running. Future work should consider expanding these evaluations to provide a more comprehensive understanding of how to design agents that behave in a more human-like manner.

## Different Interpretations of Human-Likeness

Additionally, the analysis of the free-form responses reveals that there were different interpretations of the human likeness question. Some judges relate the movement directly to human navigation in the real world. As an example, one judge comments,

> *In real life a human would almost certainly not jump down as far as the character in video A did without severely hurting themselves.*

However, others related the movement to how human players would *control* an agent in a video game. For example, another judge mentions,

> *... in Video A, the player bumps into a wall briefly before readjusting. This is something humans do when they get distracted and look away for a moment.*

**Investigation of interpretations of the question.** To investigate this disagreement, the two annotators annotate the agreement sample with which interpretation of the question the subject answered: real-world human navigation, video game navigation, and unclear. The two annotators have a high agreement for this annotation (Cohen's kappa: $\kappa = 0.94$). Based on the annotation, it is largely unclear which question the subjects were answering (40 out of 53 responses). However 11 responses refer to *video-game* navigation, while only 2 responses are clearly about real-world human navigation. We suspect that including the video game familiarity questions primed subjects to believe that the question was about video-game-specific navigation, rather than general human-like navigation. In future studies, we recommend that the study designers clarify which question is asked of participants by including an additional question that asks the participants the other interpretation of the question to provide an obvious contrast or describing the situation in which they would like the participants to envision themselves.

**Applications of different interpretations of human likeness.** Both interpretations of the question are reasonable, depending on the application. When designers seek to automate parts of the development process, such as playtesting, it is more important to create agents that appear to be *human-controlled*. In automated playtesting of games [111, 286], AI agents that act like real users would enable video game designers to expedite the development process while also alleviating the burden of game players to extensively evaluate new content. Users could give feedback only after obvious bugs, like those related to movement, have been corrected, which may increase enjoyment of the feedback process. In *shared autonomy* [6], developing agents that behave like that user would enable a more seamless integration of semi-autonomous control with user inputs. For example, we observe that the judges called out strafing as an example of what a human would do in a video game. Strafing is a tactical, sideways maneuver that would not be performed by a person navigating in the real world. Incorporating these game-specific movements would likely increase the perception of the agent being human-controlled, especially by expert players. The creation of such agents would enable players who experience disruptions, like network issues, to still play cloud games [234]. When the system detects a disruption, it can take control and begin emulating that person so that when the user can take back control, they can do so seamlessly. This option can be included for players who desire in-game assistance for other reasons, like mobility issues. Conversely, when the objective is immersion, producing more realistic navigation is the more relevant goal.

## Observational Evaluations

We employ a *third-person* Turing test where participants watched videos of the agents navigating. Although the ability to pause, rewind, and replay the videos provided a means of interrogation, it is based solely on observation, and lacks the intervention-based approach of a typical Turing test. Intervention-based approaches could include changing the camera perspective, adversarially interrupting the AI agent's intended path, and more. These forms of interaction may yield different insights.

**Limitations of interactive evaluations.** There are some downsides, however, to deploying a more interactive test, particularly at scale. Recruiting a human evaluator and a human player to interact requires their simultaneous availability for real-time feedback. One solution is in-person studies, which can be challenging to scale and deploy. For instance, at the time of this study, we could not run in-person studies due to the ongoing global pandemic or distribute our proprietary game build to remote participants. Future work could take advantage of advances in game streaming, which may enable interactive remote studies with proprietary game builds. This solution can also incorporate previous work on simultaneous recruitment of participants [23, 337]. However, constructing the architecture to incorporate these different technologies may require significant engineering effort. Importantly, previous work has demonstrated that the inclusion of more direct ways to interrogate the agent by embodying the player and agent in the same virtual space can lead to limited insight [320]. Indeed, work that included an in-game assessment of the human or bot introduced the side effect of an additional game mechanic causing some players to prioritize either gameplay or on the believability assessment [132]. This division of attention yields unreliable results, leading to other researchers adopting third-person variants of these assessments [12, 77, 292].

**A pipeline for evaluating human-like agents.** As a result, we believe that the following pipeline could be useful for evaluating human-like agents. Designers can initially deploy a third-person Turing test to evaluate the human likeness of specific behaviors, which can then be used to design diverse agents that depends on the most common beliefs of the participants. For example, agents could move more smoothly if the participant believes smooth movement to be a feature of human likeness. Players could then choose the characters that they want to interact with in the game, enabling them to tailor the game to their own subjective experience and enjoyment. This approach may offer more reliable insights into the effectiveness of the agent's design without sacrificing the integrity of the assessment process. It could also empower game players by enabling them to exert control over their experience.

## 3.8.2 Toward More General Human-Like Agents

Although the specific agent created for our study may not generalize to different games, this is a common and open challenge in the field of AI [178, 267]. Instead, we offer suggestions for using feedback from designers and players to train human-like agents more efficiently and effectively.

**Human-likeness is contextual, suggesting the need for personalization.** The differences in how more or less skilled human judges characterize human likeness suggest that different people have different interpretations of what constitutes human-like behavior. This supports the idea that the believability of NPCs in games is highly subject to the prior beliefs and expectations of the players. This finding aligns with the fundamental principle of *familiarity* [134] that centers the real-world personal experience and knowledge of the user and implicates the importance of *player-centered* design and customization [312]. Rather than producing monolithic human-like agents, we should strive to understand the beliefs of the player and tailor their experiences accordingly.

**More complex settings introduce additional subjectivity, suggesting a need to involve more groups of people in the agent design process.** When moving to more complex settings, an additional difficulty is introduced. The evaluations of human likeness become even more subjective, varying based on individual differences and cultural factors [199]. This result underscores the importance of involving diverse groups of people in the evaluation of AI agents to obtain a more comprehensive understanding of how people perceive these agents. In the context of games, this could look like utilizing participatory design methods [282] to involve game players in the design of the AI agents themselves. With the consent of the players, we could use techniques in the area of learning from human feedback [367, 57, 140], which provide additional channels for people to communicate what they want from AI agents. With these techniques, players can provide training data to the agents in the form of preferences over paired demonstrations generated by the agent, demonstrations of the desired behavior, and more. This can help to ensure that the AI agents are designed with the needs and preferences of diverse groups in mind.

**Incorporating user feedback could help alleviate manual specification burden of video game designers.** This approach can also be used to help reduce the burden on video game designers: in complex domains, it is often challenging to specify reward signals by hand [60, 169, 179]. In part, this difficulty stems from the complexity of the desired behavior: as we have shown, human-like behavior is multifaceted and necessitates optimizing over multiple objectives. Furthermore, it is sometimes challenging to write down exactly what we mean when specifying a task. For instance, how do we construct a reward signal that captures the task of *build a house in a video game in the same style as surrounding houses?* [291]. When designing a reward signal for this task, we would need to encode what counts as a house, what components are most important to emulate in the style, and which structures count as houses. A person can quickly understand the intention of this instruction, but it is challenging to make explicit this implicit understanding. We will discuss problems of this type in more detail in Chapter 7.

**Easy-to-use tools that enable quick iteration, incorporating user feedback are key to generalizing insights to more game settings.** As a result, an exciting avenue for future work involves developing more effective techniques for learning from people, evaluating user experiences of these techniques, and incorporating them into a flexible, user-friendly tool. This tool can also help extend this work to more general game settings. To more easily enable this line of work, assessments of human likeness could be incorporated into commonly-used game engines, like Unity [151]. This tool would enable game developers to easily evaluate the human likeness of their AI agents using metrics and benchmarks that have been validated in previous research. Additionally, this tool could contain libraries of pretrained *human-like* AI agents, which developers could use as a starting point for their own work. For example, developers could utilize a pretrained human-like navigation agent to perform navigation but develop their own algorithm to use for different tasks. Using this tool could save developers time and effort by enabling them to quickly and easily create more believable and engaging agents to improve the player experience.

**Human-like, tool-use agents have great potential for societal impact.**   As AI agents gain the ability to use tools and operate autonomously online, the potential for misuse grows. This issue is exacerbated when agents can act sufficiently human-like. For example, a scam agent could mimic the tone, timing, and actions of a real customer service representative, convincingly requesting sensitive information or initiating transactions. This issue is particularly problematic at scale. As a result, we will likely need robust mitigation solutions, including but not limited to designing effective user education interventions that promote public awareness of AI deception and studying contexts in which people may be more or less likely to believe AI deception.

# Part II

# Interpretable

# 4  Reducing Human Annotation Burden for Concept-Based Policies

## 4.1 Introduction

A central theme of this thesis is that RL must be approached from a human-centered perspective. The previous chapter establishes how we can train agents that exhibit behavior that is indistinguishable from a person and more deeply understand human perceptions of agent behavior. This work reveals that even in constrained environments, human-likeness is multifaceted: people hold diverse and often conflicting beliefs about what makes behavior appear human. This diversity highlights a key insight for human-centered RL: effective agent design must account for variation in human expectations, especially in interactive or collaborative settings. But aligning behavior with perceived human-likeness is only part of the challenge. People may be able to anticipate an agent's behavior, but still lack insight into the reasons behind its decision.

Due to increasing regulatory demands and broader concerns around transparency, interpretable decision-making has become a central objective. Beyond regulatory compliance, motivations behind interpretability include supporting a range of downstream, practical goals, such as enabling users to build appropriate trust, identifying errors, correcting undesired behavior, and anticipating how agents will respond in novel contexts. While much of the literature on interpretability has focused on decision-support settings, in which the human integrates knowledge from various sources to take an action, RL shifts the decision point to the agent. As a result, this chapter turns to the challenge of interpretability by making the agent's decision-making process understandable to people.

Specifically, we focus on interpretability through *concepts*. To address interpretability concerns in supervised learning, recent works have integrated human-understandable concepts into neural networks through concept bottleneck models [92, 167]. These approaches insert a bottleneck layer, whose units correspond to interpretable concepts, into a neural network. In doing so, the final decisions depend only on these concepts instead of on opaque raw inputs. By training the model both to have high task accuracy and to accurately match experts' concept labels, these models learn a high-level concept-based representation that is simultaneously meaningful to humans and useful for machine learning tasks. As an example, a concept-based explanation for a bird classification task might include a unit that encodes the bird's wing color. Concepts also offer a path to compositional generalization by enabling subtask decomposition into meaningful, reusable components [203, 341].

Figure 4.1: **LICORICE overview.** In concept-based RL, the policy first maps from states to the concepts in a bottleneck layer with $g$, and then maps from concepts to (distributions over) actions with $f$. During training, LICORICE addresses concept label efficiency concerns with three key components: i) iterative training, ii) data decorrelation, and iii) active learning.

More recently, these techniques have been applied to RL by incorporating a concept bottleneck in the policy [110, 356], so that the chosen actions are a function of the human-understandable concepts. As we establish in Section 2.1, in RL, agents learn a *policy*, which is a strategy for making sequential decisions in complex environments. Agents typically represent the policy as a neural network, as this representation tends to yield high performance [225]. However, this choice can come at a cost: such policies are challenging for stakeholders to interpret—particularly when the network inputs are also complex, such as high-dimensional sensor data. This opacity can pose a significant hurdle, especially in applications where understanding the rationale behind decisions is critical, such as healthcare [354] or finance [189]. In such applications, decisions can have significant consequences, so it is essential for stakeholders to fully grasp the reasoning behind actions to confidently adopt or collaborate on a policy.

However, a significant challenge emerges in this method's practical application: past work assumes that concept annotations are readily available during training. To learn a mapping from states to concepts, a RL agent requires concept information for every state it encounters during training, which is often millions or billions of state-action pairs. Because many domains do not automatically provide these concept labels (e.g., autonomous driving), we require a means of obtaining them. Extensively relying on human labelers is impractical, risking errors due to fatigue [204]. Using vision language models (VLMs) for automated concept extraction [237], while alleviating the human bottleneck, introduces significant API or inference costs that can be prohibitive. As a result, a key challenge is achieving sample-efficient and scalable concept learning.

In this work, we address this challenge with LICORICE (**L**abel-efficient **I**nterpretable **CO**ncept-based **R**e**I**nfor**CE**ment learning), a novel training paradigm designed to minimize the number of concept annotation queries while maintaining high task performance. LICORICE achieves state-of-the-art performance in terms of both reward and concept error, all while using substantially fewer concept labels for training. Figure 4.1 illustrates our algorithm. LICORICE includes three main innovations that address important challenges when we move from classification and regression to the RL setting.

First, LICORICE addresses the problem of concept learning on off-policy or outdated data. If concepts are learned from data collected by a random policy, the concept distribution may not reflect the data from an optimal policy. To ensure that concept learning occurs on more recent and on-policy data, LICORICE interleaves concept learning and RL training through *iterative training*: it alternately freezes the network layers corresponding to either the concept learning part or the decision-making part. Second, LICORICE addresses the problem of limited training data diversity that occurs when an agent interacts with the environment, thereby generating sequences of highly similar, temporally correlated data points. To tackle this issue, LICORICE implements a *data decorrelation* strategy to produce a more diverse set of training samples. Third, LICORICE addresses the inefficient use of annotation effort, where labeled samples may provide redundant information. To resolve this problem, we employ *disagreement-based active learning* using a concept ensemble to select the most informative data points for labeling.

To evaluate the effectiveness of LICORICE, we conduct experiments in two scenarios—perfect human annotation and VLM annotation—on five environments with image input: an image-based version of CartPole, two Minigrid environments, and two Atari environments. First, with assumption of perfect human annotation, we show that LICORICE yields both higher concept accuracy and higher reward while requiring fewer annotation queries compared to baseline methods. Second, we find that VLMs can indeed serve as concept annotators for some, but not all, of the environments. Third, we show that LICORICE is more efficient in the number of human concept labels compared with other methods, showcasing its potential to be used in lower data regimes. Next, we demonstrate that LICORICE supports (simulated) interventions by a human annotator at test-time, further supporting its potential for practical use. Finally, we conduct an ablation study and investigate how robust LICORICE is to various hyperparameter values.

**Contributions.** In summary, this chapter contributes the following.

- We are the first to formalize and investigate the problem of a limited concept annotation budget for interpretable RL.

- We introduce LICORICE, a novel training scheme that enables label efficient learning of concept-based RL policies.

- We conduct extensive experiments to show the effectiveness of LICORICE across five environments with varying budget constraints.

- We study the use of VLMs as automated concept annotators, demonstrating their effectiveness in certain environments while highlighting challenges in others.

- We show that LICORICE supports simulated test-time interventions and is robust to various hyperparameter values.

## 4.2 Preliminaries

We now provide a brief overview of concept bottleneck models and their use in the supervised learning setting, followed by their use in the RL setting.

### 4.2.1 Concept Bottleneck Models

Concept-based explanations have emerged as a common paradigm in explainable AI [256]. They explain a model's decision-making process through human-understandable attributes and abstractions. In supervised learning, concept bottleneck models [167] implement this approach using two key functions: a concept predictor $g : X \rightarrow C$, mapping inputs $X$ to interpretable concepts $C$, and an output predictor $f : C \rightarrow Y$, mapping the concept predictions $C$ to a downstream task space $Y$, such as labels for classification.

Some work leverages a so-called *soft* concept model, which means that information flow occurs through an uninterpretable channel; however, we specifically focus on *hard* concept models, which are constrained to make decisions based only on concepts. In this setting, the prediction takes the form $\hat{y} = f(g(x))$, so that the input $x$ influences the output solely through the bottleneck layer $\hat{c} = g(x)$.

### 4.2.2 Concept Policy Models

In RL, since we want the policy to be interpretable, we insert a concept bottleneck layer into the policy network, such that the policy $\pi(s)$ now consists of two functions $f, g$. Specifically, we implement it as,

$$\pi(s) \doteq f(g(s)),$$

so $\pi$ maps from states $s \in \mathcal{S}$ to concepts $c \in C$ through $g$, then to actions $a \in \mathcal{A}$ through $f$ [110, 356]. This setup allows the policy to base its decisions on understandable and meaningful features. As a result, we can use any RL algorithm that can be modified to include an additional loss function for concept prediction.

Training these models requires a dataset of state-concept-action triplets $(s, c, a) \in \mathcal{S} \times C \times \mathcal{A}$. The functions $f$ and $g$ are typically implemented as neural networks, with their parameters collectively defined as $\theta$. Previous work [356] simply combines the concept prediction loss $L^C(\theta)$ and RL loss $L^{\text{RL}}(\theta)$:

$$L(\theta) = L^{\text{RL}}(\theta) + \lambda_C L^C(\theta),$$

where the exact definitions of $L^{\text{RL}}(\theta)$ and $L^C(\theta)$ depend on the choice of RL algorithm and concept learning task. The objective is to find the optimal parameters $\theta^*$ that minimize this combined loss function with the coefficient $\lambda_C$. However, this approach requires continuous access to ground-truth concepts for training $f$, which may not always be feasible or desirable in practical RL scenarios.

## 4.3 LICORICE

The standard way of training concept-based RL agents assumes continuous access to an oracle to provide concept labels. However, this assumption is problematic due to the large annotation cost incurred by human or automated labeling efforts. To reduce the number of concept labels required for concept-based reinforcement learning, we propose LICORICE, a novel algorithm for interpretable RL consisting of three main components: iterative training, data decorrelation, and disagreement-based active learning. The full pseudocode is in Algorithm 1.

### 4.3.1 Iterative Training

As the agent's policy improves, the distribution of the visited states and their associated concepts changes. Consider a Markov decision process where states are indexed. With a random (initial) policy, the agent tends to visit small-index states near the initial state, encountering only the concepts relevant to those states. However, as the policy improves, the agent explores higher-index states, leading to a shift in both the state and concept visitation distribution. If we exhaust all queries at the beginning of training, we risk training the model only on the concepts associated with small-index states from the random policy, potentially missing important concepts that emerge as the agent explores.

We therefore propose iterative training to enable LICORICE to progressively refine its understanding of concepts as the policy improves. Iterative training consists of two parts: concept learning and behavior learning. Assume that we have access to an annotated concept dataset $\mathcal{D}_{\text{train}}$ generated by our current policy $\pi_m$. Then, *concept learning* focuses on training the concept portion of the network $g$ on this dataset $\mathcal{D}_{\text{train}}$ (line 18). *Behavior learning* involves freezing $g$ and training the behavior portion of the network $f$ with any standard RL algorithm with its associated loss $L^{\text{RL}}$ (line 19). Upon completion of behavior learning, we obtain an updated policy $\pi_{m+1}$, which serves to collect more unlabeled concept data to help train $g$ in the next iteration.

### 4.3.2 Data Decorrelation

We do not obtain ground-truth concept labels for all states when rolling out the current policy $\pi_m$. Instead, rollouts produce a dataset $\mathcal{U}_m$ of unlabeled states as candidates for querying for ground-truth concepts from an oracle. Simply collecting and labeling all encountered states would give us long chains of nearly-identical samples, undermining the diversity we need for effective learning.

To resolve this challenge, we introduce data decorrelation with two key components: a budget multiplier $\tau$ that sets $|\mathcal{U}_m| = \tau \cdot B_m$, and a per-state acceptance probability $p$. This random acceptance/rejection mechanism leverages a fundamental property of Markov processes—states become increasingly independent as their temporal distance grows according to the mixing time—allowing us to build more diverse datasets from our policy rollouts.

---

**Algorithm 1** LICORICE (**L**abel-efficient **I**nterpretable **CO**ncept-based **R**e**I**nfor**CE**ment learning)

---

1: **Input:** Total budget $B$, number of iterations $M$, sample acceptance threshold $p$, ratio $\nu$ for active learning, batch size for querying $b$, number of concept models $N$ to ensemble
2: **Initialize:** training set $\mathcal{D}_{\text{train}} \leftarrow \emptyset$, and validation set $\mathcal{D}_{\text{val}} \leftarrow \emptyset$
3: **for** $m = 1$ **to** $M$ **do**
4:      Allocate budget for iteration $m$: $B_m \leftarrow \frac{B}{M}$
5:      **while** $|\mathcal{U}_m| < \nu \cdot B_m$ **do**
6:          Execute policy $\pi_m$ to collect unlabeled data $\mathcal{U}_m$ using acceptance rate $p$
7:      **end while**
8:      Initialize iteration-specific datasets for collecting labeled data: $\mathcal{D}'_{\text{train}} \leftarrow \emptyset, \mathcal{D}'_{\text{val}} \leftarrow \emptyset$
9:      **while** $B_m > 0$ **do**
10:          Train $N$ concept models $\{\tilde{g}_i\}_{i=1}^N$ on $\mathcal{D}_{\text{train}} \cup \mathcal{D}'_{\text{train}}$, using $\mathcal{D}_{\text{val}} \cup \mathcal{D}'_{\text{val}}$ for early stopping
11:          Calculate acquisition function value $\alpha(s)$ for all $s \in \mathcal{U}_m \setminus (\mathcal{D}'_{\text{train}} \cup \mathcal{D}'_{\text{val}})$
12:          Choose $b_m = \min(b, B_m)$ unlabeled points from $\mathcal{U}_m$ according to $\arg\max_s \alpha(s)$
13:          Query for concept labels to obtain new dataset $\mathcal{D}_m \leftarrow \{(s, c)\}^{b_m}$
14:          Split $\mathcal{D}_m$ into train and validation splits and add to $\mathcal{D}'_{\text{train}}, \mathcal{D}'_{\text{val}}$
15:          Decrement budget: $B_m \leftarrow B_m - b_m$
16:      **end while**
17:      Aggregate datasets: $\mathcal{D}_{\text{train}} \leftarrow \mathcal{D}_{\text{train}} \cup \mathcal{D}'_{\text{train}}, \mathcal{D}_{\text{val}} \leftarrow \mathcal{D}_{\text{val}} \cup \mathcal{D}'_{\text{val}}$
18:      Continue training the concept network $g$ on $\mathcal{D}_{\text{train}}$, using $\mathcal{D}_{\text{val}}$ for early stopping
19:      Freeze $g$ and continue training $f$ using standard reinforcement learning training to obtain $\pi_{m+1}$
20: **end for**

---

### 4.3.3 Disagreement-Based Active Learning

Equipped with $\mathcal{U}_m$, we are now prepared to select data points for querying for concept labels. The purpose of this stage is to make use of our limited labeling budget $B_m$ to collect a labeled dataset $\mathcal{D}_m$ for training $g$.

We propose to train a concept ensemble—consisting of $N$ independent concept models—from scratch on the dataset of training points $\mathcal{D}_{\text{train}}$ that has been aggregated over all iterations (line 10). We use this ensemble to calculate the disagreement-based acquisition function $\alpha(\cdot)$, which determines whether each candidate in $\mathcal{U}_m$ ought to receive a ground-truth concept label (line 11). This function targets samples where prediction disagreement is highest among ensemble members, as these points often represent areas of uncertainty or decision boundaries where additional labeled data would be most informative.

For $\alpha(\cdot)$, we use two different formulations, depending on the concept learning task. For concept classification, we use a query-by-committee approach [287], in which we prioritize points with a high proportion of predicted class labels that are not the modal class prediction (also called the variation ratio [19]):

$$\alpha(s) \doteq 1 - \max_{c \in C} \frac{1}{N} \sum_{i=1}^{N} [\tilde{g}_i(s) = c].$$

Here, $\tilde{g}_i(s)$ is the concept prediction of the $i$-th model on state $s$. For concept regression, we directly use the estimate of variance across the concept models as a measure of disagreement:

$$\alpha(s) \doteq \sigma^2(s) = \frac{1}{N-1} \sum_{i=1}^{N} (\tilde{g}_i(s) - \mu(s))^2,$$

where $\mu(s) \doteq \frac{1}{N} \sum_{i=1}^{N} \tilde{g}_i(s)$ is the mean prediction of the $N$ concept models. After querying for $B_m$ ground-truth concept labels (line 13) using $\alpha(\cdot)$, we are prepared to continue training $g$ during the concept learning stage.

### Implementation Details

For training $g$ (line 18), we employ two types of loss functions depending on the nature of the concepts: MSE for continuous concepts and cross-entropy loss for categorical concepts. If the problem requires mixed-type concepts, we either discretize continuous attributes or simply use a mixed loss.

To train $f$ (line 19), we freeze $g$ and use PPO to continue training $f$ from the previous iterations, resulting in the final policy $\pi_{M+1}$. For complex environments that require many iterations, we observe that reinforcement learning-related parameters may get stuck in a local region in early iterations, becoming hard to optimize later. To mitigate this issue, in the last iteration, we additionally train a new reinforcement learning network $f'$ from scratch when fixing $g$ with $\pi_{M+1}$ as the anchoring policy to get the final $\pi'_{M+1}$. The anchoring policy ensures $\pi'_{M+1}$ has a similar distribution to the previous one and the annotated observations are still highly relevant. Specifically, we initialize another policy $\pi_\theta$ with the same frozen $g$ parameters and randomly initialized $f$, then train it with the standard PPO loss function and an additional KL-divergence penalty between the current policy and $\pi_{M+1}$:

$$L(\theta) \doteq L^{\text{PPO}}(\theta) + \beta \cdot \text{KL}[\pi_{M+1}(\cdot \mid s), \pi_\theta(\cdot \mid s)].$$

## 4.4 Experimental Setup

In our experiments, we investigate the following questions:

**RQ 1** Under a limited concept annotation budget, does LICORICE enable both high concept accuracy and high environment reward?

**RQ 2** How effective are VLMs as automated concept annotators when used with LICORICE?

**RQ 3** How label efficient is LICORICE compared with other methods?

**RQ 4** Does LICORICE support test-time concept interventions?

| Environment | # Concepts | Concept Type | Binarized # Concepts |
|---|---|---|---|
| PixelCartPole | 4 | Continuous | N/A |
| DoorKey | 12 | Discrete | 46 |
| DynamicObstacles | 11 | Discrete | 30 |
| Boxing | 8 | Discrete | 1480 |
| Pong | 7 | Discrete | 1848 |

Table 4.1: **Summary of environments, including their associated concepts.** Counts for the binarized version of the concepts provided where applicable to illustrate the problem size.

**RQ 5**  Which components of LICORICE contribute to its performance?

**RQ 6**  How robust is LICORICE to various hyperparameter values and other choices?

In each experiment, we run each algorithm 5 times, each with a random seed. All algorithms use PPO [284, 263] with a concept bottleneck. More implementation details and hyperparameters are in Appendix 2.

**Environments.**  We investigate these questions across five environments, each of which includes a distinct challenge and features a set of interpretable concepts describing key objects and properties. These environments are PixelCartPole [352], DoorKey [55], DynamicObstacles [55], Boxing [18], and Pong [18]. We summarize the concepts in Table 4.1, with more details in Appendix 1. These environments are characterized by their dual representation: a complex image-based input and a symbolic concept-based representation. PixelCartPole, DoorKey, and DynamicObstacles are simpler because we can extract noiseless ground-truth concept labels from their source code. In contrast, Boxing, as implemented in OCAtari [74], uses reverse engineering to extract positions of important objects from the game's RAM state. This extraction process introduces a small amount of noise. The concepts in Pong are derived from the VIPER paper [17] that also uses reverse engineering.

**Baselines.**  To our knowledge, there exist no previous algorithms that seek to minimize the number of concept labels for interpretable RL, so we implement three: Sequential-Q, Disagreement-Q, and Random-Q. In Sequential-Q, the agent spends $B$ queries on the first $B$ states encountered during the initial policy rollout. In Disagreement-Q, the agent also spends its budget at the beginning of training; however, it uses active learning with the same $\alpha(\cdot)$ as LICORICE to strategically choose the training data for annotation. In Random-Q, the agent receives $B$ concept labels at random points using a probability to decide whether to query in each state. As a budget-unconstrained baseline, we implement CPM from previous work in multi-agent reinforcement learning [356] for the single-agent setting. As mentioned in Section 4.2, CPM jointly trains $g$ and $f$, assuming unlimited access to concept labels, so it also represents an upper bound on concept accuracy.

**VLM details.** For our VLM experiments, we use GPT-4o [1]. During training, we query GPT-4o each time LICORICE requires a concept label. We include the prompt for the PixelCartPole environment here, and detail the rest in Appendix 3:

> **Prompt:** Here are the past 4 rendered frames from the CartPole environment. Please use these images to estimate the following values in the latest frame (the last one):
>   - Cart Position, within the range (-2.4, 2.4)
>   - Cart Velocity
>   - Pole Angle, within the range (-0.2095, 0.2095)
>   - Pole Angular Velocity
>
> Additionally, please note that the last action taken was [last action].
>
> Please carefully determine the following values and give concise answers one by one. Make sure to return an estimated value for each parameter, even if the task may look challenging.
>
> Follow the reporting format:
>   - Cart Position: estimated_value
>   - Cart Velocity: estimated_value
>   - Pole Angle: estimated_value
>   - Pole Angular Velocity: estimated_value

**Performance metrics.** We compare our reward to an upper bound calculated by training PPO with ground-truth concept labels (PixelCartPole: 500, DoorKey: 0.97, DynamicObstacles: 0.91, Boxing: 86.05, Pong: 21). In the first four environments, percentages (or ratios) make sense since the minimum reasonable reward is 0.[1] In Pong, a random policy has -21 reward so we first get the difference between rewards and -21 and then calculate the ratio. We additionally report the concept error (MSE for regression; $1 -$ accuracy for classification). In Boxing, following the practice of OCAtari, we consider a concept prediction correct if it is within 5 pixels of the ground truth label. In Pong, we simply use classification error. All reported numbers are calculated using 100 evaluation episodes.

---

[1]In PixelCartPole and DoorKey, all rewards are nonnegative; in DynamicObstacles, an agent can ensure nonnegative reward by simply staying in place; in Boxing, a random policy can get around 0 reward and all of the trained policies have positive reward.

Figure 4.2: **LICORICE generally achieves higher reward and lower concept error than all other budget-constrained algorithms.** The budgets for each environment are PixelCartpole: 500, DoorKey: 300, DynamicObstacles: 300, Boxing: 3000, Pong: 5000. CPM has an unlimited budget (in practice, it uses 4M, 4M, 1M, 15M, 10M labels, respectively). Despite this extreme budget difference, LICORICE achieves comparable reward and concept error. Full numerical results in Table 5, Appendix 5.

## 4.5  Results

We now present the results, answering **RQ 1 - 6** introduced in Section 5.6.

### RQ 1: Reward & Concept Accuracy Under a Limited Annotation Budget

We first seek to understand how LICORICE performs compared to the state-of-the-art in concept-based RL, CPM, which is not constrained by concept label budgets. We then compare LICORICE in a more fair manner, against budget-constrained algorithms that we describe in Section 4.4.

**LICORICE achieves similarly high reward and concept performance to the state-of-the-art budget-unconstrained baseline.** Surprisingly, LICORICE and CPM achieve nearly 100% of the maximum reward in all environments in this unfair comparison (Figure 4.2). They also achieve a similar concept error rate, only seeing a clear difference in error for the most complex environment, Pong. We emphasize that CPM is given an unlimited budget; in practice, it uses 1M-15M concept labels for all environments, which is around 2000 to $13000\times$ the budget used by LICORICE. In contrast, LICORICE operates under a strict budget constraint, with 3000 labels for Boxing and 5000 for Pong, and at most 500 concept labels for the simpler environments (PixelCartPole: 500, DoorKey: 300, DynamicObstacles: 300).

|  | $c$ **Error** | |
| **Environment** | LICORICE+GPT-4o | GPT-4o |
| --- | --- | --- |
| PixelCartPole | 0.25 | 0.24 |
| DoorKey | 0.31 | 0.30 |
| DynamicObstacles | 0.13 | 0.13 |
| Boxing | 0.98 | 0.82 |
| Pong | 0.73 | 0.87 |

Table 4.2: **Concept error comparison.** The concept error of LICORICE+GPT-4o matches that of GPT-4o alone in all environments except Boxing and Pong. This result shows that the noise from $f$ trained from VLM annotations largely stems from the VLM annotations themselves. For Boxing, the error is higher, meaning that the model is likely underfitting.

**LICORICE generally achieves higher reward and lower concept error than budget-constrained baselines.** We study all budget-constrained algorithms (LICORICE, Disagreeemnt-Q, Sequential-Q, Random-Q) under the same fixed concept labeling budget as above. Figure 4.2 shows that LICORICE outperforms all baselines in terms of reward and concept error on all but arguably the easiest environment, DynamicObstacles. In that environment, LICORICE performs similarly to the baselines. The greatest performance differences occur in PixelCartPole and DoorKey, where LICORICE achieves 100% and 99% of the maximum reward, while the second-best algorithm achieves 36% and 89%, respectively. We therefore answer **RQ 1** affirmatively: not only does LICORICE achieve the same high concept accuracy and high reward as the state-of-the-art in concept-based reinforcement learning, it also performs on-par to budget-constrained baselines in one environment and outperforms them in the other four.

## RQ 2: VLMs as Concept Annotators

We now seek to answer the question of whether VLMs can successfully provide concept labels in lieu of a human annotator within our LICORICE framework. Because using VLMs incurs costs, we operate within the same budget-constrained setting as above. We present these results in Figure 4.2.

**VLMs can serve as concept annotators for some environments.** In DoorKey and DynamicObstacles, LICORICE with GPT-4o labels achieves 83% and 88% of the maximum reward, respectively. To assess the quality of VLM-generated labels, we compare the concept error rate of our trained model against GPT-4o's labeling error, both evaluated on the same rollout observations, while GPT-4o evaluation is on a random subset of 50 observations, due to API cost constraints. Table 4.2 shows that the concept error of LICORICE is comparable to that of GPT-4o when GPT-4o labeling error is not high, which suggests that LICORICE can effectively utilize these concept labels. VLMs, particularly GPT-4o, could serve as concept annotators for certain concepts in certain environments.

| Concept Name | Value Range | GPT-4o Error |
|---|---|---|
| Agent Position x | 5 | $0.36 \pm 0.08$ |
| Agent Position y | 5 | $0.41 \pm 0.12$ |
| Agent Direction | 4 | $0.28 \pm 0.06$ |
| Key Position x | 6 | $0.20 \pm 0.03$ |
| Key Position y | 6 | $0.32 \pm 0.11$ |
| Door Position x | 5 | $0.37 \pm 0.07$ |
| Door Position y | 5 | $0.32 \pm 0.04$ |
| Door Open | 2 | $0.38 \pm 0.09$ |
| Direction Movable Right | 2 | $0.30 \pm 0.03$ |
| Direction Movable Down | 2 | $0.19 \pm 0.06$ |
| Direction Movable Left | 2 | $0.33 \pm 0.06$ |
| Direction Movable Up | 2 | $0.20 \pm 0.08$ |

Table 4.3: **Concepts and GPT-4o labeling errors in DoorKey.** All concepts are discrete. Error is 1 minus classification accuracy. The $\pm$ [value] part shows the standard deviation. GPT-4o struggles with more complex concepts, like whether the agent can move to the left (Direction Movable Left).

**VLMs struggle to provide accurate concepts in more complex environments.** In PixelCartPole, Boxing, and Pong environments, however, GPT-4o struggles to provide accurate labels. In these environments, not only does LICORICE+GPT-4o achieve less than 25% of the maximum reward, it also incurs a large concept error. The challenge with PixelCartPole is the continuous nature of the concepts and the lack of necessary knowledge of physical rules and quantities in this specific environment. Boxing is particularly challenging due to the large number of possible concept values that each of the 8 discrete concepts can take on (160 or 210 possible values). Similarly, concepts in Pong have hundreds of possible values and some physics quantities require multi-frame inference.

**GPT-4o struggles with more complex concepts, even in simpler domains.** We seek to more deeply understand the challenges that VLMs face in providing accurate concept labels. As a result, in Table 4.3, we provide more details regarding the concepts used in DoorKey, one of the simpler environments. The remaining environments are similarly decomposed in Appendix 1. We investigate the per-concept error induced by GPT-4o. We observe that GPT-4o performs best on simple positional concepts like the key's x-coordinate (0.20 error) and struggles more with relational or action-relevant concepts such as whether the door is open (0.38) or whether movement is possible in a given direction (e.g., movable left: 0.33). This within-task variation suggests that even when all concepts are discrete and visually grounded in the same scene, GPT-4o is more accurate on static object properties than on relational or affordance-based concepts that require integrating multiple visual elements. We therefore answer **RQ 2** with cautious optimism: there are indeed cases in which LICORICE could be used alongside VLMs, though not all yet. Generally speaking, a less complex environment is more likely to work well with VLMs by enabling clear and detailed labeling instructions in the prompt. However, VLM capabilities need improvement before they can fully alleviate the human annotation burden, especially in safety-critical settings.

Figure 4.3: **LICORICE more effectively makes use of varying budgets.** Reward (top) and concept error (bottom) of all algorithms for different budgets. LICORICE achieves higher reward and lower concept error. In other words, LICORICE more efficiently leverages fewer human concept labels to achieve higher reward and lower concept error. Full numerical results are in Tables 6 to 8, Appendix 5.

## RQ 3: Investigation of Budget Requirements for Performance

Given these results, we now investigate the minimum budget required for the agents to achieve 99% of the reward upper bound. To do so, we incrementally increase $B$ for each environment starting from a reasonably small budget until LICORICE reaches 99% of the reward upper bound. In Pong, we further increase $B$ until the concept error is below $0.25$. In addition to the baselines, we study LICORICE under perfect (human) annotation and noisy VLM labels in Figure 4.3.

**LICORICE more rapidly achieves high reward compared to baselines.** Across all environments, LICORICE is the most query-efficient. In three of the environments, it consistently achieves higher reward than baselines across all budget levels. In DynamicObstacles and Pong, LICORICE outperforms baselines at the smallest query budget, after which some baselines achieve comparable reward. Furthermore, except for one budget setting for one environment, LICORICE consistently achieves lower concept error than the baselines.

Figure 4.4: **Concept intervention results: LICORICE enables test-time concept intervention.** If we use a noisy concept model $\hat{g}$ and intervene on the concepts, then the performance of LICORICE steadily increases. In practice, this result indicates that LICORICE would support this form of test-time human-AI interaction.

**For LICORICE+GPT-4o, more concept labels is not always better.** In DoorKey, LICORICE+GPT-4o exhibits predictable behavior in terms of concept error: as the budget increases, the reward increases and the concept error decreases. Counterintuitively, for some environments, the concept error and reward fluctuate with more budget, likely due to the additional labeling noise introduced. We therefore provide a more nuanced answer to **RQ 3**: LICORICE is indeed label efficient across varying budgets, but the benefit is annotator-dependent.

## RQ 4: Interpretability Analysis: Test-Time Concept Interventions

A nice property of concept-based networks is the ability for people to successfully intervene on the concepts to correct them. In RL, this intervention enables the inspection of how different concepts influence the immediate decisions of the agent. We investigate not only whether LICORICE supports test-time concept intervention but also qualitatively showcase how domain experts could interact with the concept model for human-AI decision-making.

**LICORICE enables test-time concept intervention.** To validate that LICORICE supports test-time concept intervention, we simulate using a noisy concept model $\hat{g}$ with $f$ to study the impact of concept interventions on reward. For PixelCartPole, we add Gaussian noise to each predicted concept with a standard deviation of $0.2$. For the other three environments with discrete concepts, each concept label is randomly changed with probability $0.2$. Following previous work [167], we first intervene individually on concepts, setting them to ground truth, and sort the concepts in descending order of reward improvement. We then sequentially intervene on the concepts following this ordering, such that any previously intervened concepts remain intervened. According to Figure 4.4, using a noisy concept model significantly degrades reward. However, the performance increases with more interventions, meaning LICORICE supports test-time concept intervention, affirmatively answering **RQ 4**.

Figure 4.5: **Test-time intervention examples, where intervening on a single concept corrects the action.** Both sides show example interventions on the door open concept. On the left-hand side, it is first incorrectly predicted as being open, leading to the wrong action of moving forward. The correction to the door being closed changes the action to toggle, which is the optimal action to open the door. On the right-hand side, it is incorrectly predicted as being closed, so the agent attempts to toggle it open. This action is incorrect because it will close the door! The intervention changes the action to moving forward, which is the optimal action to move into the doorway toward the goal.

**Concept interventions show how the policy decisions change.** We now show how domain experts could interact with the model by intervening on the model at test-time. Specifically, we simulate a counterfactual question: *What if a concept value is incorrectly predicted?* Figure 4.5 depicts two examples in DoorKey. Here, the RL agent incorrectly predicts an important concept for decision-making, which is then adjusted by the human at test-time. Just intervening on the concept of the door being open or closed is sufficient to change the agent's behavior, both highlighting the importance of this specific concept and the ability of a user to intervene on the concepts to interrogate and understand agent behavior.

## RQ5: Ablation Study

We now conduct an ablation study of our three main contributions: iterative training, decorrelation, and disagreement-based active learning. LICORICE-IT uses only one iteration, LICORICE-DE does not perform data decorrelation, and LICORICE-AC uses the entire unlabeled dataset for querying instead of a subset chosen by active learning.

**All components of LICORICE contribute positively to its performance.** All components are critical to achieving both high reward and low concept error (Table 4.4; learning curves in Appendix 5). However, the component that most contributes to the performance differs depending on the environment. For example, compared with LICORICE, LICORICE-IT exhibits the largest reward gap for PixelCartPole; however, LICORICE-AC yields the largest reward gap for DynamicObstacles, and LICORICE-DE yields the largest gap for DoorKey. This difference may occur because the concepts in DynamicObstacles are simple enough that one iteration is sufficient, so the largest gains can be made with active learning. In contrast, PixelCartPole requires further policy refinement to better estimate on-distribution concept values, so the largest gains can be made by leveraging multiple iterations.

| Environment | Algorithm | R ↑ | c Error ↓ |
|---|---|---|---|
| PixelCartPole | LICORICE-IT | $0.53 \pm 0.26$ | $0.08 \pm 0.03$ |
| | LICORICE-DE | $0.97 \pm 0.05$ | $0.03 \pm 0.01$ |
| | LICORICE-AC | $0.91 \pm 0.10$ | $0.02 \pm 0.01$ |
| | LICORICE | $1.00 \pm 0.00$ | $0.02 \pm 0.00$ |
| DoorKey | LICORICE-IT | $0.99 \pm 0.01$ | $0.09 \pm 0.02$ |
| | LICORICE-DE | $0.78 \pm 0.08$ | $0.20 \pm 0.05$ |
| | LICORICE-AC | $0.82 \pm 0.10$ | $0.16 \pm 0.06$ |
| | LICORICE | $0.99 \pm 0.01$ | $0.05 \pm 0.01$ |
| DynamicObstacles | LICORICE-IT | $1.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | LICORICE-DE | $0.98 \pm 0.01$ | $0.00 \pm 0.00$ |
| | LICORICE-AC | $0.58 \pm 0.53$ | $0.00 \pm 0.00$ |
| | LICORICE | $0.99 \pm 0.00$ | $0.00 \pm 0.00$ |
| Boxing | LICORICE-IT | $0.98 \pm 0.04$ | $0.14 \pm 0.07$ |
| | LICORICE-DE | $1.00 \pm 0.04$ | $0.04 \pm 0.02$ |
| | LICORICE-AC | $0.99 \pm 0.01$ | $0.09 \pm 0.06$ |
| | LICORICE | $1.00 \pm 0.02$ | $0.05 \pm 0.01$ |

Table 4.4: **Ablation study results for LICORICE.** All components generally help achieve better reward and lower concept error. The benefit of each component is environment-specific, also highlighting the diversity of environments used in our work. For example, active learning is most useful for achieving high reward in DynamicObstacles, whereas data decorrelation is most useful for achieving high reward and low concept error in DoorKey.

**Including the data decorrelation module positively influences data diversity.** We hypothesize that data decorrelation enables better performance by increasing the diversity of the training data. As a result, we investigate the role of data decorrelation on state diversity by calculating the average distance to the $k$ nearest samples in the concept label space for all labeled samples. Intuitively, a larger average distance to nearest neighbors indicates samples are more spread out in the space. Let $D = \{(s_1, c_1), \ldots, (s_b, c_b)\}$ be the labeled dataset and $\text{Neighbor}(c_i, k)$ be the index set of $k$ nearest neighbors with $k = 10$ for all environments, then we calculate diversity as:

$$d = \frac{1}{bk} \sum_{i=1}^{b} \sum_{j \in \text{Neighbor}(c_i, k)} ||c_i - c_j||.$$

We present $d$ along with the reward in Table 4.5 for LICORICE and LICORICE-DE. As hypothesized, decorrelation increases diversity in 3 of the 4 environments. LICORICE achieves higher reward than LICORICE-DE. In DynamicObstacles, the state diversity and reward is similar for both LICORICE and LICORICE-DE. We believe that, in this environment, decorrelation is less critical because neighboring states are inherently more dissimilar due to the dynamic nature of the obstacles.

| Environment | LICORICE | | LICORICE-DE | |
|---|---|---|---|---|
|  | $d \uparrow$ | $R \uparrow$ | $d \uparrow$ | $R \uparrow$ |
| PixelCartPole | $\mathbf{0.12 \pm 0.00}$ | $1.00 \pm 0.00$ | $0.10 \pm 0.01$ | $0.97 \pm 0.05$ |
| DoorKey | $\mathbf{1.12 \pm 0.18}$ | $0.99 \pm 0.01$ | $0.25 \pm 0.00$ | $0.78 \pm 0.08$ |
| DynamicObstacles | $\mathbf{0.93 \pm 0.04}$ | $0.99 \pm 0.00$ | $\mathbf{0.99 \pm 0.09}$ | $0.98 \pm 0.01$ |
| Boxing | $\mathbf{7.20 \pm 0.20}$ | $0.97 \pm 0.04$ | $6.53 \pm 0.08$ | $0.97 \pm 0.01$ |

Table 4.5: **LICORICE's data decorrelation module positively influences data diversity.** Training data diversity ($d$) and reward ($R$) with and without data decorrelation. Higher $d$ is better (more diverse). For each environment, the highest $d$ is bolded. When the standard deviations overlap, we bold both. In all environments except DynamicObstacles, LICORICE with decorrelation has higher concept diversity than LICORICE-DE. For DynamicObstacles, the data diversity is similar.

| Strategy | 100 | | 200 | | 300 | |
|---|---|---|---|---|---|---|
|  | R $\uparrow$ | c Error $\downarrow$ | R $\uparrow$ | c Error $\downarrow$ | R $\uparrow$ | c Error $\downarrow$ |
| Disagreement | 0.96 | 0.05 | 0.97 | 0.01 | 1.00 | 0.00 |
| Entropy | 0.95 | 0.07 | 0.98 | 0.02 | 0.99 | 0.00 |

Table 4.6: **LICORICE is robust to the choice of acquisition strategy for active learning.** We compare disagreement- and entropy-based strategies for DynamicObstacles under different labeling budgets ($B = 100, 200, 300$). There is essentially no difference between the two strategies.

## RQ6: Robustness

LICORICE includes a few hyperparameters, described in Section 4.3 and summarized here for convenience. In data decorrelation, $\tau$ governs the size of the unlabeled dataset collected by the current policy, and $p$ controls the rate of acceptance of data points. Active learning involves a concept ensemble with $N$ models. Here, we first study the effect of varying the values of these hyperparameters, then study the effect of the specific acquisition strategy for active learning.

**LICORICE is robust to various hyperparameter values.** We study the effect of varying the values of $\tau$, $p$, and $N$ on DynamicObstacles. We try 2 values for $N$, 3 values for $p$, and 3 values for $\tau$. Figure 4.6 shows the results. Regardless of the values, LICORICE achieves $96 - 99\%$ of the maximum reward. We notice that having a larger committee ($N = 5$ vs. $N = 3$) generally tends to increase performance. This result shows that LICORICE is reasonably robust to hyperparameter values, meaning it could be more feasible for practitioners to train without worrying about severe performance degradation.

Figure 4.6: **Hyperparameter sensitivity results for LICORICE on DynamicObstacles.** Error bars represent standard deviation. Here, we group by the same $p$ values (acceptance probability). Bars of the same color share the same $\tau$ value (budget multiplier), and bars with the same pattern share the same value for $N$ (committee size). Generally, $p = 0.1$ performs better overall with greater stability. Having a larger committee ($N = 5$ vs. $N = 3$) also tends to increase performance.

**LICORICE is robust to the choice of acquisition strategy for active learning.** We now investigate the impact of the choice of acquisition strategy on both reward and concept accuracy. We implement an entropy-based disagreement measurement for the committee and study its effects on DynamicObstacles. Specifically, we implement the following vote entropy acquisition function:

$$\alpha(s_t) = -\frac{1}{K} \sum_{k=1}^{K} \sum_{j=1}^{P_k} \frac{V_{tjk}}{N} \log \frac{V_{tjk}}{N},$$

where $s_t$ is the input state, $K$ is the number of concepts we learn, $P_k$ is the number of classes for the $k$-th concept, $V_{tjk}$ is the number of votes for the $j$-th class of the $k$-th concept for state $s_t$, and $N$ is the number of committee members. Table 4.6 shows these results. Here, *disagreement* refers to the default strategy we use (see 4.3 for details). There is essentially no difference between the two strategies in terms of both reward and concept error.

## 4.6 Related Work

**Interpretable RL.** Interpretable RL has gained significant attention in recent years [102]. As we discuss in Section 2.2, there exist a variety of methods to explain the decision-making of RL agents. One prominent area uses rule-based methods—such as decision trees [298, 321], logic [75], and programs [334, 253]—to represent policies. These works either assume that the state is already interpretable or that the policy is pre-specified. Unlike prior work, our method involves *learning* the interpretable representation (through concept training) for policy learning.

**Concept learning for RL.** Due to successes in the supervised setting [63, 296, 360], concept models have recently been used in RL. Previous work [71] trains a joint embedding model between state-action pairs and concept-based explanations to expedite learning via reward shaping. Unlike LICORICE, their policy is not a strict function of the concepts, allowing our techniques to be combined to provide concept-based explanations alongside an interpretable policy. Another example, CPM [356], is a multi-agent RL approach that assumes that concept labels are continuously available during training. As we have shown, this approach uses over 1M concept labels, whereas LICORICE achieves similar performance while using 2000× fewer labels.

**Learning concepts with human feedback.** Prior research has explored leveraging human concept labels, but not for RL and without focusing on reducing the labeling burden. For instance, [172] has users label additional information about the relevance of certain feature dimensions to the concept label to facilitate concept learning for high-dimensional data. In a different vein, other work [45] develops a test-time intervention policy for querying for concept labels to improve the final prediction. These methods do not directly tackle the issue of reducing the overall labeling burden during training and could be used alongside our method to further support the goal of improving final performance. A concurrent work [293] also studies sample efficiency but mainly focus on human task demonstrations, which differs from our focus of human concept labels.

## 4.7 Discussion and Future Directions

In this work, we propose LICORICE, a novel RL algorithm that addresses the critical issue of model interpretability while minimizing reliance on continuous human annotation. We introduce a training scheme that enables RL algorithms to learn concepts efficiently from little labeled concept data. Our approach interleaves concept learning and RL training, uses an ensemble-based active learning technique to select informative data points for labeling, and uses a simple sampling strategy to decorrelate the concept data. We conduct initial experiments to demonstrate how we can leverage powerful VLMs to infer concepts from raw visual inputs without explicit labels in some environments. There are broader societal impacts of our work that must be considered, as well as several exciting opportunities for future work.

**Using VLMs as concept annotators.** Through our VLM study, we find that VLMs struggle with certain types of concepts, especially continuous variables. This limitation can impact the overall performance of concept-based models, especially in domains where continuous data is prevalent. Furthermore, VLMs still experience hallucination issues [4] and other failure cases, such as providing inaccurate counts. We believe that future work that seeks to improve general VLM capabilities and mitigate hallucinations would also help overcome this limitation. Addressing this issue could involve developing specialized techniques or using existing tools and libraries to better complement VLM capabilities. More broadly, while VLMs have the potential to significantly improve the efficiency and scalability of labeling processes, this automated labeling may risk perpetuating biases or exposing private data.

**Scaling up to even more complex settings.**  Furthermore, scaling challenges may arise from environmental complexity, such as when only a subset of given concepts are relevant or when learning certain concepts is prerequisite to gathering training data for others. Future research could explore using attention mechanisms [332] or sparse coding [240] to identify relevant concepts. Another exciting direction is the work in human-AI complementarity and learning-to-defer algorithms [231, 335] to train an additional classifier for deferring labeling to a person when the chance of an error is high. In our work, we assume a known set of concepts for the environment; future work could investigate the use of human-VLM teams to determine a reasonable concept set for the environment.

**Designing appropriate concept representations.**  Finally, designing a concept-based representation for RL remains an open challenge. Our work provides a few illustrative examples, but the exact design of these representations can significantly impact performance, often for reasons that are not entirely clear—especially when using VLMs as annotators. Prior work [71] proposed some desiderata for concepts in RL, but future work could refine these principles, especially in the face of VLM annotators. Future work could also include systematically investigating the factors that influence the effectiveness of different concept-based representations through extensive empirical studies, theoretical analyses, and the development of new design principles that guide the creation of effective concept representations. One important factor is communicating uncertainty or incorrect results to foster appropriate trust and reliance, rather than allowing stakeholders to be misled into trusting flawed decisions due to the perceived transparency of the model [159].

# 5 Extracting Interpretable Policies for Multi-Agent Coordination

## 5.1 Introduction

The previous chapter demonstrates how concept-based representations can make RL policies more interpretable by exposing the high-level factors influencing agent behavior. We show how we can do so while limiting the number of human queries necessary to train such policies. However, while concepts help illuminate what attributes the agent uses in decision-making, they do not fully address how the agent decides to act. We can think of concepts as a means to produce human-understandable features for agent decision-making, but the mapping from concepts to actions may require a neural network to achieve high reward. As a result, in this chapter, we contribute new algorithms for extracting interpretable and high-reward policies.

In particular, we focus on the multi-agent setting. In many domains of interest, teams of agents must collaborate to accomplish shared objectives while simultaneously defending against adversarial threats. As a result, multi-agent RL has emerged as a promising technique for solving these challenging problems, such as air traffic control [31], train scheduling [227], cyber defense [202], and autonomous driving [25]. In all of these settings, successful deployment depends on the ability of agents to coordinate effectively. However, despite growing interest in high-stakes applications, interpretability in multi-agent RL remains relatively understudied.

As we discuss in detail in Section 2.2, decision trees [261] offer a promising solution for representing interpretable policies. These models are generally considered to be an *intrinsically* interpretable model family [228]: sufficiently small trees can be contemplated by a person at once (simulatability), have subparts that can be intuitively explained (decomposability), and are verifiable (algorithmic transparency) [187]. In the RL setting, decision trees have been successfully used in the single-agent setting to model transition functions [309], reward functions [72], value functions [260, 324], and policies [207, 259, 271]. We show an example of such a policy in Figure 5.1. However, learning decision-tree policies that support coordinated behavior across multiple agents remains an open challenge. The coordination requirement introduces additional complexity, as each agent must consider the decisions of other agents during training but operate independently during execution.

Figure 5.1: **A decision tree policy of depth two trained with MAVIPER in the Cooperative Navigation environment.** The learned decision-tree policy successfully captures the expert's behavior of navigating to one of the landmarks. The inner nodes represent tests of the features in the environment that influence the navigation path of the agent.

We address this gap with two new algorithms. To design these algorithms, we observe that uninterpretable policies trained with multi-agent RL serve as a useful source of information about how an interpretable agent should act. We therefore propose IVIPER and MAVIPER, which combine ideas from model compression and imitation learning to learn decision-tree policies in the multi-agent setting. These are the first set of algorithms for training intrinsically interpretable decision-tree policies for multi-agent settings. IVIPER and MAVIPER work with most existing neural-network-based multi-agent RL algorithms: the policies generated by these algorithms serve as "expert policies" and guide the training of a set of decision-tree policies. Both algorithms build on VIPER [17], which extracts decision-tree policies for single-agent RL.

IVIPER builds on VIPER to accommodate multiple agents by independently training one decision-tree policy per agent. While this approach preserves individual performance, it overlooks the interdependencies that arise in coordinated tasks, which can result in poor performance when agents act jointly. For example, independently-trained agents with limited policy capacity may not prioritize the states where coordination matters most, leading to (potentially compounding) errors. MAVIPER addresses this limitation by introducing a novel centralized training process, in which the agents grow their trees jointly, layer-by-layer. To handle the partially completed trees, each agent anticipates the likely tree completion for the other agents. In addition, MAVIPER uses a novel resampling strategy of the training data to prioritize states in which coordination is particularly important.

This work bridges the gap between interpretability and coordination in multi-agent RL. By extending decision-tree policy extraction to the multi-agent setting, we introduce methods that allow each agent to act transparently while maintaining high team performance. Our proposed MAVIPER algorithm introduces a joint training scheme that enables decision-tree policies to model and anticipate teammates' behavior, overcoming the limitations of independently-learned trees. Together, these methods demonstrate that interpretable and coordinated multi-agent policies are not mutually exclusive.

**Contributions.**   In summary, this chapter contributes the following:

- We introduce IVIPER and MAVIPER, the first set of algorithms for learning interpretable multi-agent policies. MAVIPER involves a novel joint tree-growing procedure and resampling strategy that enables agents to produce coordinated behavior while preserving interpretability.

- Through experiments on three different multi-agent environments, we show that MAVIPER achieves better coordinated performance than IVIPER and the baselines.

- We demonstrate that MAVIPER produces teams that are more robust to a variety of adversaries, demonstrating its applicability to real-world settings that require robustness.

## 5.2  Background and Preliminaries

We focus on the problem of learning interpretable decision-tree policies in the multi-agent setting. We first describe the formalism of our multi-agent setting, then review the single-agent version of VIPER.

### 5.2.1  Markov Games and Multi-Agent Reinforcement Learning Algorithms

In multi-agent RL, agents act in an environment defined by a Markov game [188, 294]. A Markov game is a generalization of an MDP (see Section 2.1.1) that models $N$ interacting in a shared environment. It consists of a set of states $\mathcal{S}$ describing all possible configurations for all agents, the initial state distribution $\rho : \mathcal{S} \to [0, 1]$, and the set of actions $\{\mathcal{A}_i\}_{i=1}^N$ and observations $\{\mathcal{O}_i\}_{i=1}^N$ for each agent $i \in [N]$.

Each agent aims to maximize its own total expected return $R_i = \sum_{t=0}^{\infty} \gamma^t r_i^t$, where $\gamma$ is the discount factor that weights the relative importance of future rewards. To do so, each agent selects actions using a policy $\pi_{\theta_i} : \mathcal{O}_i \to \mathcal{A}_i$. After the agents simultaneously execute their actions **a** in the environment, the environment produces the next state according to the state transition function $P : \mathcal{S} \times \mathcal{A}_1 \times ... \times \mathcal{A}_N \to \mathcal{S}$. Each agent $i$ receives reward according to a reward function $r_i : \mathcal{S} \times \mathcal{A}_i \to \mathbb{R}$ and a private observation, consisting of a vector of *features*, correlated with the state $o_i : \mathcal{S} \to \mathcal{O}_i$.

Given a policy profile $\pi = (\pi_1, ..., \pi_N)$, agent $i$'s value function is defined as:

$$V_i^{\pi}(s) = r_i + \gamma \sum_{s' \in \mathcal{S}} P(s, \pi_1(o_1), ..., \pi_N(o_N), s') V_i^{\pi}(s'),$$

and state-action value function is:

$$Q_i^{\pi}(s, a_1, ..., a_N) = r_i + \gamma \sum_{s' \in S} P(s, a_1, ..., a_N, s') V_i^{\pi}(s'). \tag{5.1}$$

We refer to a policy profile excluding agent $i$ as $\pi_{-i}$.

---

**Algorithm 2** VIPER for Single-Agent Setting

---

**Input**: $(S, A, P, R), \pi^*, Q^*, K, M$
**Output**: $\hat{\pi}$

1: Initialize dataset $\mathcal{D} \leftarrow \emptyset$
2: Initialize policy $\hat{\pi}^0 \leftarrow \pi^*$
3: **for** $m = 1$ to $M$ **do**
4:    Sample $K$ trajectories: $\mathcal{D}^m \leftarrow \{(s, \pi^*(s)) \sim d^{\hat{\pi}^{m-1}}\}$
5:    Aggregate dataset $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}^m$
6:    Resample dataset according to loss: $\mathcal{D}' \leftarrow \{(s, a) \sim p((s, a)) \propto \tilde{l}(s)\mathbb{I}[(s, a) \in \mathcal{D}]\}$
7:    Train decision tree $\hat{\pi}^m \leftarrow \text{TrainDecisionTree}(\mathcal{D}')$
8: **end for**
9: **return** Best policy $\hat{\pi} \in \{\hat{\pi}^1, ..., \hat{\pi}^M\}$ on cross validation

---

Multi-agent RL algorithms fall into two categories: value-based [266, 303, 310] and actor-critic/policy gradient [95, 184, 193, 355]. Value-based methods often approximate $Q$-functions for individual agents in the form of $Q_i^{\pi}(o_i, a_i)$ and derive the policies $\pi_i$ by taking actions with the maximum Q-values. In contrast, actor-critic methods often follow the centralized training and decentralized execution (CTDE) paradigm [239]. They train agents in a centralized manner, enabling agents to leverage information beyond their private observation during training; however, agents must behave in a decentralized manner during execution. Each agent $i$ uses a centralized critic network $Q_i^{\pi}$, which takes as input some state information $x$ (including the observations of all agents) and the actions of all agents. This assumption addresses the stationarity issue in multi-agent RL training: without access to the actions of other agents, the environment appears non-stationary from the perspective of any one agent. Each agent $i$ also has a policy network $\pi_i$ that takes as input its observation $o_i$. In this work, we aim to follow the decentralized execution paradigm to make our method more applicable for real-world multi-agent deployment.

## 5.2.2  VIPER: Single-Agent Decision-Tree Policy Extraction

VIPER [17] is a popular imitation learning algorithm [54, 197, 209] designed to extract an interpretable decision-tree policy from a high-performing, but opaque, expert policy. VIPER extends DAgger [270] (details in Section 2.1.4) by incorporating the value function of the expert to prioritize more important states for decision-making. Importantly, VIPER is flexible enough to be used along with a high-performing deep neural network trained using any single-agent RL algorithm.

More specifically, VIPER operates over $M$ iterations. In each iteration $m$, VIPER trains a decision-tree policy $\hat{\pi}^m$; the final output is the best policy among all iterations $\hat{\pi} \in \{\hat{\pi}_1, \ldots, \hat{\pi}_M\}$. In iteration $m$, the algorithm samples $K$ trajectories $\tau_m = \{(s, \hat{\pi}^{m-1}(s)) \sim d^{\hat{\pi}^{m-1}}\}$ following the decision-tree policy trained at the previous iteration. The high-level idea is that we should collect data according to the current student's capabilities, motivating the use of $\hat{\pi}^m$ for rollouts instead of the expert policy $\pi^*$.

Then, it uses the expert policy $\pi^*$ to prescribe actions for each visited state, leading to the dataset $\mathcal{D}^m = \{(s, \pi^*(s)) \sim d^{\hat{\pi}^{m-1}}\}$ (Line 4, Algorithm 2). The idea is to use the expert to provide *corrections* to states that are on-distribution for the student. VIPER adds these relabeled experiences to a dataset $\mathcal{D}$ consisting of experiences from all previous iterations.

Let $V^{\pi^*}$ and $Q^{\pi^*}$ be the state value function and state-action value function given the expert policy $\pi^*$. Then, the quantity

$$\tilde{l}(s) = V^{\pi^*}(s) - \min_{a \in A} Q^{\pi^*}(s, a)$$

measures the cost of making a poor decision in state $s$. In other words, it measures the difference between the value of following the expert and taking the worst possible action, quantifying how critical it is to choose the correct action at state $s$. VIPER resamples points $(s, a) \in \mathcal{D}$ proportional to Section 5.2.2, focusing the capacity of the tree on regions of the state space where mistakes would be more costly This procedure results in a new, resampled dataset $\mathcal{D}'_m$.

VIPER then trains a new decision-tree policy using $\mathcal{D}'_m$ to prepare for the next round $m + 1$. Importantly, the decision-tree training uses any standard decision-tree training algorithm, such as CART [30]. Training a decision-tree with CART on the dataset $\mathcal{D}'$ constructed with this resampling procedure is in expectation the same as training a decision-tree with the loss $\tilde{l}(s, \pi)$. Algorithm 2 shows the full VIPER algorithm. Note that it is designed only for a single agent, meaning that it cannot readily handle coordinating multiple interpretable policies.

## 5.3 Approach

To address the challenges of learning interpretable policies in multi-agent RL setting, we introduce two algorithms: IVIPER and MAVIPER. These methods build upon the VIPER algorithm [17], extending it to the multi-agent domain. The key idea behind both algorithms is to extract decision-tree policies from high-performing, potentially uninterpretable expert policies, enabling transparency and verifiability in complex multi-agent tasks.

At a high level, both algorithms operate under the assumption that we have access to an expert policy profile $\pi^* = (\pi_1^*, \ldots, \pi_N^*)$ with associated Q-functions $Q^{\pi^*} = (Q_1^{\pi^*}, \ldots, Q_N^{\pi^*})$, where each $\pi_i^*$ represents a high-performing opaque policy for agent $i$, trained using any existing multi-agent RL algorithm. The goal of our approach is to distill each expert policy $\pi_i^*$ into a compact, interpretable decision tree policy $\hat{\pi}i$ that roughly approximates the expert's behavior while offering greater transparency. Importantly, both IVIPER and MAVIPER are model-agnostic in the sense that they do not depend on the specific type of expert policy architecture. This generality allows our algorithms to serve as drop-in interpretable replacements across a range of multi-agent learning settings.

In the following sections, we first describe IVIPER, which serves as a relatively straightforward extension of the original VIPER algorithm to the multi-agent case. We then introduce MAVIPER, which incorporates additional coordination mechanisms during training to improve joint performance across agents.

---

**Algorithm 3** IVIPER in Multi-Agent Setting

---

**Input**: $(X, A, P, R), \pi^*, Q^{\pi^*} = (Q_1^{\pi^*}, ..., Q_N^{\pi^*}), K, M$
**Output**: $\hat{\pi}_1, ..., \hat{\pi}_N$

1: **for** i=1 to N **do**
2:      Initialize dataset $\mathcal{D}_i \leftarrow \emptyset$ and policy $\hat{\pi}_i^0 \leftarrow \pi_i^*$
3:      **for** $m = 1$ to $M$ **do**
4:          Sample $K$ trajectories: $\mathcal{D}_i^m \leftarrow \{(x, \pi_1^*(o_1), ..., \pi_N^*(o_N)) \sim d^{\hat{\pi}_i^{m-1}, \pi_{-i}^*}\}$
5:          Aggregate dataset $\mathcal{D}_i \leftarrow \mathcal{D}_i \cup \mathcal{D}_i^m$
6:          Resample dataset according to loss:
             $\mathcal{D}_i' \leftarrow \{(x, \overrightarrow{a}) \sim p((x, \mathbf{a})) \propto \tilde{l}_i(x) \mathbb{I}[(x, \mathbf{a}) \in \mathcal{D}_i]\}$
7:          Train decision tree $\hat{\pi}_i^m \leftarrow \text{TrainDecisionTree}(\mathcal{D}_i')$
8:      **end for**
9:      Get best policy $\hat{\pi}_i \leftarrow \text{BestPolicy}(\hat{\pi}_i^1, ..., \hat{\pi}_i^M, \pi_{-i}^*)$
10: **end for**
11: **return** Best policies for each agent $\hat{\pi} = (\hat{\pi}_1, ..., \hat{\pi}_N)$

---

## 5.4 IVIPER

Motivated by the practical success of single-agent RL algorithms in the multi-agent setting [22, 205], we extend single-agent VIPER to the multi-agent setting by independently applying the single-agent algorithm to each agent, with a few critical changes described below. This algorithm can be viewed as transforming the multi-agent learning problem to a single-agent one, in which other agents are folded into the environment. Algorithm 3 shows the full IVIPER pseudocode.

### 5.4.1 Handling Non-Stationarity

First, we ensure that each agent has sufficient information for training its decision-tree policy. Observe that, if we apply VIPER independently, each agent has its own dataset $\mathcal{D}_i$ of training tuples (initialized in line 2, Algorithm 3). When populating the dataset in line 4, we have a choice: each tuple can consist only of the information relevant to that individual agent (i.e., for each item $j$, $\mathcal{D}_i[j] = (o_i, a_i)$), or of the global information available for all agents (i.e., for each item $j$, $\mathcal{D}_i[j] = (x, \pi_1^*(o_1), \ldots, \pi_N^*(o_N))$).

Recall that VIPER works with any underlying multi-agent RL algorithm. To align with the CTDE setting, many multi-agent actor-critic methods leverage a per-agent centralized critic network $Q_i^\pi$, as defined in Equation (5.1). This function is used to resample the dataset in the VIPER algorithm, meaning we need to ensure that each agent's dataset $\mathcal{D}_i$ has not only its observation and actions, but also the complete state information $x$—which consists of the observations of all of the agents—and the expert-labeled actions of all of the other agents $\pi_j^*(o_j) \forall j \neq i$. By providing each agent with the information about all other agents, we avoid the stationarity issue that arises when the policies of all agents are changing throughout the training process.

## 5.4.2 Data Collection and Relabeling

Second, we account for important changes that emerge from moving to a multi-agent formalism. When we collect trajectories for training each agent's decision-tree policy, observe that the data is drawn from the distribution $d^{\hat{\pi}_i^{m-1}, \pi_{-i}^*}$ induced by agent $i$'s policy at the previous iteration $\hat{\pi}_i^{m-1}$ and the expert policies of all other agents $\pi_{-i}^*$ (line 4). This distribution is induced because of the *independent* nature of this algorithm. In other words, we cannot perform rollouts with the decision-tree policies of the other agents because there is no communication between the decision-tree agents during this training process.

As a result, we only relabel the action for agent $i$ because the other agents choose their actions according to $\pi^*$. It is equivalent to treating all other expert agents as part of the environment and only using decision-tree policy for agent $i$.

## 5.4.3 Dataset Resampling

Third, we incorporate the actions of all agents when resampling the dataset to construct a new, weighted dataset (line 6, Algorithm 3). If the multi-agent RL algorithm uses a centralized critic (Equation (5.1)), then we resample points according to:

$$p((x, a_1, ..., a_N)) \propto \tilde{l}_i(x)\mathbb{I}[(x, a_1, ..., a_N) \in \mathcal{D}_i], \tag{5.2}$$

where,

$$\tilde{l}_i(x) = V_i^{\pi^*}(x) - \min_{a_i \in \mathcal{A}_i} Q_i^{\pi^*}(x, a_i, \overrightarrow{a}_{-i})|_{\mathbf{a}_{-i}=\pi_j^*(o_j)\forall j\neq i}. \tag{5.3}$$

Crucially, we include the actions of all other agents in Equation (5.3) to select agent $i$'s minimum Q-value from its centralized state-action value function.

When applied to value-based methods, IVIPER is more similar to single-agent VIPER. In particular, in line 4, Algorithm 3, it is sufficient to only store $o_i$ and $\pi_i^*(o_i)$ in the dataset $\mathcal{D}_i^m$, although we still must sample trajectories according to $\hat{\pi}_i^{m-1}$ and $\pi_{-i}^*$. In line 6, we use $\tilde{l}(x) = V_i^{\pi^*}(s) - \min_{a_i \in \mathcal{A}_i} Q_i^{\pi^*}(o_i, a_i)$ from single-agent VIPER, removing the reliance of the loss on a centralized critic.

### Limitations

This approach works well if i) we only want an interpretable policy for a single agent in a multi-agent setting or ii) agents do not need to *coordinate* with each other. When coordination is needed, this algorithm does not reliably capture coordinated behaviors, as each decision-tree is trained independently without consideration for what the other agent's resulting decision-tree policy will learn. This issue is particularly apparent when trees are constrained to have a small maximum depth, as is desired for interpretability.

---

**Algorithm 4** MAVIPER (Joint Training)

---

**Input:** $(\mathcal{X}, A, P, R), \pi^*, Q^{\pi^*} = (Q_1^{\pi^*}, \ldots, Q_N^{\pi^*}), K, M$

**Output:** $(\hat{\pi}_1, \ldots, \hat{\pi}_N)$

1: Initialize dataset $\mathcal{D} \leftarrow \emptyset$ and policies $\hat{\pi}_i^0 \leftarrow \pi_i^*$ for all $i \in N$
2: **for** $m = 1$ to $M$ **do**
3:     Sample $K$ trajectories: $\mathcal{D}^m \leftarrow \{(x, \pi_1^*(o_1), \ldots, \pi_N^*(o_N)) \sim d^{(\hat{\pi}_1^{m-1}, \ldots, \hat{\pi}_N^{m-1})}\}$
4:     Aggregate dataset: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}^m$
5:     For each agent $i$, resample: $\mathcal{D}_i \leftarrow \{(x, \vec{a}) \sim p((x, \vec{a})) \propto \tilde{l}_i(x)\mathbb{I}[(x, \vec{a}) \in \mathcal{D}]\}$ for all $i \in N$
6:     Jointly train decision trees: $(\hat{\pi}_1^m, \ldots, \hat{\pi}_N^m) \leftarrow \textsc{TrainJointTrees}(\mathcal{D}_1, \ldots, \mathcal{D}_N)$
7: **end for**
8: **return** best set of agents
9: $\hat{\pi} = (\hat{\pi}_1, \ldots, \hat{\pi}_N) \in \{(\hat{\pi}_1^1, \ldots, \hat{\pi}_N^1), \ldots, (\hat{\pi}_1^M, \ldots, \hat{\pi}_N^M)\}$

---

## 5.5 MAVIPER

To address the issue of coordination, we propose MAVIPER, our novel algorithm for centralized training of coordinated multi-agent decision-tree policies. For expository purpose, we describe MAVIPER in a fully cooperative setting, then explain how to use MAVIPER for mixed cooperative-competitive settings.

At a high-level, MAVIPER trains all of the decision-tree policies, one for each agent, in a centralized manner. It jointly grows the trees of each agent by predicting the behavior of the other agents in the environment using their anticipated trees. To train each decision-tree policy, MAVIPER employs a new resampling technique to find states that are critical for its interactions with other agents. Algorithm 4 shows the full MAVIPER algorithm. Because we explicitly account for the anticipated behavior of other agents in both the predictions and the resampling probability, we hypothesize that MAVIPER will better capture coordinated behavior. Specifically, MAVIPER is built upon the following extensions to IVIPER that aim at addressing the issue of coordination.

### 5.5.1 Collecting Data With Updated Decision-Tree Policies

First, we observe that only using the decision-tree policy corresponding to a single agent in each iteration $m$, as in Section 5.4.2, will yield miscoordinated policies. This approach assumes that the rollout distribution used to train each agent can be generated by mixing one learned policy with fixed expert policies for the others. However, unless each student policy exactly replicates its corresponding expert, the resulting rollout distribution will diverge from the one induced when all agents follow their own learned policies. In other words, in general, $d^{(\hat{\pi}_i^{m-1}, \pi_{-i}^*)} \neq d^{(\hat{\pi}_1^{m-1}, \ldots, \hat{\pi}_N^{m-1})}$. This mismatch is particularly problematic in settings that require tight coordination, since agent behaviors trained under inconsistent distributions may fail to compose effectively. This distribution shift introduces cascading errors that degrade overall coordination.

To address the distribution mismatch, MAVIPER uses the decision-tree policies $(\hat{\pi}_1^{m-1}, \ldots, \hat{\pi}_N^{m-1})$ from the last iteration to perform rollouts and collect new data (line 3, Algorithm 4). This approach more closely follows the original intent of the DAGGER algorithm: to gather data from the state distribution that the learned policies are likely to visit at test time. Using the learned policies during rollouts helps ensure that the dataset reflects the actual trajectories agents will encounter. This choice also has the benefit of producing a shared dataset $\mathcal{D}$ among all agents to aid in promoting coordination.

## 5.5.2 Resampling Based on Average-Case

Second, notice that the IVIPER loss (Equation (5.3)) treats the other agents as stationary experts. This assumption is problematic, as there is no guarantee that the learned decision-tree policies will be optimal. As a result, the assumption that the decision-tree policies of the other agents will align with the corresponding experts is unlikely to be true in practice.

To mitigate this issue, MAVIPER modifies this resampling probability to consider the expected impact of an agent's decision when the other agents may not be fixed to an expert. This choice enables us to explicitly factor in uncertainty over the possible actions of other agents. Specifically, for each agent $i$, MAVIPER defines the corresponding loss $l_i(x)$ as the expected performance gap between the expert action and the worst-possible action that agent $i$ could take, averaged over the possible actions of the other agents. Then we establish the new loss function for MAVIPER (which serves as the basis for the resampling probability) as:

$$\tilde{l}_i(x) = \mathbb{E}_{\mathbf{a}_{-i}} \left[ Q_i^{\pi^*}(x, \pi_i^*(x), \mathbf{a}_{-i}) - \min_{a_i \in \mathcal{A}_i} Q_i^{\pi^*}(x, a_i, \overrightarrow{a}_{-i}) \right]. \tag{5.4}$$

Intuitively, this quantity serves as a proxy for how critical it is for agent $i$ to match the expert at that state, on average. High values of $\tilde{l}_i$ indicate states where choosing a poor action can result in significant value loss, regardless of how others act. Averaging over the actions of other agents accounts for uncertainty in their behavior, as it is unlikely that the other agents will act exactly like the corresponding expert (as is assumed by IVIPER). This choice should make the resulting policy more robust to imperfect coordination.

## 5.5.3 Predicting the Actions of Other Agents

Third, we make the following observation with respect to the training procedure of IVIPER. IVIPER essentially grows the decision-trees depth-first, meaning that it implicitly trains each student under the assumption that the others will exactly mimic the corresponding expert policies. We know in practice this outcome is unlikely to occur, given the differences in representational capacities of the decision trees and experts. MAVIPER instead introduces a joint training procedure, detailed in Algorithm 5.

---

**Algorithm 5** TRAINJOINTTREES($\mathcal{D}_1, \ldots, \mathcal{D}_N$)

---

**Output:** Joint policies $(\hat{\pi}_1, \ldots, \hat{\pi}_N)$

1: **Function** TRAINJOINTTREES($\mathcal{D}_1, \ldots, \mathcal{D}_N$)
2: Initialize decision trees $\hat{\pi}_1^m, \ldots, \hat{\pi}_N^m$
3: **repeat**
4:     Grow one more level for agent $i$: $\hat{\pi}_i^m \leftarrow \text{BUILD}(\hat{\pi}_1^m, \ldots, \hat{\pi}_N^m, \mathcal{D}_i)$
5:     $i \leftarrow (i + 1) \mod N$
6: **until** all trees have reached the maximum allowed depth
7: **return** $\hat{\pi}_1^m, \ldots, \hat{\pi}_N^m$

8: **Function** BUILD($\hat{\pi}_1^m, \ldots, \hat{\pi}_N^m, \mathcal{D}_i$)
9: **for** each data point $(x, \vec{a}) \in \mathcal{D}_i$ **do**
10:     // Will each agent's projected final decision-tree predict correctly?
11:     $v_j \leftarrow \mathbb{I}[\text{PREDICT}(\hat{\pi}_j^m, x) = a_j]$ for all $j \in [1, N]$
12:     // Keep point only if many agents predict correctly
13:     **if** $\sum_{j=1}^{N} v_j <$ threshold **then**
14:         $\mathcal{D}_i \leftarrow \mathcal{D}_i \setminus \{(x, \vec{a})\}$
15:     **end if**
16: **end for**
17: Compute best next feature split for $\hat{\pi}_i^m$ using $\mathcal{D}_i$
18: **return** $\hat{\pi}_i^m$

19: **Function** PREDICT($\hat{\pi}_j^m, x$)
20: Use $x$ to traverse $\hat{\pi}_j^m$ to leaf node $l(x)$
21: Train projected final tree $\hat{\pi}_j' \leftarrow \text{TrainDecisionTree}(\mathcal{D}_j)$
22: **return** $\hat{\pi}_j'.\text{predict}(x)$

---

Specifically, MAVIPER removes this assumption by jointly training the decision trees for all agents. Instead of fixing other agents' policies during training, MAVIPER grows all trees in parallel. At each step, a single agent expands its tree by one level, and control rotates across agents in a round-robin fashion. This alternating expansion procedure ensures that the policies evolve together, so each agent's tree is shaped in light of its peers' current behavior.

This joint training setup introduces a new challenge: during training, agents must make decisions about how to grow their trees without knowing what actions the other agents will ultimately take. Because the other trees are still under construction, their final behavior is uncertain. This uncertainty makes it difficult to determine which states are worth training on. To tackle this issue, we introduce a prediction module to the decision-tree training process, as shown in the Predict function. The goal of the prediction module is to predict the actions that the other decision-trees $\{\hat{\pi}_j\}_{j \neq i}$ might make, given their partial observations. Since the trees are not complete at the time of prediction, we use the output of another decision-tree trained with the full dataset associated with that node for the prediction. This idea is illustrated in Figure 5.2.

Figure 5.2: **MAVIPER includes a prediction module for promoting joint accuracy**. When jointly growing the decision-tree policies, each agent does not know which action the other agents will likely take. As a result, each agent essentially hypothesizes what the resulting trees of the other agents will look like, then uses the anticipated trees to gather the likely actions the other agents will take. Data points with high joint errors are then removed from the dataset.

Following the intuition that the correct prediction of one agent alone may not yield much benefit if the other agents are wrong, we use this prediction module to remove all data points whose proportion of correct predictions is lower than a predefined threshold. We then calculate the splitting criteria based on this modified dataset and continue iteratively growing the tree. The goal of this procedure is to increase the *joint* accuracy of all agents; in other words, we aim to prioritize states where coordination is likely to succeed under the current partial policies.

## Extending to Mixed Cooperative-Competitive Settings

In some mixed cooperative-competitive settings, agents in a team share goals and need to coordinate with each other, but they face other agents or other teams whose goals are not fully aligned with theirs. In these settings, MAVIPER follows a similar procedure but applied on a per-team basis. More specifically, for a team $Z$, the `Build` and `Predict` function is constrained to only make predictions for the agents in the same team.

Equation (5.4) now takes the expectation over the joint actions for agents outside the team and becomes:

$$\tilde{l}_i(x) = \mathbb{E}_{\mathbf{a}_{-Z}}\left[ Q_i^{\pi^*}(x, \pi_Z^*(x), \mathbf{a}_{-Z}) - \min_{a_i \in \mathcal{A}_i, i \in Z} Q_i^{\pi^*}(x, \overrightarrow{a}_Z, \overrightarrow{a}_{-Z}) \right]. \tag{5.5}$$

## 5.6 Experimental Setup

This section establishes the experimental setup we use to evaluate how well MAVIPER and IVIPER agents perform in a variety of environments. Because the goal is to learn high-performing yet interpretable policies, the quality of the trained policies is empirically evaluated in three multi-agent environments: two mixed competitive-cooperative environments and one fully cooperative environment. We evaluate the agents based on the performance of the decision-tree policies because the goal is to deploy these policies in place of the expert ones. We organize our experiments around the following research questions:

**RQ 1** Individual performance: How well do MAVIPER and IVIPER perform when we replace a single neural-network policy with a corresponding decision-tree policy?

**RQ 2** Joint performance: How well do MAVIPER and IVIPER perform when we replace a team of neural-network policies with the corresponding decision-tree policies?

**RQ 3** Robustness: How well do MAVIPER and IVIPER perform when we consider a variety of opponent policies?

**RQ 4** Ablation Study: Which components contribute most to the performance of MAVIPER?

**Baselines.**    We compare to two baselines that train decision-tree policies, Fitted Q-Iteration and Imitation DT. In Fitted Q-Iteration, we iteratively approximate the Q-function with a regression decision tree [91]. We discretize states to account for continuous state values. More details in Appendix 1.3. We derive the policy by taking the action associated with the highest estimated Q-value for that input state. In Imitation DT, each decision-tree policy is directly trained using a dataset collected by running the expert policies for multiple episodes. No resampling is performed. The observations for an agent are the features, and the actions for that agent are the labels.

**Additional details.**    We test the aforementioned algorithms under the following experimental conditions. Since small decision trees are considered interpretable, we constrain the maximum tree depth to be at most 6. The expert policies used to guide the decision-tree training are generated by MAD-DPG [193][1]. We detail the hyperparameters and the hyperparameter-selection process in Appendix 1.2. We train a high-performing MADDPG expert, then run each decision-tree-learning algorithm 10 times with different random seeds. We evaluate all policies by running 100 episodes. Error bars correspond to the $95\%$ confidence interval. Our code is available through our project website: `https://stephmilani.github.io/maviper/`.

---

[1]We use the Pytorch [250] implementation `https://github.com/shariqiqbal2810/maddpg-pytorch`.

## Environments and Evaluation Metrics

We evaluate the quality of the trained policies in three multi-agent particle world environments [193], described below. Episodes terminate when the maximum number of timesteps $T = 25$ is reached. We measure how well the decision-tree policies perform in the environment because our goal is to *deploy* these policies, not the expert ones. We choose the primary performance metric based on the environment (detailed below), and we also provide results using expected return as the performance metric in Appendix 2. Since small decision trees are considered interpretable, the maximum depth is constrained to be at most 6.

**Physical deception.**    In this environment, a team of $N$ defenders must protect $N$ targets from one adversary. One of the targets is the true target, which is known to the defenders but not to the adversary. For our experiments, $N = 2$. Defenders succeed during an episode if they split up to cover all of the targets simultaneously; the adversary succeeds if it reaches the true target during the episode. Covering and reaching targets is defined as being $\epsilon$-close to a target for at least one timestep during the episode. We use the defenders' and the adversary's success rate as the primary performance metric in this environment.

**Cooperative navigation.**    This environment consists of a team of $N$ agents, who must learn to cover all $N$ targets while avoiding collisions with each other. For our experiments, $N = 3$. Agents succeed during an episode if they split up to cover all of the targets without colliding. Our primary performance metric is the summation of the distance of the closest agent to each target, for all targets. Low values of the metric indicate that the agents correctly learn to split up.

**Predator-prey.**    This variant involves a team of $K$ slower, cooperating predators that chase $M$ faster prey. There are $L = 2$ landmarks impeding the way. We choose $K = M = 2$. We assume that each agent has a restricted observation space mostly consisting of binarized relative positions and velocity (if applicable) of the landmarks and other agents in the environment. See Appendix 1.1 for full details. Our primary performance metric is the number of collisions between predators and prey. For prey, lower is better; for predators, higher is better.

## 5.7  Results

In this section, we present the experimental results organized by research question (RQ1–RQ4). Across all evaluations, we compare our methods (MAVIPER and IVIPER) against established baselines, using the performance metrics detailed in the experimental setup.

(a) Physical Deception

(b) Cooperative Navigation

(c) Predator-prey

Figure 5.3: **Individual MAVIPER-trained agents perform at least as well as IVIPER and the baselines across all maximum depths for all environments.** Here, we analyze the individual performance ratio: the relative performance when only one agent adopts decision-tree policy and all other agents use expert policy. Higher is better. Error bars represent the 95% confidence interval. In some instances, with sufficient representational capacity, MAVIPER-trained agents can recover 85% or higher of the expert performance.

## RQ 1: Individual Performance Compared to Experts

We evaluate whether the performance of the individual decision-tree policies can match that of the expert neural network policies. As a result, we analyze the performance of the decision-tree policies when only one agent adopts the decision-tree policy while all other agents use the expert policies. Given a decision-tree policy profile $\hat{\pi}$ and the expert policy profile $\pi^*$, if agent $i$ who belongs to team $Z$ uses its decision-tree policy, then the individual performance ratio is defined as:

$$\frac{U_Z(\hat{\pi}_i, \pi^*_{-i})}{U_Z(\pi^*)},$$

where $U_Z(\cdot)$ is team $Z$'s performance given the agents' policy profile (reflecting that our primary performance metrics are defined at the team level). A ratio of 1 indicates parity with expert performance. Ratios above 1 can arise because we compare the performance of the decision-tree and expert policies in the environment, not the *similarity* of the policies. We now present the results using this metric.

**Individual MAVIPER-trained agents perform at least as well as other methods.**    We first establish that individual MAVIPER agents perform at least as well as other methods (see Figure 5.3). In the physical deception environment, individual IVIPER defenders also outperform the baselines for all maximum depths. For the adversary role in the physical deception environment and in the cooperative navigation environment, MAVIPER and IVIPER agents perform similarly to at least one baseline when replacing the adversary and for the cooperative navigation environment. This result indicates that the correct strategy may be simple enough to capture with a less-sophisticated algorithm. In predator-prey, the most notable performance difference occurs for the predator agents. When the maximum depth is 2, MAVIPER achieves near-expert performance. When the maximum depths are 4 and 6, MAVIPER and IVIPER agents achieve similar performance and significantly outperform the baselines. The preys achieve similar performance across all algorithms. We suspect that the complexity of this environment makes it challenging to replace even a single prey's policy with a decision tree.

**MAVIPER achieves a reasonable performance ratio at the maximum depth of** 6 **across all environments.**    Specifically, we want to investigate whether, given sufficient capacity, MAVIPER achieves good performance. MAVIPER agents achieve a performance ratio above $0.75$ in all environments with a maximum depth of 6. The same is true for IVIPER, except for the adversaries in physical deception. That means decision-tree policies generated by IVIPER and MAVIPER lead to a performance degradation of less than or around 20% compared to the less interpretable neural-network-based expert policies. In contrast, even with the maximum depth, the baselines struggle on some of the environments (e.g., Fitted Q-Iteration struggles to achieve a performance ratio of even $0.25$ for the adversary on Physical Deception. These results show that IVIPER and MAVIPER generate reasonable decision-tree policies and outperform the baselines overall when adopted by a single agent.

## RQ 2: Joint Performance Compared to Experts

As we establish, a crucial aspect in multi-agent environments is agent coordination, especially when agents are on the same team with shared goals. To ensure that the decision-tree policies capture this coordination, we analyze the performance of the decision-tree policies when all agents in a team adopt decision-tree policies, while other agents use expert policies. We define the joint performance ratio as:

$$\frac{U_Z(\hat{\pi}_Z, \pi^*_{-Z})}{U_Z(\pi^*)},$$

where $U_Z(\hat{\pi}_Z, \pi^*_{-Z})$ is the utility of team $Z$ when using their decision-tree policies against the expert policies of the other agents $-Z$. This metric captures any coordinated performance degradation compared to the experts. We now present the results using this metric.

(a) Physical Deception

(b) Cooperative Navigation

(c) Predator-prey

Figure 5.4: **MAVIPER teams yield better-coordinated policies than the two baselines across all maximum depths and all environments.** This plot shows the joint performance ratio; that is, the relative performance when all agents in a team adopt decision-tree policy and other agents use expert policy. Higher ratios are better. Error bars represent the 95% confidence interval.

**MAVIPER teams significantly outperform teams trained by the two baselines for all maximum depths for all environments.** Figure 5.4 shows the mean joint performance ratio for each team, averaged over all trials. MAVIPER agents also significantly outperform IVIPER agents for 7 out of 11 environment-depth combinations. This result shows that MAVIPER agents better capture coordinated behavior, even when we introduce another cooperating agent (Figures 5.4a and 5.4b). Figure 5.4b shows that MAVIPER agents significantly outperform all other algorithms in the cooperative navigation environment for all maximum depths. These results indicate that MAVIPER better captures the coordinated behavior necessary for a team to succeed in different environments.

Figure 5.5: **MAVIPER defenders most commonly split importance across the two targets, leading to the correct behavior.** We show this by visualizing the feature importances of the features used by the decision-tree policies for the two defenders in the physical deception environment. We note that the actual features used in the environment are the relative positions of that agent and the labeled feature. We shorten the labeled feature for ease of reading. Darker squares (higher values) correspond to higher feature importance.

**MAVIPER defenders achieve better coordination by dividing feature attention.** Figure 5.4a shows that MAVIPER defenders outperform IVIPER and the baselines, indicating that it better captures the coordinated behavior necessary to succeed in this environment. In this environment, it is critical for the two defenders to split up to cover the two targets: if both agents cover the same target, the adversary can immediately infer which is the correct target to attack. We hypothesize that this superior performance arises because MAVIPER defender agents effectively split their "attention" between the two targets, enabling them to cover both simultaneously. To investigate this hypothesis, we examine the normalized average feature importances of the decision-tree policies (depth 4) for both IVIPER and MAVIPER over 5 of trials. Plotted in Figure 5.5, the analysis reveals that each MAVIPER defender (top) predominantly focuses on attributes associated with one of the two targets. More specifically, defender 1 focuses on target 2 and defender 2 focuses on target 1. In contrast, both IVIPER defenders (bottom) focus mainly on the attributes associated with the *goal* target. This overlapping focus in feature space not only eliminates the chance that IVIPER agents learn the correct coordinated covering behavior but also renders them more vulnerable to the adversary, as it becomes easier to infer the correct target.

## RQ 3: Robustness to Different Opponents

We investigate the robustness of the decision-tree policies when a team using decision-tree policies plays against a variety of opponents in the mixed competitive-cooperative environments. For this set of experiments, we choose a maximum depth of 4. Given a decision-tree policy profile $\hat{\pi}$, a team Z's performance against an alternative policy profile $\pi'$ used by the opponents is: $U_Z(\hat{\pi}_Z, \pi'_{-Z})$. We consider a broad set of opponent policies $\pi'$, including the policies generated by MAVIPER, IVIPER, Imitation DT, Fitted Q-Iteration, and MADDPG.

| Environment | Team | MAVIPER | IVIPER | Imitation DT | Fitted Q-Iteration |
|---|---|---|---|---|---|
| Physical Deception | Defender | **0.77 ± 0.01** | 0.33 ± 0.01 | 0.24 ± 0.03 | 0.00 ± 0.00 |
| | Adversary | **0.42 ± 0.03** | 0.41 ± 0.03 | **0.42 ± 0.03** | 0.07 ± 0.01 |
| Predator-prey | Predator | **2.51 ± 0.72** | 1.98 ± 0.58 | 1.14 ± 0.28 | 0.26 ± 0.11 |
| | Prey | 1.76 ± 0.80 | **2.16 ± 1.24** | **2.36 ± 1.90** | 1.11 ± 0.82 |

Table 5.1: **MAVIPER agents are more robust to possible attackers compared with all other algorithms.** We report mean team performance and standard deviation of decision-tree policies for each team, averaged across a variety of opponent policies. The best-performing algorithm for each agent type is shown in **bold**. Higher value indicates better performance against a variety of opponents. Lower standard deviation indicates less variation in performance across different opponent policies.

**MAVIPER defenders are generally more robust to possible attackers compared with all other algorithms.** We first investigate the mean team performance averaged over all opponent policies (Table 5.1). For physical deception, MAVIPER defenders outperform all other algorithms, with a gap of 0.44 between its performance and the next-best algorithm, IVIPER. This result indicates that MAVIPER learns coordinated defender policies that perform well against various adversaries. MAVIPER, IVIPER, and Imitation DT adversaries perform similarly on average, with a similar standard deviation, which supports the idea that the adversary's desired behavior is simple enough to capture with a less-sophisticated algorithm. For predator-prey, MAVIPER predators and prey outperform all other algorithms. The standard deviation of the performance of all algorithms is high due to this environment's complexity.

### Decomposed Results Based on Environment

We present the full robustness results for the predator-prey and physical deception environments. For space reasons, we only report the average over the 100 trials; however, we only label the best-performing agent of each type in either red or blue if the 95% confidence intervals do not overlap, unless otherwise mentioned. We exclude MADDPG from this calculation, since we know that MADDPG agents will outperform all other agent types, and we are mostly interested in how well the *decision tree* policies perform.

**For predator-prey, MAVIPER predators are strictly more robust than all other agents, and prey are generally more robust.** Table 5.2 shows the full decomposed results for both predator and prey across all algorithms. MAVIPER predators are strictly more robust than all other agents (except MADDPG) to different types of prey. MAVIPER prey are the most or second most robust to different types of predators. In this environment, predator coordination is more critical, as predators must strategically catch the prey. The prey, on the other hand, does not require much coordination, which explains the Imitation DT prey's robustness by imitating the action of the single-agent expert. However, MAVIPER prey are still generally more robust than all other agents to different types of predators.

|  | Prey | | | | |
|---|---|---|---|---|---|
| **Predator** | MAVIPER | IVIPER | Imitation DT | Fitted Q-Iteration | MADDPG |
| MAVIPER | (2.28, 2.28) | (3.49, 3.49) | (2.41, 2.41) | (3.01, 3.01) | (1.37, 1.37) |
| IVIPER | (1.95, 1.95) | (2.46, 2.46) | (2.17, 2.17) | (2.44, 2.44) | (0.88, 0.88) |
| Imitation DT | (1.32, 1.32) | (1.17, 1.17) | (1.18, 1.18) | (1.40, 1.40) | (0.61, 0.61) |
| Fitted Q-Iteration | (0.46, 0.46) | (0.30, 0.30) | (0.24, 0.24) | (0.18, 0.18) | (0.14, 0.14) |
| MADDPG | (2.78, 2.78) | (3.36, 3.36) | (5.82, 5.82) | (4.98, 4.98) | (2.54, 2.54) |

Table 5.2: **Robustness results of decision-tree agents on predator-prey.** Results are presented as: average number of touches in an episode. Higher is better for predator, and lower is better for prey. Excluding MADDPG, the best-performing prey (lowest in value) for each predator type is in blue and the best-performing predator (highest in value) for each prey type is in red.

|  | Defender | | | | |
|---|---|---|---|---|---|
| **Adversary** | MAVIPER | IVIPER | Imitation DT | Fitted Q-Iteration | MADDPG |
| MAVIPER | (0.42, 0.76) | (0.45, 0.33) | (0.45, 0.23) | (0.37, 0.01) | (0.40, 0.93) |
| IVIPER | (0.39, 0.78) | (0.45, 0.32) | (0.40, 0.23) | (0.38, 0.00) | (0.43, 0.92) |
| Imitation DT | (0.40, 0.79) | (0.42, 0.34) | (0.46, 0.26) | (0.38, 0.01) | (0.46, 0.92) |
| Fitted Q-Iteration | (0.07, 0.77) | (0.06, 0.33) | (0.07, 0.19) | (0.08, 0.00) | (0.08, 0.79) |
| MADDPG | (0.71, 0.76) | (0.77, 0.32) | (0.77, 0.26) | (0.58, 0.00) | (0.62, 0.90) |

Table 5.3: **Robustness results of decision-tree agents on physical deception.** Results are presented as: (adversary success ratio, defender success ratio). Higher is better. Excluding MADDPG, the best-performing defender for each adversary type is in blue and the best-performing adversary for each defender type is in red.

**For physical deception, MAVIPER defenders are the most robust than all agents to different types of adversaries.** Table 5.3 shows the full results. Interestingly, MAVIPER, IVIPER, and Imitation DT adversaries all perform similarly. Indeed, they often do not achieve performance that is statistically significant from one another, as measured by the 95% confidence interval. However, we still highlight the best-performing adversary in red to more easily show the attained performance. Note that MADDPG adversaries can occasionally achieve success greater than around 0.50, which means that these adversaries can take advantage of some information about the defenders to correctly choose the target to visit. In contrast, adversaries trained with any of the decision-tree-learning algorithms never achieve greater than 0.50, which indicates that they may need a more complex representation to capture important details about the defenders.

Figure 5.6: **Ablation study for MAVIPER for a maximum depth of** 4. MAVIPER (No Prediction) does not utilize the predicted behavior of the anticipated decision-trees of the other agents to grow each agent's tree. MAVIPER (IVIPER Resampling) uses the same resampling method as IVIPER. Both ablated changes contribute to the improvement of MAVIPER over IVIPER when considering the joint performance of the defenders.

## RQ 4: Ablation Study

As we establish in Section 5.5, MAVIPER improves upon IVIPER with a few critical changes. First, we utilize the predicted behavior of the anticipated decision-trees of the other agents to grow each agent's tree. Second, we alter the resampling probability to incorporate the average Q-values over all actions for the other agents. To investigate the contribution of these changes to the performance, we run an ablation study with a maximum depth of 4 on the physical deception environment. We report both the mean independent and joint performance ratios for the defender team in Figure 5.6, comparing MAVIPER and IVIPER to two variants of MAVIPER without one of the two critical changes. Specifically, MAVIPER (No Prediction) does not utilize the predicted behavior of the anticipated decision-trees of the other agents to grow each agent's tree. MAVIPER (IVIPER Resampling) uses the same resampling method as IVIPER.

**The effect of the MAVIPER changes emerges primarily for coordination.** Interestingly, due to the high variance in the performance, all algorithms are effectively equivalent for the adversaries. Recall that there is only one adversary, so coordination is not needed. This result means that the benefits of the MAVIPER changes emerges primarily when coordination is needed. When investigating the individual performance for the defenders, the most significant algorithmic change for this metric seems to be the prediction module. When investigating the joint performance for the defenders, both ablated changes contribute to the improvement over IVIPER.

## 5.8  Related Work

**Interpretable single-agent RL.**   Most work on interpretable RL is in the single-agent setting [219]. We first discuss techniques that directly learn decision-tree policies. CUSTARD [321] extends the action space of the Markov decision process to contain actions for constructing a decision tree, i.e., choosing a feature to branch a node. Training an agent in the augmented Markov decision process yields a decision-tree policy for the original Markov decision process while still enabling training using any function approximator, like neural networks, during training. By redefining the Markov decision process, the learning problem becomes more complex, which is problematic in multi-agent settings where the presence of other agents already complicates the learning problem. A few other works directly learn decision-tree policies [91, 207, 328] for single-agent RL but not for the purpose of interpretability. Further, these works have custom learning algorithms and cannot utilize a high-performing neural-network policy to guide training.

**Post-hoc interpretable RL.**   VIPER is considered a post-hoc decision-tree-learning method [17]; however, we use it to produce *intrinsically* interpretable policies for deployment. MOET [331] extends VIPER by learning a mixture of decision-tree policies trained on different regions of the state space. The resulting policy is a linear combination of multiple trees with non-axis-parallel partitions of the state. We find that the performance difference between VIPER and MOET is not significant enough to increase the complexity of the policy structure, which would sacrifice interpretability.

**Interpretable multi-agent RL.**   Despite increased interest in interpretable single-agent RL, work on interpretable multi-agent RL is limited. One line of work generates explanations from non-interpretable policies. Some work uses attention [143, 185, 230] to select and focus on critical factors that impact agents in the training process. Other work generates explanations as verbal explanations with predefined rules [343] or Shapley values [131]. The most similar line of work to ours [160] approximates non-interpretable multi-agent RL policies to interpretable ones using the framework of abstract argumentation. This work constructs argument preference graphs given manually-provided arguments. In contrast, our work does not need these manually-provided arguments for interpretability. Instead, we generate decision-tree policies.

## 5.9  Discussion and Future Directions

This chapter introduces novel techniques for developing interpretable policies that maintain performance. Specifically, we envision assisting stakeholders in understanding the decision-making of *teams* of AI agents, a problem previously unaddressed in the literature. To do so, we provide the first algorithms for extracting interpretable decision-tree policies in the multi-agent setting— IVIPER and MAVIPER—and show that MAVIPER-trained agents are well-suited to achieve high reward with limited policy capacity in a variety of domains.

IVIPER enables independent and parallelizable training, making it a good choice when only one agent must be interpretable. With MAVIPER, we address the challenge of achieving high reward in domains that require coordination among agents in a team. We validate that MAVIPER effectively captures coordinated behavior by showing that teams of MAVIPER-trained agents outperform the agents trained by IVIPER and several baselines. We achieve this result by introducing a novel centralized training regime, in which we explicitly account for the impact of the decisions and policies of the other agents when growing each agent's tree policy. We further address the deployment challenge of robustness to potential adversaries by showing that our algorithmic choices enable MAVIPER to generally produce more robust agents than the other decision-tree-learning algorithms.

## Future Directions

**Combining with automated feature selection algorithms.** We also note that our algorithms can work in some environments where the experts and decision-trees are trained on different sets of features. Since decision-trees can be easier to learn with a simpler set of features, future work includes augmenting our algorithm with an automatic feature selection component that constructs simplified yet still interpretable features for training the decision-tree policies. This means that we could leverage our approach for efficiently querying for concept labels alongside VIPER-style algorithms to enable policy extraction with agents that learn directly from images to produce a concept mapping to a downstream decision-tree policy.

**Studying the temporal and evolving nature of RL explanations.** Furthermore, we would like to study the comprehensibility of these policies in context, especially as decision-making unfolds over time. How do people evaluate explanations for temporally extended behavior? What explanation modalities, such as natural language rationales, visualizations, or comparisons to reference policies, could be included alongside or in lieu of the tree policies to best support users in tracking and anticipating agent behavior?

## Acknowledgments

# Part III

# ALIGNED

# 6 Aligning Reinforcement Learning Policies with Human Feedback

## 6.1 Introduction

A key dimension of human-centered RL is alignment. Alignment broadly refers to ensuring that an AI system's behavior reflects human intentions, preferences, and values. Traditionally, alignment has focused on outcomes: designing agents that *do* what humans want, even in novel situations. However, as AI agents are deployed for decision-making, some domains demand alignment not only with outcomes (i.e., *behavioral* alignment) but in thought process (i.e., *decision-making* alignment). By decision-making alignment, we mean that agents should make decisions in ways that reflect the trade-offs, principles, and reasoning that humans consider appropriate. This chapter addresses this need by focusing on the challenge of aligning the internal decision logic and structure of policies using human feedback in the form of preferences.

Specifically, this chapter considers decision-making alignment in the context of interpretable RL, where the policy structure is exposed and can be scrutinized. As we discuss in Section 2.2, there exist many types of interpretable models—such as decision trees [114, 206, 220, 260, 321] or rule lists—that can support the increasing regulatory and ethical demands for transparency and oversight of AI agents. The key advantage of these structures is that they expose the reasoning process itself, in a way that is faithful to the true decision-making process of the agent. Because these models make explicit how decisions are reached, they offer a substrate for aligning not just what the agent does but how it decides.

While interpretable policies enable this new form of alignment, achieving interpretability is itself nebulous and context-specific. In practice, many desirable properties emerge from the overall structure of a policy, not from its individual components. For example, among a set of interpretable policies with similar performance, some may better support fairness, be easier to audit by being more concise, or offer clearer hooks for human intervention. This occurrence creates a challenge: evaluating or enforcing these global properties often cannot be reduced to constraints on local parts. For instance, limiting the total number of decision nodes is a global constraint: it cannot be enforced by independently constraining each branch, since adding complexity in one part of the policy affects the capacity available elsewhere. Likewise, fairness is generally measured in expectation, across all paths. As a result, identifying and selecting interpretable policies requires reasoning about the structure as a whole.

Figure 6.1: **PASTEL overview.** PASTEL is a novel algorithm for aligning agent-decision making consisting of two main parts: i) preference elicitation and learning and ii) interpretable policy generation. In preference elicitation and learning, users provide feedback on candidate interpretable models. During interpretable policy generation, the current preference estimate guides generation of interpretable models that better align with user preferences.

Standard RL approaches excel at incrementally improving local portions of a policy, but they are not designed to capture preferences over the policy as a whole. As a result, we cannot rely on these methods to optimize for objectives that depend on the policy as a whole. To address this limitation, we propose a framework that learns user preferences over both policy attributes, and incorporates them directly into the learning objective to guide the search toward policies that better reflect human priorities. These attributes reflect key desiderata for these policies, which can be established from both regulatory guidance (e.g., the EU AI Act [5]) and human-centered design (e.g., the ease of tracing decisions for oversight and the clarity of policy outputs). This representation allows us to define user or institutional preferences as utility functions over the policy space, enabling explicit reasoning about tradeoffs and alignment. This compact representation also enables the direct usage of human feedback during *training* or test-time adaptation to develop better-aligned policies.

This chapter establishes a novel learning framework to develop interpretable agents that align with user preferences. Our approach explicitly learns underlying user preferences over policy characteristics and generates policies that are optimized for those preferences. A key insight underlying our method is representing policies as feature vectors that capture deployment-relevant attributes, such as decision-tree depth and performance. This representation enables us to compactly model preferences over the space of policies themselves, not just their outputs (e.g., trajectories), By learning in this policy feature space, we can accommodate potentially infinite candidate policies by modeling preferences over these attributes rather than explicitly counting the win rates of individual policies. This idea also produces interpretable preference models that end users or decision makers could audit, inspect, and reason about. Users can inspect not only the final policy but also the learned preference structure that guided the selection.

We make the following contributions. We propose *Preference Aligned Selection of Trees via Evolutionary Learning* (PASTEL), the first algorithm to leverage preference feedback for interpretable RL training. This algorithm consists of three main parts. First, we introduce a simple but effective evolutionary algorithm to guide training toward decision-tree policies that are better tailored toward preferences of end users. Second, we maintain a population of preference estimates to more informatively select the policies for eliciting user feedback. Third, we improve efficiency by pruning the set of candidate decision-tree policies for querying by maintaining the Pareto frontier. We demonstrate the effectiveness of our method through experiments on two RL environments using synthetic but empirically justified preference data. We show that our approach not only yields policies that are better aligned with user preferences but is also more sample efficient in the number of user queries, decreasing the burden on human users. By bridging the gap between training agents and evaluating their explanations, we believe our work opens new avenues for developing more interpretable, user-centered reinforcement learning agents.

**Contributions.** In summary, we contribute the following:

- We develop the first formal framework for preference-based interpretable RL that enables learning user preferences over policy characteristics beyond task performance.

- We develop PASTEL, the first algorithm that incorporates preference feedback directly into interpretable policy training to generate policies aligned with user preferences.

- We provide empirical validation demonstrating that PASTEL produces more preference-aligned policies compared with standard techniques.

## 6.2  Problem Statement

We first formulate the problem of aligning the decision-making process of an agent with human feedback. Specifically, we incorporate user feedback into interpretable RL to identify the most preferred policy. We propose to model the problem as an online, iterative process in which a learner interacts with a user by eliciting pairwise comparisons between candidate policies. To capture practical constraints on the amount of feedback that people can provide, the interaction proceeds over a fixed query budget, $M$.

A key challenge is that preferences over policies are often non-decomposable: users care about global properties of behavior—such as simplicity—that cannot be easily reduced to step-wise rewards. As a result, it is unclear how to integrate these properties into a standard reward function defined over local states and actions. In contrast, when preferences are only over transitions (i.e., local action choices), a model of a user's preferences can be directly used within the RL framework (using e.g., Equations (2.4) and (2.5) from Section 2.1.4 to incrementally improve a policy. In the rest of this section, we formally describe the preference elicitation framework, the preference feedback model, and the learning objective of our problem formulation.

**Preference elicitation framework.** At each iteration $m \in [M]$, the learner can query the user to obtain more information about their preferences. Standard bandit settings assume users can provide precise numerical reward specifications; however, users may struggle to articulate this feedback consistently, as it requires communicating exact values and tradeoffs. To address this issue, pairwise comparisons have emerged as a popular supervision signal [20, 40, 53, 57, 246, 252]. Rather than assigning absolute scores to items [174], users compare two alternatives $i$ and $j$, expressing a preference $i \succ j$ if item $i$ is preferred. We adopt this form of feedback in our framework. Specifically, the learner can present the user with a pair of policies $q_m = (\hat{\pi}_m^i, \hat{\pi}_m^j)$ selected from a set $\hat{\Pi}_m$. The user provides stochastic feedback $O_m \in \{0, 1\}$, realized as $o_m$, with the probability of selection determined by their underlying utility function $u$.

**Preference feedback model.** Human preferences are often context-dependent and noisy [153, 245]. Rather than assuming users always choose the utility-maximizing option, we adopt a stochastic response model. This formulation captures stochasticity in human preferences while remaining simple to estimate via maximum likelihood, as we will later see. We follow the standard Terry-Plackett-Luce model [28, 196] with a rationality coefficient [175, 288]. Specifically, we assume that a user's stochastic response $o$ to pairwise queries of policies $q = (\hat{\pi}^i, \hat{\pi}^j)$ is based on their underlying linear utility function $u$ [59, 183, 277]. Given two policies $\hat{\pi}_i, \hat{\pi}_j$ the probability that $\hat{\pi}^i$ is preferred to $\hat{\pi}^j$ is:

$$\mathbb{P}(\hat{\pi}^i \succ \hat{\pi}^j | \theta) = \frac{1}{1 + \exp\left(\beta\left(u(\hat{\pi}^j; \theta) - u(\hat{\pi}^i; \theta)\right)\right)}, \tag{6.1}$$

where $\beta$ is an inverse temperature parameter governing the stochasticity of user responses and $\boldsymbol{\theta} \in \mathbb{R}^d$ is a latent parameter vector representing user preferences. We can see that this formulation is extremely similar to the model in Equation (2.4), except we define utilities over policies instead of trajectories, and we include a $\beta$ parameter to further model the stochasticity level of the responses.

**Learning objective.** The learner's goal is to identify a high-utility policy using a bounded number of pairwise queries. Specifically, after $M$ queries, the learner outputs a recommended tree $\hat{\pi}_M$. We define the objective as selecting a policy that maximizes the user's true utility:

$$\max_{\hat{\pi} \in \hat{\Pi}} u(\hat{\pi}; \boldsymbol{\theta}), \tag{6.2}$$

where $\hat{\Pi}$ is the set of candidate policies. Since $u(\cdot; \boldsymbol{\theta})$ is not directly observable, the learner must infer which policy to recommend based only on the observed comparisons.

# 6.3  PASTEL

How can we find interpretable policies that align with what users actually want? In nearly all preference-based learning settings—such as movie or product recommendation—the learner selects from a fixed pool of items. The goal is then to efficiently identify the most preferred item using a limited number of pairwise comparisons. In these settings, it is infeasible or impossible for a learner to create new items: we (probably) cannot generate wholly new movies or build entirely new products based on users preferences. As a result, learning focuses purely on ranking existing options.

## Overview

In contrast, our insight is that our setting more closely resembles RL from human feedback for policy optimization: we can indeed control the policies that are generated. This flexibility enables us to explore a much richer space of candidate policies, including those that may not have been present at the outset. However, it introduces an additional challenge: given that the policy space is potentially infinite due to this generation ability, we need an efficient and compact way to determine the best policy. Rather than maintaining something like a per-policy win rate, we instead propose to estimate user preferences over a compact space. This choice enables us to efficiently estimate user preferences, then use these preferences to not only choose the best policy but also guide the generation of new candidate policies for querying.

Leveraging these insights, we propose PASTEL, an algorithm that interleaves preference learning and policy generation. In each iteration $m \in [M]$, the algorithm updates its estimate of the user's preferences based on pairwise comparisons and then uses this estimate to guide the generation of new interpretable policies. By maintaining a diverse population of policies over time, the algorithm both improves preference estimates and enables increasingly aligned policy generation. To promote diversity, we propose to maintain an ensemble of preference estimates, such that each of the estimates can be used to generate a new set of policies that align with that estimate. For learning the preferences, each comparison is costly, so we focus user queries on Pareto-optimal policies that capture non-dominated trade-offs across the features and that are likely to improve our preference estimate. After a fixed number of queries, the highest-utility policy under the final estimate is returned to the user.

**Contributions.**   In summary, compared to standard preference learning approaches, PASTEL introduces several key components: (1) an ensemble of preference estimates to maintain diversity of the policy population; (2) an evolutionary algorithm (EA) that generates new interpretable policies guided by the current preference estimates; (3) a preference model that represents policies as feature vectors, enabling generalization across similar policies; (4) a Pareto frontier filtering step that focuses the query selection on non-dominated policies; and (5) an informed query selection strategy that selects comparisons expected to be most informative for refining the preference model. We now describe each component in detail.

### 6.3.1 Preference Modeling Over Interpretable Policies

When searching for the best interpretable policy, the space of candidate policies is often large but structured. For example, a shallow decision tree might achieve optimal performance with coarser decision boundaries, while a deeper decision tree achieves this same performance with more fine-grained rules. Evaluating and estimating utility for each policy independently ignores these shared attributes. As a result, it is desirable to use this information to share information *across* policies. We propose to represent each policy using a feature vector that encodes key interpretable attributes to enable a learning algorithm to generalize across policies by learning a preference model over features, rather than over discrete policies.

**Featurize policy.** Let $\hat{\Pi}$ denote the set of all possible policies of a specific form and $\mathcal{E}$ represent the set of all possible environmental or contextual information. Let $\phi : \hat{\Pi} \times \mathcal{E} \to \mathbb{R}^d$ be a function that encodes a policy $\hat{\pi}$ and the environmental information into a $d$-dimensional feature vector $\mathbf{f}_{\hat{\pi}} \in \mathbb{R}^d$. Because users evaluate policies by examining observable properties, like the structure and resulting behavior, we assume that the feature mapping $\phi(\cdot)$ captures all aspects of the policy that contribute to the user preference evaluation. In other words, in this work, we assume that the mapping is known and sufficient to describe the users' preferences, such that,

$$\mathbf{f}_{\hat{\pi}} = \phi(\hat{\pi}, \mathcal{E}). \tag{6.3}$$

These features may include structural attributes of the policy, performance metrics, and explanation characteristics, which must be well-defined computations on the interpretable policy.

**Reparameterize utility model.** We now need to handle the user's utility function. Observe that, given $\mathbf{f}_{\hat{\pi}}$ and $\hat{\boldsymbol{\theta}}$, we can simply rewrite $u(\hat{\pi}; \hat{\boldsymbol{\theta}})$ as a linear combination of the policy features, such that:

$$u(\hat{\pi}; \boldsymbol{\theta}) = u(\phi(\hat{\pi}, \mathcal{E}); \boldsymbol{\theta}) = \mathbf{f}_{\hat{\pi}}^{\top} \hat{\boldsymbol{\theta}}. \tag{6.4}$$

This reparameterization preserves the original linearity assumption, so we operate in a condensed feature space without changing the model class. An appropriate choice of $\phi(\cdot)$ can capture potentially non-linear and non-differentiable preferences over policies. Then, we derive the recommendation as:

$$\hat{\pi}_M = \mathrm{argmax}_{\hat{\pi} \in \hat{\Pi}} \hat{u}_m(\hat{\pi}; \hat{\boldsymbol{\theta}}_M), \tag{6.5}$$

where $\hat{\boldsymbol{\theta}}_M$ is the estimate at iteration $M$ and $\hat{\Pi}_M$ is the set of policies considered up to iteration $M$.

**Benefits.** This approach brings two benefits. First, modeling this shared structure enables inference of the quality of previously-unobserved policies from observed ones with similar features. This generalization property enables more sample-efficient learning, which is ideal when each query represents an online human annotation request. Second, these features consist of human-understandable properties, meaning we can analyze trade-offs among the features throughout the learning process. For example, if a loan approval policy has features for approval rate and fairness across demographics with learned weights $\hat{\boldsymbol{\theta}} = [0.3, 0.8]$, respectively, then the company prioritizes fairness over maximizing approvals.

---

**Algorithm 6** PASTEL

---

1: $\hat{\boldsymbol{\theta}}_0^1, \ldots, \hat{\boldsymbol{\theta}}_0^N \leftarrow$ INITIALIZEESTIMATES()                    *// Initialize ensemble of preference estimates*
2: $\pi^* \leftarrow$ TRAINRLEXPERT()
3: $\Pi_0 \leftarrow$ GENERATEINITIALPOPULATION($\pi^*$)
4: $\mathbf{f}_{\hat{\pi}} \leftarrow$ EXTRACTFEATUREVECTORS($\Pi_0$)
5: **for** $m = 1, \ldots, M$ **do**
6:     $\Pi_{m-1} \leftarrow$ PARETOFILTER($\Pi_{m-1}$)                    *// Filter Pareto-dominated options*
7:     $q_{m-1} \leftarrow$ SELECTQUERY($\Pi_{m-1}, \hat{\boldsymbol{\theta}}_{m-1}^1, \ldots, \hat{\boldsymbol{\theta}}_{m-1}^N$)                    *// According to Equation* (6.6)
8:     $o_{t-1} \leftarrow$ QUERYUSER($q_{m-1}$)
9:     $\hat{\boldsymbol{\theta}}_m^1, \ldots, \hat{\boldsymbol{\theta}}_m^N \leftarrow$ UPDATE($\hat{\boldsymbol{\theta}}_{m-1}^1, \ldots, \hat{\boldsymbol{\theta}}_{m-1}^N, q_{m-1}, o_{m-1}$)                    *// Aßs in Equation* (6.7)
10:     **for** each $\hat{\boldsymbol{\theta}}_m^i$ **do**
11:         $\Pi_m^i \leftarrow$ RUNEVOLUTIONARYALGORITHM($\hat{\boldsymbol{\theta}}_N^i, \Pi_0$)                    *// Generate policies using Algorithm 7*
12:     **end for**
13:     $\Pi_m \leftarrow$ AGGREGATEPOPULATIONS($\{\Pi_m^i\}_{i=1}^N$)
14: **end for**
15: Calculate utility $\hat{u}_M$ under each $\hat{\boldsymbol{\theta}}_M$
16: **return** $\hat{\pi}^* = \arg\max \Pi_M$ according to $\hat{u}_M$                    *// Use ensemble to vote on best policy*

---

## 6.3.2 Preference Elicitation and Learning

We now introduce our preference elicitation and learning process, designed to learn user preferences over decision-tree policies. Algorithm 6 provides an overview of our algorithm. At a high-level, LICORICE works by maintaining an *ensemble* of preference estimates, which are used to guide the creation of new candidate policies for querying users. To focus queries on more informative policies, we *filter* the set of candidate policies by eliminating those that are dominated in all dimensions. To support this goal, we *select* queries using a novel selection strategy that aims to maximize information gained in parameter space. Finally, we *update* the ensemble of preference estimates using the user feedback with the goal of moving the preference estimates closer to the ground-truth preference. We now describe each of the aforementioned steps in more detail.

**Maintain ensemble of preference estimates.** We aim to initialize a diverse set of guesses about the true user preference to enable efficient exploration of the preference space (Algorithm 6, line 1). Inspired by the Query by Committee paradigm in active learning [287], we use an *ensemble* of preference estimates $\hat{\boldsymbol{\theta}} = \{\hat{\boldsymbol{\theta}}^1, \ldots, \hat{\boldsymbol{\theta}}^N\}$. Each $\hat{\boldsymbol{\theta}}$ can be thought of as a hypothesis about the user preference. Intuitively, multiple guesses about the true preference enables us to use prior knowledge with initialization. We also leverage these diverse guesses by coupling each estimate with its own EA, such that the output of the EA maximizes the current guess. Furthermore, we use these guesses during the policy selection step, in which we employ a voting mechanism where each estimate $\hat{\boldsymbol{\theta}}^i$ contributes equally to the decision, breaking ties uniformly at random (Algorithm 6, line 16). Over time, because these estimates are trained on the same data, they will eventually converge on the same policy, but the initial diversity promotes exploration.

**Filter data by the Pareto frontier.** Given that users have limited time, we want to ensure that the feedback we gain from each comparison is informative. To maximize the value of each user interaction, we avoid presenting policies that are dominated in all dimensions, as repeated selections of clearly superior policies provide little information about how users trade off different feature dimensions. Instead, we propose filtering the policy set based on the Pareto frontier (Algorithm 6, line 6). Given a dataset $D$ in a multi-dimensional space and a corresponding preference direction for each dimension, we identify the points that comprise the Pareto frontier:

$$\mathcal{P} = \{x \in D \mid \nexists y \in D, \text{ such that } y \text{ dominates } x\},$$

where a point $y$ dominates $x$ if it is at least as good in all dimensions and strictly better in at least one dimension. PASTEL then chooses from this reduced set of trees for preference elicitation. This approach not only ensures that users are presented with only the most promising options but can also dramatically reduce the computational complexity of the selection process. In many practical cases, the size of the Pareto frontier grows sublinearly with respect to the size of the dataset $k$. For instance, in two-dimensional problems, it follows a logarithmic relationship [21]. This reduces the complexity from the naive $O(Md^2)$ to $O(M \log^2 d)$, so we can explore the use of more computationally intensive techniques for identifying the most informative pairs for comparison.

**Select queries for feedback using uncertainty.** We introduce a novel approach to identify the most informative query, a pair of items $(\mathbf{f}_{\hat{\pi}_j}, \mathbf{f}_{\hat{\pi}_k})$ from a set $\mathcal{F}$, given a collection of preference estimates $\{\hat{\boldsymbol{\theta}}^i\}_{i=1}^N$. We seek to identify the examples for which we have the most uncertainty over the entire ensemble. Specifically, in Algorithm 6, line 7, we compute hypothetical gradient updates for both possible outcomes for each parameter. As a result, for each potential pair and each $\hat{\boldsymbol{\theta}}^i$, the algorithm computes two hypothetical updates:

$$\hat{\boldsymbol{\theta}}^i_{(j)} = \hat{\boldsymbol{\theta}}^i - \eta \nabla \mathcal{L}(\mathbf{f}_{\hat{\pi}_j} \succ \mathbf{f}_{\hat{\pi}_k} | \hat{\boldsymbol{\theta}}^i) \text{ and}$$
$$\hat{\boldsymbol{\theta}}^i_{(k)} = \hat{\boldsymbol{\theta}}^i - \eta \nabla \mathcal{L}(\mathbf{f}_{\hat{\pi}_k} \succ \mathbf{f}_{\hat{\pi}_j} | \hat{\boldsymbol{\theta}}^i),$$

where $\eta$ is the learning rate and $\nabla \mathcal{L}$ is the gradient of the logistic loss. We then compute the cosine similarity $\cos(\hat{\boldsymbol{\theta}}^i_{(j)}, \hat{\boldsymbol{\theta}}^i_{(k)})$ for each $\theta^i$ to measure the degree of alignment in the parameter space. When the cosine similarity is low, the two potential updates point in opposite directions, indicating more uncertainty about the outcome. Consequently, the algorithm selects the pair $(\mathbf{f}_{\hat{\pi}_j}, \mathbf{f}_{\hat{\pi}_k})$ that minimizes the average cosine similarity across all $\hat{\boldsymbol{\theta}}^i$:

$$(\mathbf{f}_{\hat{\pi}_j}, \mathbf{f}_{\hat{\pi}_k}) = \underset{(\mathbf{f}_{\hat{\pi}_j}, \mathbf{f}_{\hat{\pi}_k})}{\arg\min} \frac{1}{N} \sum_{i=1}^N \cos(\hat{\boldsymbol{\theta}}^i_{(j)}, \hat{\boldsymbol{\theta}}^i_{(k)}). \tag{6.6}$$

This approach aims to maximize the expected information gain by choosing pairs that lead to the most orthogonal (i.e., least similar) updates across the ensemble of preference estimates.

**Update preference estimate with feedback.** We have not yet discussed how we leverage the collected preference data to update the preference estimates. We adopt a logistic regression model with stochastic gradient updates, such that we can update our estimates in an online manner. In each iteration $m$, we receive (noisy) feedback $o_m$ about $q_m$ from the user based on Equation (6.1). This model induces a likelihood over observed pairwise comparisons, where more preferred items should receive higher predicted utility. Maximizing this likelihood leads directly to the logistic regression objective. As a result, PASTEL (Algorithm 6, line 9) updates each $\hat{\boldsymbol{\theta}}$ as:

$$\hat{\boldsymbol{\theta}}_{m+1} = \hat{\boldsymbol{\theta}}_m + \alpha \cdot \left( o_m - \sigma \left( \hat{\boldsymbol{\theta}}_m^\top (\mathbf{f}_{\hat{\pi}} - \mathbf{f}'_{\hat{\pi}}) \right) \right) \cdot (\mathbf{f}_{\hat{\pi}} - \mathbf{f}'_{\hat{\pi}}), \tag{6.7}$$

where $\alpha$ is the learning rate, $\sigma(\cdot)$ is the logistic function, $o_m \in \{0, 1\}$ is the realization of the binary label derived from the user's feedback, and $\mathbf{f}_{\hat{\pi}} - \mathbf{f}'_{\hat{\pi}}$ is the difference between the feature vectors of the compared decision-tree policies. Observe that, although we update each $\hat{\boldsymbol{\theta}}^1, \ldots, \hat{\boldsymbol{\theta}}^N$ using the same preference data, the initial guesses of these parameters enables diversity in the preference space.

### 6.3.3 Interpretable Policy Generation

We now focus on generating new policies to align with each member of the ensemble. The goal is to create these structures such that they can be checked against (and possibly added to) the current Pareto frontier for further elicitation. Because many popular interpretable policy structures, such as decision trees, are not differentiable, our algorithm must produce new decision trees without relying on gradients.[1]

We introduce a novel module that leverages the power of evolutionary algorithms (EAs) to generate new decision-tree policies. EAs [14] do not require the objective function to be differentiable. They work by evolving a population of solutions over generations to optimize a given objective function, called the fitness function. This approach is called in Algorithm 6, line 11, and we provide the additional details of this subroutine in Algorithm 7. Full details are in Appendix 1.

**Generate policies with EA.** In PASTEL, we call the EA for each $\hat{\boldsymbol{\theta}}^1, \ldots, \hat{\boldsymbol{\theta}}^N$, such that the fitness function for the associated EA evaluates the quality of each individual in the population with $\hat{\boldsymbol{\theta}}^i$. Each EA runs for $G$ generations to iteratively modify the policy population. Within a generation $g$, each policy $\hat{\pi} \in \Pi_g$ is scored by the utility $\mathbf{f}_{\hat{\pi}}^\top \hat{\boldsymbol{\theta}}_m^i$ (line 4). The next population $\Pi_{g+1}$ is first formed by *selection*, retaining parents with probability proportional to fitness. Then, new policies are formed using *crossover* (recombining selected pairs to exploit complementary structure) and *mutation* (randomly perturbing offspring to sustain diversity). These mechanisms can be thought of as balancing exploitation of known high-quality policies with exploration of promising regions of the policy space. The procedure then outputs $\Pi_m^i = \Pi_G$, from which the highest-fitness policy can be extracted. The new policy is compared with the set of decision-tree policies in the Pareto frontier and added to the set if it is non-dominated.

---

[1] When we can update with gradients, one can directly leverage work that uses preference feedback for RL [15, 57].

---

**Algorithm 7** RunEvolutionaryAlgorithm

---

1: **Input:** Preference vector $\hat{\boldsymbol{\theta}}_m^i$, initial population $\Pi_0$, crossover rate $\rho$
2: **Output:** New population $\Pi_m^i$
3: **for** $g = 0, \ldots, G - 1$ **do**
4:     Evaluate fitness of each policy in $\Pi_g$ using $\mathbf{f}_{\hat{\pi}}^\top \hat{\boldsymbol{\theta}}_m^i$
5:     Initialize next population $\Pi_{g+1}^i \leftarrow \emptyset$
6:     **while** $|\Pi_{g+1}^i| < |\Pi_g^i|$ **do**
7:         $(p_1, \mathbf{f}_{p_1}), (p_2, \mathbf{f}_{p_2}) \leftarrow$ SelectParents$(\Pi_g^i)$
8:         Generate a random number $a \in [0, 1]$
9:         **if** $a < \rho$ **then**
10:             $p_1, p_2 \leftarrow$ Crossover$(p_1, p_2)$
11:         **end if**
12:         Mutate$(p_1)$, Mutate$(p_2)$
13:         $\mathbf{f}_{p_1}, \mathbf{f}_{p_2} \leftarrow$ ExtractFeatureVectors$(p_1, p_2)$
14:         Add $(p_1, \mathbf{f}_{p_1}), (\mathbf{f}_{p_2}, p_2)$ to $\Pi_{g+1}^i$
15:     **end while**
16: **end for**
17: **return** $\Pi_m^i = \Pi_G$

---

**Expedite run-time.** Because we envision this system running in real-time with real users providing preference feedback, we introduce a method to speed up one of the more computationally-expensive parts of our method: estimating the reward of a policy during feature extraction (line 13). With a large number of candidate policies to precisely evaluate at each round of the EA, this estimation step quickly becomes burdensome. As a result, we use a UCB-inspired bound [301] to adaptively pick policies for roll-out. Furthermore, notice that each per-estimate EA run is completely independent, meaning we can enjoy the benefit of parallelization, and wall-clock time per call of the EA is calculated by the per-estimate EA run that takes the longest.

## 6.4  Experimental Setup

We evaluate the performance of PASTEL on two different RL environments. Because decision trees are generally considered interpretable, we conduct a *functionally-grounded* evaluation [82], in which we simulate user preferences over features. A crucial advantage of testing our method in these simulated environments is that we can evaluate how well our model is able to recover the preferences of the "ground truth" users, which provides us with insights about how our methods could perform with real users. In our experiments, we seek to answer the following questions:

  **RQ1** Preference Alignment: Does PASTEL produce more preference-aligned trees? Does it do so in a query-efficient manner?

| CartPole-v1 | | | PotholeWorld-v0 | | |
|---|---|---|---|---|---|
| Feature | Type | Values | Feature | Type | Values |
| Reward | Performance | $[0, 500]$ | Reward | Performance | $[-170, 49.95]$ |
| Depth | Structural | $\{0, \dots, 10\}$ | Depth | Structural | $\{0, \dots, 10\}$ |
| Num. leaves | Structural | $\{1, \dots, 1024\}$ | Num. leaves | Structural | $\{1, \dots, 1024\}$ |
| Feature used | Explanation | $\{0, 1\}$ | Action taken | Explanation | $\{0, 1\}$ |

Table 6.1: **Features used in the preference vectors for the two environments.** Here, "Feature used" corresponds to the inclusion of a specific feature in the state space as a feature that is tested in the internal node of the DT.

**RQ 2** Human-AI Interaction: What benefits could PASTEL have in terms of user interactions?

**RQ 3** Noise Robustness: How robust is PASTEL to preference noise?

**RQ 4** Ablation Study: Which components of PASTEL contribute most to its performance?

**Environments.** We select CartPole-v1 and PotholeWorld-v0 as our test environments not only because a reward-optimal policy can be represented as a tree but also for their complementary characteristics. CartPole-v1, a classic control problem, offers a simple, well-understood domain with a low-dimensional state space (cart position, velocity, pole angle, and angular velocity) and binary actions (push left or right), serving as an excellent benchmark for basic RL capabilities. In contrast, PotholeWorld-v0 [321] presents a more complex, driving-inspired scenario where an agent navigates lanes to avoid obstacles. We augment PotholeWorld-v0 with a controllable parameter governing lane reward trade-offs. In PotholeWorld-v0, the state is simply the current position, with three possible actions corresponding to lane choices. As a result, the agent must learn through experience to navigate around the potholes.

**Features.** For our feature vectors (computed as in Equation (6.3)), we test with $d = 4$ features for each environment. To evaluate the method, we choose at least one feature from each type (performance-based, structural, and explanation characteristics). Table 6.1 shows the features used.

**Preference generation.** To generate the values for the underlying user preference $\boldsymbol{\theta}$, we use a vector scaling technique [87]. Specifically, we linearly combine the basis vectors, with each vector maximizing only one preference. This setup enables the exploration of explicit trade-offs among these preferences by adjusting the scaling factor for each vector entry $\alpha$, such that $\sum_{i=1}^{d} |\boldsymbol{\theta}_i| = 1$. The scaling factor is 0.25, resulting in 35 different $\theta$ values per environment. We run each $\boldsymbol{\theta}$ with 3 random seeds for each experiment.

**Policy set initialization.** To create the initial policy population, we use a modified version of VIPER [17], an imitation learning algorithm for constructing reward-maximizing decision-tree policies by leveraging a neural network expert $\pi^*$. More details about VIPER in Section 5.2.2. To obtain a broader variety of trees, we introduce a depth randomization procedure and use all trees found instead of only reward-maximizing ones. More details of this initialization procedure are in Appendix 1.

**Baselines.** We compare with two baselines: VIPER [17] and Randomized Dueling Preference Selection (RDPS). Because VIPER does not incorporate preferences in its learning process, the best tree is chosen based on return only. The RDPS baseline is inspired by randomized dueling bandit algorithms used in preference-based learning. The algorithm maintains the current best policy. At each iteration $m$, the algorithm compares the current best with a randomly selected challenger. After eliciting pairwise preference feedback, the winner of the duel becomes the updated best explanation. This process is repeated for $M$ queries. Note that this algorithm does not attempt to estimate the ground-truth preference model, which could make it more susceptible to preference noise.

**Performance metrics.** Our main performance metric is the *alignment* score, which we compute using $\hat{u}m(\hat{\pi}; \boldsymbol{\theta})$ with respect to the ground-truth user preference. We normalize the alignment score using the maximum and minimum possible alignment scores of a policy of that structure. We obtain these values by running the EA for each ground-truth preference vector, $\boldsymbol{\theta}$, and its inverse, $-\boldsymbol{\theta}$, respectively. Doing so enables us to obtain the practical maximum and minimum alignment scores for that preference vector and environment. The normalized values are in $[0, 1]$.

## 6.5 Results

### RQ 1: Preference Alignment

To evaluate whether PASTEL produces better-aligned decision-tree policies, we first examine performance under an *easy* noisy regime by setting $\beta = 10$ in Equation (6.1). We assess policy alignment quality at the maximum training horizon $M$ across all environments. We then analyze two key dimensions over the training horizon: sample efficiency in terms of the user queries and the ability of PASTEL to recover the underlying user preference structure.

**PASTEL produces better aligned policies.** As we establish, the goal of the algorithm is to output the best possible policy after a predetermined number of $M$ queries. As a result, we evaluate the quality of the decision-tree policies that each approach recommends at $M = 40$ queries for PotholeWorld-v0 and $M = 20$ queries for CartPole-v1. As shown in Figure 6.2, PASTEL produces substantially more preference-aligned trees than both VIPER and RDPS in both environments. These findings show that PASTEL can indeed produce more preference-aligned decision-tree policies. This result indicates that personalizing policies to support end users in their goals and to align with their preferences is not only possible, but also not solved by existing techniques.

Figure 6.2: **Comparison of normalized alignment scores.** Error bars correspond to standard error. By explicitly incorporating preferences over policies and using those preferences to generate new policies, PASTEL creates more aligned interpretable policies than the baselines in both environments.

**PASTEL more efficiently produces better-aligned policies.** Because our goal is to not overburden a user, we aim to *efficiently* output highly-aligned policies. As a result, we also investigate the quality of the policies output by each algorithm with increasing $M$. Specifically, in Figure 6.3, we plot the quality over time, starting from $M = 1$ for both environments, and ending with $M = 40$ queries for PotholeWorld-v0 and $M = 20$ queries for CartPole-v1. In both environments, PASTEL generally more rapidly outputs higher-quality policies in fewer user queries. Interestingly, for PotholeWorld-v0, RDPS happens to choose strong initial guesses, but it struggles to improve the alignment of the policy with more queries. It is soon overtaken by PASTEL in around 8 queries. This result indicates that PASTEL can more effectively make use of the human feedback in terms of both policy selection and generation.

**PASTEL can generate more aligned policies because it is iteratively refining its preference estimate.** We now seek to understand what enables PASTEL to generate more aligned policies. To do so, we calculate the cosine similarity between the ground-truth preference vector $\theta$ and each of the ensemble members. Specifically, over successive preference queries in CartPole-v1 and PotholeWorld-v0, we first sort all ensemble members by the cosine similarity, then average along that axis across all runs and preference vectors. Figure 6.4 plots this result. In both domains every curve climbs monotonically: a handful of early comparisons eliminate blatantly wrong weightings, pushing even the least-similar estimate from around 0.30 to around 0.55 within only 10 queries, while the best estimate quickly exceeds 0.80. After this rapid phase the slopes flatten and the trajectories asymptote. This plateau is expected: once the algorithm has learned a *sufficient* direction to rank the current Pareto set of tree policies, additional queries offer diminishing discriminative power because the remaining candidate ensemble members differ only in fine-grained trade-offs that never alter which policy is chosen.

## RQ 2: Human-AI Interaction

We now illustrate the benefits of a few of the modeling choices introduced by PASTEL. We first show the envisioned interactions that users could have with the entire system. We then show how a few design choices contribute to the comprehensibility of the system.

Figure 6.3: **Normalized alignment scores with varying** $M$**.** Error bars correspond to standard error. Compared with the baseline, PASTEL more effectively utilizes a small number of human preference queries to produce more aligned interpretable policies. Interestingly, for PotholeWorld-v0, RDPS happens to choose strong initial guesses, but it struggles to improve the alignment of the policy with more queries.

**Modeling the factors that contribute to preference alignment enables better understanding of tradeoffs.**    One key distinguishing factor of PASTEL is that we represent each policy as a compact subset of features. This representation is in contrast to uninterpretable latent modeling approaches and approaches that bypass these factors to simply output a point estimate. For example, RDPS only outputs the current best policy, which does not enable the end user to intuitively understand the contribution of each feature to the policy that is selected. We visualize how a user can interact with the learned preference models and policy features in Figure 6.5. This figure shows an example from a PASTEL run in CartPole-v1 (an example for PotholeWorld-v0 is shown in Appendix 3). The top part of the figure showcases that the user can investigate which of the $\hat{\boldsymbol{\theta}}$ voted for each of the $\hat{\pi}$, in addition to the overall winner. The bottom part highlights that the user can further decompose the assessment into pairwise comparisons over features to see the directions of the $\hat{\boldsymbol{\theta}}$ for each combination.

Figure 6.4: **Similarity of the learned preference estimates to the true preference.** We calculate the cosine similarity of the 3 learned preference estimates in the ensemble for each environment. We do so by first ordering them in terms of least to most similar, then averaging over runs along those axes. This plot shows that, in both environments, PASTEL is indeed refining and improving its estimate of the true user preference.

**Implementation details enable our system to potentially be used in real time.** Because we envision that this system could potentially be used to elicit from downstream users, we prioritize efficiency in our approach. As a result, each EA run only takes around 1 minute of wall clock time. This speedup is because we implement not only algorithmic details to enable more efficient calculations but also parallelize where possible. Consequently, with well-spaced parallelized runs of the EA, each individual run could be completed in less than 10 minutes. As a point of comparison, RDPS baseline takes around 1 minute to complete 1 run of 20 queries. In contrast, PASTEL would take around 5 minutes for 20 queries; however, it results in substantially more aligned policies.

## RQ 3: Robustness to Preference Noise

We investigate whether PASTEL is robust to preference noise in CartPole-v0 for three different noise regimes: easy, medium, and hard. We achieve these regimes by setting $\beta$ in Equation (6.1) to values that result in progressively *more* probabilistic disparity between the best and the worst trees in the original set of trees. In other words, there exists a $1 - \epsilon$ probability of choosing the best policy over the worst policy. As we increase in difficulty, we set $\beta$ such that $\epsilon$ is closer to $0.5$ than in the easy setting. Specifically, we set $\beta = 10$ in the easy noisy regime, $\beta = 30$ in the medium regime, and $\beta = 50$ in the hard regime.

**PASTEL is more robust to preference noise.** Figure 6.6 shows that PASTEL finds better preference-aligned trees in all three noise regimes. Note that VIPER does not learn preferences and therefore is resistant to noise (it will deterministically pick the reward-maximizing option). Unsurprisingly, this resistance is beneficial in cases where more preference weight is placed on reward. Compared to RDPS, PASTEL enjoys both higher preference alignment and lower variance. In fact, in both the medium and hard regimes, VIPER selects policies that are more preference-aligned than RDPS. This result indicates that PASTEL is more robust to noise.

| $\widehat{\pi}$ | Reward | Depth | Num. Leaves | Uses Feature | Score $(\widehat{\theta}_1)$ | Score $(\widehat{\theta}_2)$ | Score $(\widehat{\theta}_3)$ |
|---|---|---|---|---|---|---|---|
| 1 | 500 | 2 | 3 | 1 | 2.18 | 3.11 | 2.52 |
| 2 | 217 | 1 | 2 | 0 | -0.73 | 0.29 | -0.24 |
| 3 | 9 | 0 | 1 | 0 | -1.01 | -0.13 | -0.69 |



Figure 6.5: **Example illustration of how one can leverage and inspect the learned preferences and attributes.** This figure shows an example from a PASTEL run in the CartPole-v1 environment. The top table shows the different values of the features (Reward, Depth, etc.) that correspond to each policy $\widehat{\pi}$ in the Pareto frontier. It also shows the overall score according to each $\widehat{\theta}$, as well as the overall $\widehat{\pi}$ that is selected as the best (policy 1, in green) and the ones that were not selected (policies 2 and 3, in pink). The three plots at the bottom show the directions of the $\widehat{\theta}$ for different combinations of feature comparisons.

## RQ 4: Ablation Study

Table 6.2 presents the results of our ablation study of the PASTEL algorithm. Recall that the PASTEL consists of four main algorithmic changes: the use of a custom EA to generate new policies, iteratively generating new policies and refining the preference estimates, filtering the set of policies using the Pareto frontier, and selecting user queries for feedback using uncertainty estimates. Respectively, the ablations correspond to: PASTEL-DTEA, PASTEL-ITER, PASTEL-PFF, and PASTEL-CSS. Overall, we find that different components of PASTEL have different, environment-dependent contributions to the overall performance. We now discuss each ablation in more detail.

**Generating new policies most impacts the overall performance.** We first discuss the effect of removing the EA component (PASTEL-DTEA). In other words, PASTEL-DTEA learns the underlying user preference by only querying from the policies in the initial generation. At the end of the $M$ iterations, PASTEL-DTEA simply outputs the best policy from the initial set. We find that, compared with PASTEL, PASTEL-DTEA results in a notable drop in performance—especially in CartPole-v1, where the score decreases from $0.99 \pm 0.00$ to $0.89 \pm 0.02$. This result indicates that the additional ability to *generate* new policies based on preference estimates most substantially impacts the ability to produce better-aligned policies.

Figure 6.6: **PASTEL creates and finds better aligned policies than the baselines in all three noise regimes.** Ideally, we want an algorithm that achieves a high normalized alignment score while minimizing the variance (standard error).

**Interleaving preference learning and policy generation helps in some cases.** The non-iterative variant (PASTEL-ITER) first exhausts the $M$ preference queries, then runs the EA. In contrast, PASTEL interleaves preference learning and policy generation. As a result, this ablation tests whether generating new policies throughout the algorithm is more useful than first learning the preferences and then using the preference estimate to guide the EA. PASTEL-ITER exhibits a significant performance degradation in PotholeWorld-v0; in contrast, PASTEL-ITER maintains comparable performance in CartPole-v1. In CartPole-v1, the initial VIPER–extracted tree set is already sufficient for learning a reasonable preference estimate, so expanding the set of candidate policies throughout the learning process is not necessary for extracting additional preference information. However, in PotholeWorld-v0, the initial VIPER-extracted trees provide poor coverage of high-quality regions; without interleaving, the system spends its entire feedback budget comparing sub-optimal policies, so the eventual EA search is guided by a potentially lower-quality preference model. By contrast, the alternating scheme continuously injects improved trees, steering both the preference model and subsequent evolutionary steps toward promising areas and yielding higher final utility.

**Filtering by the Pareto frontier provides a modest benefit.** We now discuss the impact of not filtering by the Pareto frontier (PASTEL-PFF). Compared with PASTEL, PASTEL-PFF leads to a slight alignment reduction in both environments, indicating that Pareto filtering contributes modestly but consistently to overall performance. In addition to the modest performance improvement, one of the benefits of this filtering process is enabling users to reason about a smaller subset of alternatives, as outlined in Section 6.5.

| Algorithm | PotholeWorld-v0 | CartPole-v1 |
|---|---|---|
| PASTEL | $0.98 \pm 0.00$ | $0.99 \pm 0.00$ |
| PASTEL-DTEA | $0.93 \pm 0.00$ | $0.89 \pm 0.02$ |
| PASTEL-ITER | $0.72 \pm 0.05$ | $0.99 \pm 0.00$ |
| PASTEL-PFF | $0.97 \pm 0.01$ | $0.96 \pm 0.01$ |
| PASTEL-CSS | $0.97 \pm 0.01$ | $0.99 \pm 0.00$ |

Table 6.2: **Ablation study: average normalized alignment scores for all PASTEL ablations.** Error is standard error. For PotholeWorld-v0, all components of PASTEL contribute to its performance. However, CartPole-v1 is simpler, so the major benefit stems from the EA to generate more aligned policies (PASTEL vs. PASTEL-DTEA).

**PASTEL is relatively robust to the choice of elicitation strategy.** We now ablate the choice of elicitation strategy. Specifically, we compare the designed preference elicitation strategy used by PASTEL with random sampling (PASTEL-CSS). Interestingly, PASTEL-CSS leads to only a minor performance drop in both environments compared with PASTEL. This robustness is particularly encouraging given that many preference elicitation methods in the literature employ a "warm start" approach, beginning with random queries to broadly explore the preference space before transitioning to more targeted, strategic querying [277]. PASTEL could benefit from such hybrid approaches in future work.

## 6.6 Related Work

**Interpretable machine learning.** Interpretable machine learning has become increasingly important as AI agents are deployed for more complex and sensitive use cases [46, 257]. The need for interpretability arises from the desire to understand, appropriately trust, and effectively manage these systems. The most relevant line of work to ours is in interpretable RL. Although many structures have been proposed to use in lieu of neural networks for RL policies, including domain-specific programming languages [334] and human-friendly prototypes [161], we focus on decision-tree policies due to the general consensus that they are human-understandable [187, 273]. One challenge with decision trees is their discrete structure is not immediately amenable to gradient-based training. Although previous work has applied gradient-based training to decision-tree policies by introducing non-linearities in the splits [298], the conversion of the soft tree [144] to a standard DT results in severe performance loss, making this approach inapplicable in our setting. Moreover, prior work trains decision-tree policies only from the perspective of reward maximization [17, 321]. Instead, our approach compares policies based on policy-level properties. Reward-based metrics are one class of properties, so our method aligns policies with both reward and policy-wide interpretability metrics.

**Learning from human feedback.** As we discuss in detail in Section 2.1.4, there is growing interest in the area of learning from human feedback [109, 121, 307] due to the difficulty of specifying a concrete objective that captures the full extent of what people want agents to do. Early work collects preferences over agent trajectories to train a reward model $r_\theta$, and the resulting cumulative reward $R(\tau; \theta) = \sum_t r_\theta(s_t, a_t)$ is used to model preference probabilities via the Bradley-Terry model [57]. Similar techniques have been applied to language models, where preferences over two *completions* $y_i, y_j$ of queries $x$ are used to train a reward model $R(x, y)$. The model is then fine-tuned using reinforcement learning [246] with the reward model. In contrast, we estimate a utility function over interpretable policies using policy-level features, and use it to actively guide the generation of new policies that align with user preferences. This utility function need not be decomposable into a sum over transition utilities.

**Preference elicitation and learning.** Preference elicitation and learning have been extensively studied across various disciplines [43, 64, 180, 191, 344]. The objective of most work in this area is to determine a full or partial ranking over alternatives by employing methods such as pairwise comparisons [29, 90, 194] or asking individuals to rank their top choices among all or a subset of alternatives [83, 195, 305, 364]. Our work builds on these works, but it shifts the focus from eliciting rankings to learning the underlying preference model parameters to aid in *generating* preference-aligned interpretable policies. Our approach shares similarities with Bayesian Experimental Design [44, 264], where the objective is to learn information about the underlying user preference. However, we integrate preference elicitation and learning with an emphasis on aligning models with the preferences of downstream users.

## 6.7 Discussion and Future Directions

This chapter proposes a novel alignment paradigm, decision-making alignment. Different from behavioral alignment, this paradigm focuses on building agents that employ decision-making processes that align with user preferences. Toward this goal, this chapter also advances both a new model and a novel learning algorithm, PASTEL. To scope our approach, we focus on representing policies in an inherently interpretable manner through decision trees, which means that users can interpret the reasoning process and provide feedback. From the interpretability perspective, we are also the first to incorporate user feedback directly into the interpretable policy generation process.

PASTEL contributes multiple algorithmic innovations with the goal of efficiently producing more aligned policies. To enable generalization across similar policies in a potentially-infinite policy space, PASTEL represents policies as feature vectors over critical attributes. To promote diversity of the policy population, it maintains an ensemble of preference estimates that guide the generation of new policies and the selection of policies for querying. To generate new policies, we develop an EA that is guided by the current preference estimates. Before querying users, PASTEL employs a Pareto frontier filtering step that focuses the query selection on non-dominated policies. Finally, PASTEL uses an informed query selection strategy that selects comparisons expected to be most informative for refining the preference model.

We demonstrate that our method not only produces decision-tree policies that better align with the underlying preferences but also does so efficiently, in 40 queries or fewer. We showcase how our approach could facilitate human-AI interaction through the more interpretable preference model and policy feature representation. We then show that PASTEL is more robust to preference noise, meaning it could be used in settings where users have more uncertainty about their alignment goals. Our work lays the groundwork for future research on developing more user-centered, interpretable RL agents whose reasoning process aligns with user preferences.

We also note that this general problem structure is relatively unexplored in the literature. With the advances in generative modeling, there may indeed come a time when movie recommendations can use user preferences to generate movies that align with these preferences. As a result, important questions around sample efficiency, learning from implicit signals, and more, are open areas of study.

## Future Directions

**Accommodating additional structures.**   Although we study our approach in the context of decision-tree policies, our framework is general enough to accommodate policies of various forms and architectures [129, 161, 334, 198]. The key insight that interleaving preference elicitation with policy evaluation can improve alignment extends beyond any specific policy representation. Through our design choices—including focusing on the pairwise comparison setting, efficient learning through query reduction and algorithm speed—we hope that our approach could be integrated into real settings where users must quickly converge on a policy for e.g., decision support.

**Navigating trade-offs for deployment.**   Our approach offers flexibility through two distinct deployment settings: 1) post-hoc policy selection, in which this estimate is used to choose among a pool of trained interpretable policies, and 2) in-the-loop policy training, where preference feedback is elicited and incorporated directly during policy learning. Post-hoc selection allows users to navigate tradeoffs across interpretable policies and choose one aligned with their values, without requiring modifications to the underlying training procedure. In-the-loop policy training enables simultaneous refinement of preference estimates and the generation of new policies to align with those estimates. However, there is a trade-off: as we have shown, the EA is critical for generating the more aligned policies. As a result, post-hoc selection may be faster; however, without a diverse initial pool of policies, users may not fully realize the alignment gains that could be made with the in-the-loop training method.

**Relaxing preference modeling assumptions.**   We assume static and linear user preferences, which may not accurately reflect preferences in real-world applications. Although these are common assumptions, future work could model and adapt to evolving user preferences over time. Furthermore, interpretability is inherently subjective; it can vary significantly depending on the user's background, experiences, and domain knowledge [365]. Although we aim to capture *task-specific* factors that influence a user's judgment of a policy, future work could incorporate *user-specific* factors.

**Involving real end users.** In future studies, moving from synthetic user feedback to incorporating a wider range of user input and adopting an *application-grounded* evaluation approach [82] could provide additional insights. For example, real-world data could inform the latent factors that influence user preferences. Involving end users would also be beneficial in designing the frontend of the system to enable users to interact with it in a way that is beneficial for them. The means of presentation of such policies is important, and understanding what works well for users in a particular context is critical before the deployment of a system [295].

**Expanding to large reasoning models.** Although in our work, we focused on producing inherently interpretable decision-tree policies, the basic ideas could be applied to work with large reasoning models. Recent work in reasoning with large language models—particularly via chain-of-thought prompting— focuses on generating human-interpretable rationales [346]. However, much of this research remains limited to post-hoc explanation, without guarantees that the reasoning traces reflect faithful internal computation or support reliable decision-making [323]. Furthermore, it remains unclear how to evaluate the quality of these reasoning traces, or how they should be optimized to support real-world use. One direction is to frame reasoning as an interactive process and optimize traces based on their downstream effects. For example, reasoning could be refined through human feedback in a similar manner to PASTEL.

# 7 Learning from Human Feedback for Fuzzy Tasks

## 7.1 Introduction

In the previous chapter, we establish *decision-making* alignment and provide an algorithm for aligning RL agent policies with human preferences. We show that it is possible to efficiently improve the alignment of the decision-making process of these agents by incorporating user feedback into the policy generation process. In many applications, however, we additionally care about *behavioral* alignment: how the agent acts in the world. This chapter turns to the problem of behavioral alignment, which is especially challenging in tasks where desirable behavior is hard to formalize. Many real-world tasks in which AI agents will be deployed have this property, making this area important to study.

However, traditional RL is not well equipped to handle problems with so-called fuzzy objectives. In traditional RL, an agent learns how to act using reward based on an explicitly-defined reward signal [311]. This reward signal is often carefully designed by domain experts to communicate the intended goal for the agent to accomplish. Precisely specifying this form of feedback programmatically requires designers to *a priori* enumerate all potential outcomes or constraints on how they would like the task to be completed. This enumeration is difficult to achieve in practice, and the resulting reward signals often fall short at correctly specifying the designer's intent [170]. To address this challenge, researchers have explored the idea of incorporating alternative channels for communicating information about the desired behavior of the agent. This class of techniques is generally called *learning from human feedback* (LfHF) [57, 147]. The goal of LfHF is to utilize the feedback modalities most likely to result in an agent acting according to human-desired specifications.

The complexity of this approach has been exemplified in our MineRL BASALT competition series [291], which is the first to evaluate fine-tuning from human feedback techniques in the sequential decision-making setting. BASALT leverages the open-ended and flexible Minecraft environment to promote the development of LfHF algorithms. TThis competition series has provided a platform for training and fine-tuning agents capable of solving fuzzy tasks that lack well-defined reward signals. However, despite two years of competition, no agents have yet matched human performance levels [213], emphasizing the need to establish the BASALT tasks as a consistent, standardized benchmark. In this paper, we clarify the BASALT benchmark and present concrete evaluation recommendations toward the goal of consistency.

Figure 7.1: **The BASALT challenge.** Parts shaded in green show the components of the included datasets. The goal of the challenge is to develop agents that can successfully learn to perform fuzzy tasks using human feedback. Toward this goal, we contribute human demonstrations data and human pairwise comparisons of people and AI agents attempting to complete BASALT tasks. We hope the comparison data not only streamlines assessments by enabling fewer comparisons to determine the ranking of a new agent on the leaderboard but also establishes the challenge of proper reward modeling for fuzzy tasks. We also contribute code to streamline further benchmarking.

To support this benchmark we introduce the BASALT Evaluation and Demonstrations Dataset (BEDD), an open and accessible dataset for learning to solve fuzzy tasks from human feedback. As shown in Figure 7.1, this dataset consists of three main ingredients: the Demonstrations Dataset, the Evaluation Dataset, and supporting code for utilizing and analyzing the data. The Demonstrations Dataset consists of over 26 million image-action pairs from 14,000 videos of labeled Minecraft gameplay of human players completing the BASALT tasks.

To facilitate the evaluation of AI agents using real human judgments, we present the Evaluation Dataset, derived from the most recent BASALT competition [154]. The Evaluation Dataset consists of over 3,000 dense pairwise human evaluations of videos of various agents performing the BASALT tasks. By dense, we mean that each evaluation includes a comparison of the relative task-completion performance of the agents and at least four additional questions, such as which agent was more human-like. This results in 27,905 comparison points between 17 different agents. This dataset also includes a natural language response justifying why an agent was selected as being the better of the two. The provided data functions as a leaderboard, offering researchers the ability to compare their newly-developed algorithm against various agents without redoing all costly human evaluations from scratch.

To help other researchers use the Demonstrations Dataset and the Evaluation Dataset, we present a streamlined codebase. With this codebase, one can train a new model from the demonstration dataset and evaluate it against the provided leaderboard. Alongside the presentation of these datasets, we conduct a detailed analysis of the data to guide algorithm development and evaluation. We release the code and detailed documentation for others to perform such analyses. We hope that this codebase will assist others with quickly developing and evaluating LfHF algorithms to spur further progress toward agents that are better aligned with human intent.

**Contributions.** In summary, we make the following contributions:

- We clarify the formalization of the BASALT benchmark, positioning it as a standardized testbed for learning from human feedback on fuzzy, reward-free tasks.

- We introduce the BASALT Evaluation and Demonstrations Dataset (BEDD), a new dataset containing over 26 million state-action pairs from human demonstrations and more than 3,000 dense human evaluations of agent behavior.

- We release supporting infrastructure for training and evaluating agents on the BASALT tasks, including tools for leaderboard evaluation, agent comparison, and analysis of behavioral alignment failure modes.

- We conduct an extensive empirical analysis to help researchers better understand alignment challenges in open-ended environments and use BEDD effectively in algorithm development.

## 7.2 The MineRL BASALT Benchmark

We first provide an overview of the MineRL BASALT benchmark, consisting of a task suite and evaluation framework for learning from human feedback, illustrated in Figure 7.1. We establish this benchmark in prior work [291] but include details here for completeness. The benchmark uses Minecraft, a videogame that provides a rich and complex environment in which to define different tasks. The states are pixel observations; the actions are regular keyboard and mouse actions, closely following how humans play the game. For example, the agent must navigate the crafting menus using a mouse (Figure 7.2).

This section presents the BASALT task specifications and evaluation methodology. We begin with task descriptions, followed by the evaluation protocol using TrueSkill ratings and human judgment. We then examine the challenge of developing automated evaluation systems to complement human assessment. Additional evaluation guidelines appear in Appendix 1.2.

Figure 7.2: **Two example images of the pixel observations provided by the MineRL environment used by BASALT.** The agent receives pixels as observations and must use regular keyboard and mouse actions, including navigating the crafting menu (right-hand screenshot) with a mouse to make items.

## 7.2.1  Tasks

Because our goal is to promote work on LfHF algorithms that can eventually be applied to hard-to-specify problems of interest in the real world, we aim to incentivize general solutions rather than algorithms that are specific to a given task. As a result, we choose Minecraft as our environment. Minecraft is extremely open-ended: people pursue a wide variety of goals within the game, and the game mechanics provide a hierarchical and combinatorial space of possible tasks that agents can be trained to perform. By designing tasks involving high-level concepts that cannot be easily identified programmatically, we can recreate the situation we often face in real-world settings: reward functions cannot be used simply because they are too hard to write down.

We design a suite consisting of four reward-free tasks and one rewarded task. The reward-free tasks are BuildVillageHouse, CreateVillageAnimalPen, MakeWaterfall, and FindCave ; the rewarded task is ObtainDiamondShovel. Because reward-free tasks are challenging to evaluate, we include the rewarded task [213] to enable quick iteration. All tasks are accompanied by a Gym environment [33] and either a simple English-language task description or a reward function to indicate the desired agent behavior. The *reward-free* tasks have the language task specification; the *rewarded* tasks are accompanied by a reward function. We now explain the tasks in detail.

**BuildVillageHouse.** In BuildVillageHouse, the agent spawns in a village, which is a naturally-occurring structure in Minecraft that contains houses and other buildings, villagers, and other structures. The style of the buildings depends on the biome in which the village is generated. The agent must build a new house for itself in the same style as the other village houses; however, it must not damage the village. To assist agents with achieving this task, we provide the agent with the materials required to build a house in a variety of village types.

**CreateVillageAnimalPen.**    In CreateVillageAnimalPen, the agent must build an animal pen next to an existing house. It must then corral a pair of farm animals (chickens, cows, sheep, or pigs) into the pen. The agent spawns in a village. In Minecraft, players perform this task to breed the animals to create a steady source of food. We provide the agent with materials to build the pen. Because different animals are lured with different types of food, we also equip the agent with the food to lure any of the animals.

**MakeWaterfall.**    To complete MakeWaterfall, the agent must create a waterfall and subsequently take a picture of it. Because taking a photograph is not a supported task in Minecraft, we implement it by having the agent throw a snowball. We interpret the moment that the agent throws it as its photograph. The challenge with this task is it involves both the evaluation of task completion and human aesthetic preferences. We provide the agent with the tools needed for efficiently moving around the often-treacherous biome, a snowball for taking a picture, and two buckets of water to construct the waterfall. The agent also spawns in an extreme hills biome to enable easier task completion.

**FindCave.**    The goal of FindCave is for the agent to discover a naturally-generated cave. The agent spawns in a plains biome and must explore the surrounding area to find a cave.

**ObtainDiamondShovel.**  This task is more challenging than ObtainDiamond [116], a well-established benchmark in RL, since it requires an additional crafting step after obtaining the diamond. An agent receives reward each time it obtains the next required resource in the crafting tree towards a diamond shovel.

## 7.2.2 Evaluation Protocol

The BASALT benchmark establishes a standardized evaluation framework that addresses key challenges in assessing AI agents on complex, open-ended tasks. BASALT agents undergo evaluation on designated hold-out test environments, generating multiple video demonstrations of task completion attempts. For each task, we then play agents against one another in matches. In a single match, we select two agents and one test environment, then show a human judge the two corresponding videos. The judge determines which agent performs the specified task more effectively. Given this dataset, we compute scores for each agent using TrueSkill [130].

**TrueSkill.**  We choose TrueSkill because, in addition to estimating the relative skill of an agent, it provides an uncertainty estimation. As a result, we also use TrueSkill to select matches, so that we can maximize the expected increase in information gained per match. To obtain a stable ranking of agents, one can simply run this protocol until the scores are reasonably confident. Using TrueSkill not only enables more efficient human query usage for evaluating new agents but also a principled evaluation framework for assessing BASALT agents. For a detailed discussion of TrueSkill vs. alternative approaches, please refer to Appendix 1.1 for a detailed discussion.

**Multi-task performance aggregation.** To determine overall performance across BASALT's multiple tasks, we normalize TrueSkill scores within each reward-free task before computing aggregate measures. For a score $x$ on a given task, we compute the normalized score as: $\frac{x-\mu}{\sigma}$, where $\mu$ is the mean and $\sigma$ is the standard deviation of the scores. (Note that here $\mu$ and $\sigma$ are *not* the means and standard deviations estimated by TrueSkill, but are rather the empirical mean and standard deviation calculated across participant scores.) We can then take the average of the normalized scores to determine an overall score. The normalization ensures that no one task unduly drives the final scores, while they remain sensitive to the magnitude of the differences in scores between different agents on a task. Unfortunately, this scheme has a major downside: if all submissions behave similarly on a task, then $\sigma$ will be very small, and tiny differences between agents can be amplified into large differences in overall score. To mitigate this issue, we require the denominator to be at least 1, such that a score of $x$ maps to $\frac{x-\mu}{\max(\sigma,1)}$.

### Toward Better Reward Models for Automated Evaluations

Although human evaluation provides the ground truth for BASALT assessment, the associated costs and time requirements present practical barriers for iterative algorithm development. Automated evaluation systems have shown promise across various machine learning domains [52, 77], offering rapid feedback during the development cycle. We do not propose automated evaluation as a replacement for human judgment. Rather, automated methods serve as a complementary tool that enables researchers to conduct preliminary assessments and iterate quickly on algorithm designs before committing to expensive human evaluation. This two-stage approach allows practitioners to filter and refine their approaches efficiently.

To support automated evaluation research, BASALT includes the ObtainDiamondShovel task, which provides a concrete reward function that can serve as training data for reward modeling techniques [141]. Additionally, we release the Evaluation Datasetcontaining human evaluation data to enable development of automated systems that can predict human preferences. In addition to the overall performance assessment, researchers could utilize the fine-grained assessments, as well as the natural language justifications. The development of reliable automated evaluation remains an open research challenge. We anticipate that improved automated methods will emerge as researchers apply machine learning techniques to predict human evaluative judgments, potentially using the human evaluation data provided through BASALT as training supervision.

## 7.3  Benchmarking AI Agents on BASALT

To assist with the development of LfHF algorithms and automated evaluations, we implement and share a codebase with two major contributions. First, the code contains an example of training a LfHF algorithm with the shared data in the Demonstrations Dataset. Second, we include the tools for performing the evaluations presented in this work. The code is a Python-installable library, which allows the functionality to be imported into other codebases for use in research.

**Training example.**   The training example provides tools to train a behavior cloning model on top of the Video PreTraining (VPT) [16] model using the `imitation` [103] library. VPT is a large foundation model that can complete various tasks in Minecraft, but it is difficult to fine-tune on new tasks due to its size. Inspired by the success of an imitation-learning approach[1] in the BASALT 2022 competition, we use behavior cloning as the base algorithm with the rich embeddings from the VPT model as input. Motivated by the original VPT results, we remove the no-op actions from the demonstration dataset. The code supports different VPT model variants for benchmarking. The final output is a set of videos to use in evaluations against the shared recordings of other agents in the Evaluation Dataset.

**Human evaluation platform.**   Setting up human evaluations requires significant effort. As a result, we also share our platform for conducting human evaluations. The platform runs as a webpage with a Python webserver backend. The collected data is stored either locally in a SQLite database or remotely for larger-scale deployments. The platform includes a flexible API to add or remove agents from the set of comparisons. After a simple setup, researchers can point human evaluators to a specific URL to provide answers. We provide examples of the form that the human evaluators will see in Appendix 4.

**Analysis tools.**   We release code for generating TrueSkill rating figures and all analyses from Sections 7.5 and 7.6 using the demonstrations and evaluation datasets. This process provides an end-to-end workflow: train a method with Demonstrations Dataset, then evaluate it using Evaluation Dataset.

## 7.4  BASALT Evaluation and Demonstrations Dataset

We now introduce BEDD, our extensive dataset of human demonstrations and algorithm evaluations. This dataset consists of the following components:

- The Demonstrations Dataset, a set of 13,928 videos (state-action pairs) demonstrating largely successful task completion attempts of the reward-free tasks,

- The Evaluation Dataset, a set of 3,049 dense pairwise comparisons of algorithmic and human agents attempting to complete the BASALT tasks, and

- The code for utilizing and analyzing these datasets for developing LfHF algorithms (some details in Section 7.3).

We introduce Demonstrations Dataset, then Evaluation Dataset. For a full datasheet [101], please see Appendix 2.

| Task | Videos | Episodes | Hours | Size | Episode length, s | Success % |
|------|--------|----------|-------|------|-------------------|-----------|
| FindCave | 5,466 | 5,466 | 91 | 165GB | 60 | 93% |
| MakeWaterfall | 4,230 | 4,176 | 97 | 175GB | 84 | 98% |
| CreateVillageAnimalPen | 2,833 | 2,708 | 89 | 165GB | 119 | 95% |
| BuildVillageHouse | 1,399 | 778 | 85 | 146GB | 391 | 92% |
| Total | 13,928 | 13,128 | 361 | 651GB | 99 | 95% |

Table 7.1: **High-level demonstration data statistics decomposed by task.** For each task, we report the number of videos (Videos), the number of episodes (Episodes), the corresponding number of demonstration hours (Hours), the size of the data (Size), the average episode length in seconds (Episode Length, s), and the success rate. A demonstration counts as a success if the player manually ended the episode instead of dying or timing-out. Each task has at least a 92% success rate, meaning that the demonstrations are of high quality overall.

## 7.4.1 Demonstrations Dataset: Completing Reward-Free Tasks

The Demonstrations Dataset for developing new methods consists of 361 hours (26 million image-action pairs) of human demonstrations of the reward-free BASALT tasks. This data consists of labeled trajectories, both with high-resolution image observations and keyboard and mouse actions for each frame. In total, this is 651 GB of data. Table 7.1 decomposes the high-level data statistics by task. More details about this dataset are in Appendix 3.

**Demonstration structure.** Each demonstration consists of a trajectory $\tau = [s_1, a_1, \ldots, s_T, a_T]$, or a sequence of state-action pairs, where $T$ is the trajectory length. These pairs are contiguously sampled at every Minecraft game tick (20Hz). Each state consists of the $640 \times 360$ RGB frame from the perspective of the player (see Figure 7.2). Each action consists of two parts $a = [K, M]$, where $K$ is all keyboard interactions, $M$ is all "mouse" interactions (change in view, pitch, and yaw), mimicking the native human control interface of Minecraft. This dataset serves as a starting point for using *demonstrations* as a form of feedback to train agents.

## 7.4.2 Evaluation Dataset: Evaluating BASALT Agents

The Evaluation Dataset contains 3,049 pairwise comparisons of different algorithms, produced from 273 hours of human labeling effort by 65 unique MTurk workers. This dataset enables the construction of leaderboards for comparing LfHF algorithms across multiple performance dimensions. All evaluation data is distributed as a single JSON file. Table 7.2 decomposes this dataset by task.

---

[1] https://github.com/shuishida/minerl_2022

| Task | Comparisons | Hours | Words in Response |
|------|-------------|-------|-------------------|
| FindCave | 722 | 60 | 27,948 |
| MakeWaterfall | 682 | 56 | 26,437 |
| CreateVillageAnimalPen | 914 | 81 | 32,768 |
| BuildVillageHouse | 731 | 76 | 26,917 |
| Total | 3,049 | 273 | 114,070 |

Table 7.2: **High-level Evaluation Dataset statistics decomposed by task.** We report the total number of agent-agent comparisons (Comparisons), human labor hours (Hours), and total number of words used in the natural-language justifications of selecting a specific agent as the best one (Words in Response).

**Dataset structure.**   Each evaluation entry in the dataset contains six core components: (i) the names of the two agents used in the comparison, (ii) the corresponding videos shown to the human judge, (iii) an answer of which player is better overall (Left, Right, Draw), (iv) a natural-language justification explaining the judge's reasoning behind this choice, (v) answers to at least one task-specific direct question about concrete achievements by the players (Left, Right, or Both), and answers to 4-7 task-specific comparative questions, such as which agent was more human-like (Left, Right, Draw, N/A).

**Dataset scale and scope.**   When aggregating across all evaluation dimensions, including the primary performance judgment, the dataset encompasses 27,905 individual comparisons. This multi-dimensional approach captures aspects of algorithmic performance that aggregated automated metrics often miss, as they tend to only assess the overall performance. With this choice, we aim to provide a more complete understanding of agent capabilities.

**Quality control and data filtering.**   We implement systematic quality control procedures. The dataset includes documentation of anonymized MTurk worker identifiers whose responses failed to meet established standards, such as providing identical responses across all tasks or demonstrating insufficient engagement with evaluation criteria. Prior to conducting the analyses presented in this paper, we filter the data to exclude these responses. However, we provide all responses in our public data release.

**Agents included.**   The dataset includes 17 existing approaches: the top 13 teams from the 2022 BASALT competition, a behavioral cloning baseline, a random agent providing a naive performance reference, and two human experts (authors of this paper) establishing human-level performance benchmarks. This diverse collection allows researchers to position new developments within the broader landscape of state-of-the-art approaches by leveraging the existing comparisons to establish relative performance rankings.

**Reproducibility.**   We provide all details needed to reproduce these evaluations in Appendix 4. By releasing the data and code, we aim to reduce the time and financial costs associated with large-scale human evaluation studies while maintaining the benefits of human-centered performance assessment.

# 7.5 Analysis: Demonstrations Dataset

We now analyze the Demonstrations Dataset. Because the BASALT tasks lack concrete reward functions, evaluating the progress of agents on these tasks is challenging. This difficulty necessitates the use of informative proxy measures. As a result, when analyzing the Demonstrations Dataset, we focus on defining proxy measures that may be useful for understanding the data or tracking training progress. We include these results both as a way to understand this dataset and as values to monitor during agent training to estimate agent performance. However, it is crucial to avoid using these metrics as direct optimization targets, as they can be easily exploited. In this section, we describe these proxies and present the results of our analysis.

## 7.5.1 Analysis

**Characterizing the difficulty of the reward-free tasks.** We first seek to understand the relative difficulty of the four reward-free tasks. Given the experience of the contracted data collectors with Minecraft, we believe that the length of the demonstration is a reasonable proxy for task difficulty. Using demonstration length as a proxy for task difficulty, we note that BuildVillageHouse is likely the most challenging task, even for humans: each video lasts around 6.52 minutes on average, while the next most time-consuming task, CreateVillageAnimalPen, takes an average of around 1.98 minutes (Table 7.1). By this metric, the easiest task is FindCave (1 minute). In practice, one could use this metric to assess training progress or likelihood of task completion. If episodes terminate at a rate that deviates substantially from the average, that may signal worse-quality behavior. In contrast, if an agent takes a similar amount of time to the average to complete a task, then that could signal better behavior (but is clearly not definitive without additional signals, such as lack of death of the agent).

**Determining when an agent could be underperforming.** We also want to understand when an agent—human or AI—may be underperforming on the task. Because FindCave requires navigation to find a cave, an agent that remains stationary is likely unsuccessful. Similarly, because CreateVillageAnimalPen and BuildVillageHouse require the construction of objects, agents that fail to place any blocks likely do not succeed at the task. As a result, we employ the number of steps with an active movement key as a proxy for the distance traveled within the tasks. We also use the count of right mouse button clicks as a proxy for the number of blocks placed. The only other actions performed by clicking the right mouse button are using a crafting table or a chest, of which there are relatively few per episode. These proxies are depicted in the per-episode distributions within the Demonstrations Dataset in Figure 7.3. The number of right mouse button clicks is highest at around 180 in BuildVillageHouse and between 20-40 in MakeWaterfall and CreateVillageAnimalPen. These align with expectations, given the amount of building required to complete the tasks. The number of movement actions has the highest variability in FindCave. This task also has the shortest time-out at 3,600 steps. Perhaps a better proxy would be the average number of movement actions per step: FindCave should have the highest number, as the task mainly consists of moving around.

(a) Right mouse button clicks

(b) Movement key presses

Figure 7.3: **The distributions of the number of right mouse button clicks and movement key presses across tasks in the Demonstrations Dataset.** These metrics act as proxies for the number of blocks placed and distance traveled per episode, respectively. Because BuildVillageHouse takes the longest to complete, it is reasonable that it requires more right mouse button clicks and movement key presses.

## 7.6  Analysis: Evaluation Dataset

We now present an analysis of the Evaluation Dataset. We first analyze the overall dataset, then focus on a few main comparisons: the top two algorithms from the BASALT 2022 competition (GoUp and UniTeam), the behavioral cloning baseline (BC-Baseline), the two human experts (Human1 and Human2), and a random agent (Random). This subset of data corresponds to 394 of 3,049 total comparisons or nearly 34 hours of human labeling effort.

### 7.6.1  Overview of AI Agents

We begin by briefly describing the two top-performing AI agents included in the evaluation: GoUp and UniTeam. The remaining agents include a behavioral cloning baseline (BC-Baseline), a random policy, and two human demonstrators (Human1, Human2).

**First-place: GoUp.**   GoUp leverages the insight that all of the tasks consist of the same flow—the agent walks around, searches for a target, then solves the task.[2] The approach trains several classifiers and object detection models to identify the targets in each task. To source training data, they manually label images collected from the expert videos provided by the competition. Taking the AnimalPen task as an example, the solution contains a fine-tuned VPT model for moving the agent, a fine-tuned YOLOv5 [326] detector for detecting the types and locations of animals, a fine-tuned MobileNet [138] for identifying the location of fence placement, and a finite state machine to control the flow of all the components.

---

[2]Open-source code for GoUp: `https://github.com/gomiss/neurips-2022-minerl-basalt-competition`.

| Agent | FindCave | MakeWaterfall | AnimalPen | BuildVillageHouse | Average |
|---|---|---|---|---|---|
| GoUp | 0.31 | 1.21 | 0.28 | 1.11 | 0.73 |
| UniTeam | 0.56 | −0.10 | 0.02 | 0.04 | 0.13 |
| Human2 | **2.52** | **2.42** | 2.46 | **2.34** | **2.43** |
| Human1 | 1.94 | 1.94 | **2.52** | 2.28 | 2.17 |
| BC-Baseline | −0.43 | −0.23 | −0.19 | −0.42 | −0.32 |
| Random | −1.80 | −1.29 | −1.14 | −1.16 | −1.35 |

Table 7.3: **Leaderboard: normalized TrueSkill scores.** We bold the player (human or AI) with the highest score for each task, as well as the highest score overall. Human players substantially outperform the best-performing AI agent across all tasks.

**Second-place: UniTeam.** This solution involves a search-based behavioral cloning approach, which aims to reproduce an expert's behavior by copying relevant actions from relevant situations in the demonstration dataset [201] .[3] They define a situation as a subset of consecutive frames and actions within a recorded trajectory, which they encoded using the VPT network. To find the most similar situations to the current one, they use a pre-trained VPT model to produce latent representations. They search for the nearest embedding point in the VPT latent space to find the reference situation. They assume that retrieved situations represent an optimal solution to a specific past situation. As a result, they simply copy the corresponding actions. After each timestep, they update both the current and reference situations. Because the reference and current situations diverge over time, they measure the $L_1$ similarity between the situations at each timestep. Once the distance between trajectories exceeds some fixed threshold (or after 128 timesteps), they perform a new search.

## 7.6.2  Performance Analysis

Here, we seek to understand whether AI agents can successfully complete fuzzy tasks at a level that matches or exceeds human players. As a result, we first present the overall performance of the agents. Then, we turn to the additional task-specific questions to gain a deeper understanding of which facets of the tasks AI or human agents find challenging. We note that this analysis is conducted over the 3,049 total overall comparisons.

**There is still a large gap between overall human and AI agent capabilities for fuzzy tasks.** We present the agent-based overview of the evaluations in Table 7.3. The human agents substantially outperform the AI agents. In particular, there is a large gap between both human players and the best-performing AI agent both overall and when decomposed by task. This result emphasizes the difficulty of BASALT for AI agents.

---

[3]Open-source code for UniTeam: `https://github.com/fmalato/basalt_2022_submission`.

(a) FindCave

(b) MakeWaterfall

(c) BuildVillageHouse

(d) AnimalPen

Figure 7.4: **Comparison of baseline solutions, top BASALT 2022 competition solutions, and humans on additional questions.** Bars represent average score, and error bars represent standard error. Higher is better: it means that agent exhibited more of that factor according to human judges.

**The tasks that are the most challenging for AI and people may differ.** We now leverage both our proxy measures from the analysis in Section 7.5 to evaluate whether there is any trend of difficulty with respect to which tasks AI agents perform worse on and which humans find challenging. Interestingly, for both agents, there seems to be little to no correlation, as the GoUp agent receives similarly low scores for the easiest and hardest tasks for the human demonstrators. However, we note that these are all proxies, so we caution against any strong conclusions.

**People still outperform AI agents along all factors.** We now investigate the responses to the direct and comparative questions. For succinctness, we summarize each question into a single factor. For example, we use "More Human Like" to capture the question, "Which player seemed more human-like (rather than a bot or computer player)?". Figure 7.4 presents these results. Importantly, no AI agents outperform people on any of the factors, indicating that there is still plenty of improvement to be made. When analyzing future algorithms, we suggest preserving this decomposition to both validate that the responses are sensible and obtain a fine-grained understanding of agent capabilities. To validate sensible responses, we suggest checking the scores for the *random* agent: except for causing the least harm, this agent should score poorly compared to the other agents.

| Task | Length | | Sentiment | | |
|------|--------|--|-----------|--|--|
| | Characters | Words | Positive | Neutral | Negative |
| FindCave | $210.46 \pm 3.82$ | $38.71 \pm 0.71$ | 79.64% | 6.65% | 13.71% |
| MakeWaterfall | $217.96 \pm 4.18$ | $38.76 \pm 0.79$ | 76.10% | 7.33% | 16.57% |
| CreateVillageAnimalPen | $197.83 \pm 3.22$ | $35.85 \pm 0.60$ | 56.67% | 10.83% | 32.49% |
| BuildVillageHouse | $205.73 \pm 4.23$ | $36.82 \pm 0.77$ | 62.79% | 9.58% | 27.63% |

Table 7.4: **Details about the justification provided for choosing a particular agent as the best one, decomposed per task.** We calculate the per-task average length of the response (characters and words) along with the standard error. We also report the percent of positive, neutral, and negative sentiments of these responses. This is a more detailed view of the Sentiment columns in Table 7.3.

**Decomposing the evaluation criteria enables interesting analyses beyond overall performance.** Finally, we highlight a few interesting findings that are enabled by our decomposed and extensive set of evaluation criteria. Figure 7.4b (MakeWaterfall) reveals that, although Team GoUp's algorithm can create waterfalls at a rate more similar to the human players, it struggles along all other criteria, including choosing a good location and taking a high-quality photograph. Only looking at the binary success/failure condition for creating a waterfall would ignore these important nuances in the algorithm's behavior. As another example, Figure 7.4a (FindCave) suggests that, while Team GoUp's algorithm still struggles to find caves, it can reasonably search for and navigate to areas that are likely to have caves. This finding suggests that the performance bottleneck may be the cave detection system employed by this approach.

### 7.6.3  Natural Language Reasoning Analysis

Recall that, in addition to providing a binary choice of the preferred agent, the human judges also justify their decision using natural language. We envision the utility of this data as coming from two directions: i) we can more deeply understand human perceptions and preferences regarding fuzzy task completion and ii) we can potentially leverage this data (if high-quality) as a data source to feed back into AI agents to improve performance. In this section, we analyze these responses. We investigate both the lengths of the responses and the sentiments. Specifically, to understand the general perceptions of the human judges, we categorize each response into positive, neutral, or negative sentiment [192]. We present task-level and agent-level views of both of these analyses.

**Task-Level Analysis**

We now present the *task-level* analysis of the natural language rationale. Table 7.4 presents the overall results of the length and sentiment analysis, which is computed over the entire dataset. After discussing these in detail, we provide a qualitative examination of sampled justification texts.

| Comparison | Chi-square |
|---|---|
| FindCave vs MakeWaterfall | 2.69 |
| **FindCave** vs CreateVillageAnimalPen | 98.49 |
| **FindCave** vs BuildVillageHouse | 52.31 |
| **MakeWaterfall** vs CreateVillageAnimalPen | 66.37 |
| **MakeWaterfall** vs BuildVillageHouse | 30.50 |
| CreateVillageAnimalPen vs BuildVillageHouse | 6.35 |

Table 7.5: **Chi-square test results comparing sentiment distributions between tasks.** Significant differences in sentiment distribution are emphasized with a highlight . The task with the more positive sentiment is shown in bold. Overall, human judges exhibit more positive sentiment when evaluating FindCave and MakeWaterfall .

**Justification length is generally consistent across tasks, with some exceptions.** The human judges generally respond using a similar number of words across the different tasks, meaning that the judges are fairly consistent. On average, responses for FindCave and MakeWaterfall are slightly longer than those for CreateVillageAnimalPen , while BuildVillageHouse falls in between. CreateVillageAnimalPen contains more quantitative additional questions than the other tasks (4 vs. 2, the next highest). For example, the human judges identify which agent penned the correct amount and type of animals, among other things. These additional questions may enable the judges to spend less effort explaining their choice since many of the factors were already captured by other questions.

**Sentiment differs substantially by task and aligns with task complexity.** The human judges display different levels of sentiment in their responses, depending on the task.[4] Table 7.5 shows the results. Tasks such as FindCave and MakeWaterfall elicit the highest rates of positive sentiment (approximately 77–80%), while CreateVillageAnimalPen and BuildVillageHouse receive a substantially higher proportion of negative sentiment (27–32%). These differences likely reflect perceived difficulty: more complex tasks that involve multi-step construction and stricter criteria are more likely to result in agent failure and evaluator frustration. Judges appear more positive when agents perform well on shorter, more constrained tasks.

**When analyzing by task, shorter responses tend to accompany more negative evaluations.** Combining these two observations, we find that tasks with more negative sentiment also tend to elicit shorter justifications. For example, CreateVillageAnimalPen , which receives the lowest proportion of positive sentiment, also has the shortest responses on average. This pairing may reflect a shift in annotator behavior: when agents fail visibly or clearly, judges may feel less need to elaborate. In contrast, agents that make partial progress may prompt longer, more detailed commentary.

---

[4]We conduct a chi-square test of independence to examine the relationship between task and sentiment classification. The relation between these variables is significant, $X^2(6, N = 3049) = 132.21, p < .001$.

| | Length | | Sentiment | | |
|---|---|---|---|---|---|
| **Agent** | Characters | Words | Positive | Neutral | Negative |
| Random | $196.80 \pm 4.93$ | $35.41 \pm 0.92$ | 63.20% | 8.31% | 28.49% |
| Human1 | $207.06 \pm 5.66$ | $37.50 \pm 1.05$ | 91.55% | 2.11% | 6.34% |
| Human2 | $214.87 \pm 6.11$ | $38.86 \pm 1.12$ | 91.77% | 2.88% | 5.35% |
| BC-Baseline | $201.72 \pm 5.30$ | $36.25 \pm 0.97$ | 64.96% | 9.00% | 26.03% |
| GoUp | $212.86 \pm 5.55$ | $38.67 \pm 1.02$ | 73.77% | 7.53% | 18.70% |
| UniTeam | $201.77 \pm 4.85$ | $36.35 \pm 0.92$ | 65.68% | 9.38% | 24.94% |

Table 7.6: **Details about the justification provided for choosing a particular agent as the best one.** We calculate the per-agent average length of the response (characters and words) along with the standard error. We also report the percent of positive, neutral, and negative sentiments of the responses.

**Justification examples.** We now present an example of justification text provided for each of the tasks. We choose these examples by splitting the data by task, then randomly sampling a response. For FindCave, an example justification is,

> *They explored a lot, and seemed to go to more new places than before even the water. Left player went AFK.*

An example justification response from MakeWaterfall is,

> *The player on the left couldn't even manage to climb the side of a mountain, the one on the right didn't finish the tasks but at least they seemed to know how to navigate and move around.*

An example justification text from CreateVillageAnimalPen is,

> *the left player completed all the requirements, while the right player did nothing and still behaved like a bot*

Finally, an example justification text from BuildVillageHouse is,

> *The right player was better overall, being able to build a house, but he failed to use the appropriate type of block.*

These justifications demonstrate that judges are not only paying attention to high-level task outcomes but are also referencing specific aspects of Minecraft gameplay such as exploration, building mechanics, and task criteria. Natural language feedback in such evaluations encode fine-grained signals that could be leveraged for behavior modeling or diagnostic training signals.

| Comparison | Chi-square |
|---|---|
| **GoUp** vs UniTeam | 6.37 |
| **GoUp** vs BC-Baseline | 7.50 |
| **GoUp** vs Random | 10.44 |
| GoUp vs **Human1** | 34.10 |
| GoUp vs **Human2** | 31.22 |
| UniTeam vs BC-Baseline | 0.15 |
| UniTeam vs Random | 1.33 |
| UniTeam vs **Human1** | 62.97 |
| UniTeam vs **Human2** | 56.94 |
| BC-Baseline vs Random | 0.60 |
| BC-Baseline vs **Human1** | 64.77 |
| BC-Baseline vs **Human2** | 58.76 |
| Random vs **Human1** | 68.25 |
| Random vs **Human2** | 62.44 |
| Human1 vs Human2 | 0.53 |

Table 7.7: **Chi-square test results comparing sentiment distributions between agent pairs.** Highlighted rows indicate significant differences. The agent with significantly more positive sentiment is shown in bold. Overall, human judges display significantly more positive sentiment when one of the two human players is involved in a comparison. The sentiment is also roughly correlated with performance (although not significantly), with sentiments toward GoUp being generally more positive than to UniTeam, BC-Baseline, or Random.

## Agent-Level Analysis

We now present the *agent-level* analysis of the natural language rationale. Table 7.6 presents the overall results of the length and sentiment analysis. We compute these values using only the comparisons with the other presented AI agents, not the full dataset.

**Justification length is generally consistent across agents, with some exceptions.** The human judges generally respond using a similar number of words across the different agents, meaning that the judges are fairly consistent. On average, the MTurk workers use both the most words and characters (38.86 words, 214.87 characters) when describing their rationale when the comparison involves Human2 (who is also the most competent overall). In contrast, when the comparison includes the Random agent (who is also the least competent overall), the MTurk workers use the least words and characters (35.41 words and 196.80 characters) when explaining their rationale. We believe that this may be due to the relative competency of the agents: because it is easier to identify the Random agent as less skilled and the humans as more skilled, the human judges may require less justification for their selection.

**Sentiment differs significantly only in comparisons involving human agents.** Any comparisons that include either of the human agents, except for when they are pitted against one another, exhibit significant differences in sentiment distribution. In particular, the sentiment toward comparisons that include human players are significantly more positive[?]. We present the results of the statistical analysis in Table 7.7. These judgments also correlate with stronger task performance, as the responses with the highest positive sentiment are for the two human agents (who perform the best overall) and the least positive sentiment are for the Random agent (who performs the worst overall). In contrast, no significant sentiment differences emerge between any pair of AI agents, regardless of differences in performance.

**Human judges are generally positive when evaluating the agents.** Overall, the sentiment across all responses skews positive. This result could be due to the nature of the task: Minecraft players may fundamentally enjoy watching agents play Minecraft, even if they do so poorly. This result could also be due to social influence, since the judges are aware that the videos are produced by different players for ranking in the competition.

**Shorter responses tend to accompany more negative evaluations.** Similar to the task-based decomposition, the length and sentiment of the responses positively correlates. For example, Random, which receives the lowest proportion of positive sentiment, also has the shortest responses on average. Additionally, this agent also exhibits the lowest performance overall. These results are intuitive: when agents fail clearly, judges may feel less need to elaborate and their justification may focus more on what the agent did wrong rather than right. In contrast, agents that make even partial progress may receive more detailed commentary on their successes.

## 7.7 Related Work

**Learning from human feedback.** Learning from human feedback has become a crucial research direction in machine learning [11, 57, 291], aiming to leverage human expertise to improve algorithm performance and generalization. Various algorithms have been proposed to integrate human feedback into the learning process, often using techniques such as imitation learning [136], inverse reinforcement learning [35, 350], and reward modeling [307]. Most commonly, these algorithms are evaluated in standard RL benchmarks [18, 61, 85, 315]. However, most of these benchmarks do not have the property that human feedback is *critical* for task identification, as in BASALT. In some Atari games, if an agent does anything other than the intended gameplay, it dies and resets to the initial state, so pure curiosity-based agents perform well [36]. In contrast, BASALT tasks require human feedback or data to identify and complete tasks.

**Minecraft for machine learning.** Minecraft is a valuable platform for machine learning research [124, 186, 302] due to its complex, dynamic environment and customizable game engine. As a result, numerous competitions and benchmarks [108, 122, 149, 278] have been developed to improve AI capabilities. These often include tasks with concrete or programmatic reward functions, such as learning to collect diamonds in a sample-efficient way [116], multi-agent learning of cooperative and competitive tasks [254], and more [94]. In contrast, we emphasize tasks that are designed to be hard to specify through a reward signal. Other work focuses on natural language as a specific modality for communicating intent [107, 165, 166, 235]; instead, we aim to promote the development of algorithms that generally learn from human feedback, including natural language. Some techniques require complete game state information (knowledge of blocks and items around them) [340, 366]; however, our benchmark only permits agents to access observations, not game state, to promote learning from pixels, which is more generalizable to other tasks.

**Comparison to MineDojo.** Perhaps the most similar benchmark to ours is MineDojo [94], which emerged during the BASALT competitions and contributed a large dataset of Minecraft gameplay with the aim of developing generally-capable embodied agents. The focus of MineDojo was to provide a massive dataset scraped from the internet; in contrast, we focus on curating a smaller set of high-quality demonstrations and evaluations. Our demonstration data was produced by experienced Minecraft data collectors using a consistent Minecraft version and settings. In contrast, the MineDojo data, although plentiful, contains many videos with streamer overlays and non-Minecraft parts, different texture packs, mod packs, and more. Although this diversity may be beneficial in some settings, previous work had to filter out such data before it could be useful for training [16]. Another critical difference is evaluation: MineDojo provides only binary success or failure criterion; instead, our recommended evaluation and resulting Evaluation Dataset involves human judgments across a range of quantitative (e.g., Found Cave) and qualitative (e.g., Style Matching) criteria. Finally, BASALT is the first to establish a benchmark and associated competition for training and benchmarking agents that leverage human feedback via fine-tuning.

## 7.8 Discussion and Future Directions

We propose BEDD, a large and accessible dataset to facilitate algorithm development for the BASALT benchmark on learning from human feedback. This benchmark is the first to evaluate fine-tuning from human feedback for embodied agents. Our dataset, BEDD, consists of two key parts: Demonstrations Dataset and Evaluation Dataset. To our knowledge, this dataset is the largest of its type: one that supports both training and evaluating LfHF agents on tasks with hard-to-specify reward functions. Our results suggest that there is ample room for improvement in learning from human feedback. With the inclusion of our accessible code for benchmarking and evaluating agents, we hope that our contributions will encourage the development of more effective approaches for both learning from human feedback and evaluating these techniques in the future.

## Future Directions

**More detailed investigation of human factors.**   In the assessment of the agents, there are even more fine-grained details and important questions about how to leverage this information to further improve agent capabilities. As one example, we could leverage the rationales to provide additional information to the agent about which behavior is preferred [359]. As another example, the additional questions could be integrated to train agents that are both human-like and performant. Continuing with the human likeness example, future work investigating human likeness may want to substitute some of the non-human-like comparative questions with more detailed questions about human likeness to gain a deeper understanding of assessments of human likeness in this setting. This analysis could also be combined with recent work on the Automated Navigation Turing Test (ANTT) [77] and our work in Chapter 3 on the Human Navigation Turing Test (HNTT) but applied to a much more complex setting: completing fuzzy tasks in Minecraft.

**Toward capable and aligned AI agents for real-world tasks.**   Our hope is that this is just the first step along a path to general, capable AI agents that are aligned with human intent. Through this benchmark and datasets, we aim to promote research advancing LfHF techniques, rather than reward functions, to solve sequential decision-making tasks. We see the potential for real-world impact as coming from three main sources. First, such techniques can enable us to build AI agents for tasks without formal specifications. Second, LfHF methods greatly expand the set of properties we can incorporate into our AI systems, allowing more effective governance of AI systems. Finally, by learning from human feedback, we can train AI systems that optimize for *our* goals, making progress on the value alignment problem.

## Acknowledgments

Part IV

# DISCUSSIONS AND CONCLUSIONS

# 8 Discussion

While AI is rapidly advancing in terms of capabilities, many of these advances do not incorporate the human or sequential context in which agents will act. This thesis makes several contributions of new AI methodology that prioritizes those considerations, addressing problems that emerge from current or future deployment of such agents.

In particular, we make the following contributions. Toward the goal of building agents that exhibit more intuitive behavior, we develop a human-like agent and conduct an extensive empirical investigation of human perceptions, revealing findings that could be generalizable to developing agents that exhibit intuitive behavior more broadly. We then turn to interpretability. We contribute multiple algorithms for training agents with interpretable policies to facilitate deployment in domains with medium to high risks. Specifically, we uncover a critical limitation in training concept-based interpretable policies for RL: the reliance on human-in-the-loop concept annotations. We develop the first algorithm that addresses this problem, enabling the training of performant and interpretable concept-based policies with substantially fewer human labels. We also develop the first algorithm for training interpretable decision-tree policies that promotes coordination to maintain high performance on multi-agent tasks. Finally, we make several contributions to alignment. We introduce a new alignment space (decision-making-space) and provide an algorithm that achieves better-aligned policies in this new space. Finally, we show how current techniques are insufficient to train agents to exhibit aligned behavior for hard-to-specify tasks.

Together, these contributions reflect a broader vision: building AI agents from a human-centered perspective. This perspective acknowledges that real-world deployment introduces challenges that are often overlooked in standard formulations: objectives are hard to specify, feedback is noisy and diverse, and success depends not only on performance but on whether agents behave in ways that are intuitive, interpretable, and aligned with human expectations. By rethinking how we build and study agents from the viewpoint of human-centered RL, this thesis demonstrates how algorithmic and evaluative choices can be shaped by the demands of real-world deployment.

To showcase the breadth of opportunities for human-centered RL, we discuss directions for future work that extend beyond our three core pillars of creating intuitive, interpretable, and aligned agents. We then discuss challenges that emerge across all areas of human-centered RL. We then highlight bidirectional opportunities for RL problems to inform human-centered AI questions (RL for HCAI) (e.g., how do people's perceptions of AI agents unfold over multiple steps?) and for human-centered challenges to inform RL agent design (e.g., personalizing policies for diverse users) (HCAI for RL). These steps will help ensure that AI agent developments are in service of, and not in opposition to, human flourishing.

# 8.1 The Landscape of Human-Centered RL

Although this thesis focuses on developing agents that are intuitive, interpretable, and aligned, human-centered RL encompasses many additional dimensions that merit investigation. We distinguish these dimensions as specific properties or outcomes that we desire from our agents. In the following, we elaborate on several such properties, including robustness, controllability, privacy preservation, and usability.

**Robustness.** AI agents must maintain desirable behavior when encountering out-of-distribution situations stemming from varying users and edge cases. In the context of human-AI interaction, these distribution shifts may differ from those assumed in related settings, like adversarial robustness. Adversarial robustness typically assumes that attackers aim to cause misclassification (or missteps) of the agent [37], but human-centered robustness assumes that users are generally well-intended. This assumption is not always true, but we categorize handling malicious use as more of a security issue [76], which is a separate but important topic. As a result, we may need a new language and set of assumptions to make progress toward this goal.

**Controllability.** RL agents can already surpass human performance in some domains [123, 299, 300]. As this domain space broadens, the standard assumption that humans can easily supervise their behavior breaks down. We now must grapple with a key question: how can humans retain meaningful oversight and control of agents they do not fully understand? Achieving this kind of controllability presents various challenges. For example, safe interruptibility [244] ensures agents do not develop incentives to resist interruptions [120] but assumes idealized conditions for supervision. Similarly, approaches like learned intervention policies or blocking models [280] reduce human burden but tend to generalize poorly in novel states. These issues are amplified in settings where agents must explore, as these actions may elicit unnecessary user interventions.

**Privacy preservation.** The personalization of agents is a double-edged sword. On one hand, personalized agents are better-suited to take actions on behalf of users. On the other hand, accumulating large amounts of personal data presents privacy concerns. User tolerance of data collection to achieve greater degrees of personalization may be context-specific [47]. However, despite concerns around privacy, people in practice still tend to trade privacy for convenience [99], suggesting that policy-level interventions and guidelines may be necessary to navigate these trade-offs.

**Usability.** Agents must also be usable. Current approaches stress the importance of learning effective prompting techniques, but these ignore that different users may require different interfaces, such as elderly individuals [69]. And most studies of usability focus on chat interfaces [48, 27] As a result, these investigations are limited in terms of both the scope of users and interaction mechanisms. We may need to take a counterintuitive approach to ensure that the tools are first aligned with users. Another approach is to open up the communication channel, enabling agents to provide examples, questions, and other forms of feedback ("By X, I thought you meant Y") to create a shared language toward usability.

## 8.2 Common Challenges Across Human-Centered RL

As we establish, human-centered RL is a broad and interdisciplinary field, encompassing goals like transparency, alignment, collaboration, and robustness. Despite this diversity, several foundational challenges consistently emerge. A core challenge is *objective specification*: human goals are fuzzy, contextual, and difficult to codify precisely. In response, agents must increasingly rely on *multi-modal learning* to integrate partial and ambiguous signals across modalities like language, vision, and actions. Furthermore, open-ended objectives and diverse user needs make *scalability* a core challenge for deploying these agents in practice. To make such agents usable and steerable at scale, we also need better *tooling* that allows users and developers to provide feedback, inspect behavior, and intervene when necessary. Finally, *evaluation and benchmarking* remain open and underappreciated challenges: standard metrics often fail to capture human-centered objectives, and benchmarks quickly saturate without reflecting interaction quality, generalization, or long-term alignment with human intent. We now examine each of these challenges in greater depth, highlighting their unique difficulties and the opportunities they present for advancing human-centered RL.

**Objective specification.** We have seen the failure modes of misspecified objectives. Users become addicted to social media, LLMs generate convincing but inaccurate text. For example, some argue that recommender systems cause harm, such as by exacerbating political polarization and promoting extremist content [222] (though such concerns have been disputed [243]). These concerns persist despite efforts by companies to change recommender systems to optimize subjective wellbeing [223, 308], suggesting that the issue is that we do not know how to optimize for subjective wellbeing. Learning from human feedback provides one solution by circumventing the issue of precise objective specification; however, we must grapple with a key trade-off. Richer feedback can provide more information, but only if we have already defined the dimensions along which such input is expressed. This challenge introduces opportunities to rethink optimization objectives from a human-centered perspective and to leverage implicit feedback signals without requiring people to directly articulate them directly.

**Multi-modal learning.** Current AI agents perceive the world through a limited scope. However, people learn and provide feedback through various modalities [248]. To teach a dance, we might physically demonstrate a move; to promote deeper reflection, we might ask open-ended questions; to learn math, we might practice problems that share the same underlying concept but promote generalization. These different modalities of interaction are shaped by both the task and the communicative strengths of the person involved. To support human-centered goals, AI agents must be able to leverage not only rich inputs in the form of image, video, and audio, but also implicit and explicit human signals that occur over extended interactions within the context of sequential decision-making. Advances in multi-modal learning provide the tools needed to combine these diverse input types to support more naturalistic human-AI interaction, enable generalization across contexts, and make better use of the signals people already provide [247].

**Scalability.**    How do we ensure that human-centered approaches remain effective as agents become more capable and deploy at larger scales? To integrate diverse and multi-modal feedback that aligns with varying objectives, we need to develop efficient methods to collect human feedback across diverse populations, automated evaluation techniques that preserve human-centered principles, and systems that can generalize across different cultural contexts and user groups. There exist many techniques for eliciting human feedback in the HCI literature. These techniques, such as think aloud protocols, have remained understudied in the AI literature [58]. For example, LLM summarization techniques could distill the main points automatically and feed back into agents. These studies can be conducted across different contexts, tasks, and populations as people increasingly interact with these agents at different scales.

**Tooling.**  As agents become more complex and interactive, a central challenge is designing tools for end users and model developers that facilitate novel interaction modalities and scalable oversight. For example, users may want to isolate particular behavioral instances and intervene on that behavior [150], which raises exciting questions about scalability of interpretability techniques to models with billions and trillions of parameters. Furthermore, there exist many properties that an RL agent may estimate (e.g., world models [333, 349]), which a user may want to understand. Tools that enable interrogation of these components during human-AI interaction would help support human understanding of AI decision-making by offering a more complete picture. Model developers would also benefit from large-scale debugging tools that enable better attribution of failure cases or enable them to trace back failures in the execution stack.

**Evaluation and benchmarking.**    Current benchmarks for RL often focus on narrow performance metrics that do not capture human-centered properties. Furthermore, the domains used vary widely, spanning toy control [362], games [241], real healthcare data [26], and more. The usage of such disparate domains makes comparing algorithms challenging, especially when human-centered properties are often context-dependent. To make progress toward concrete and context-dependent tasks, we believe it is critical to evaluate the quality of agents using benchmarks with standardized environments and metrics. Although traditional machine learning benchmarks suffer from important issues, such as construct validity with respect to the problem being addressed [265], we believe that benchmarks, and competitions based off of these benchmarks [61, 116], can propel research forward if appropriate scope and focus is provided to an important, carefully-selected problem. For these, it may be necessary to collaborate with domain experts and HCI researchers to ensure that the benchmark is appropriately defined and scoped. Much of HCI work argues that systems must be studied in context with real users. But many studies of AI do not involve agents and, instead, involve studying how people interact with LLMs or chatbots. On the other hand, RL benchmarks make no consideration for the real context in which these agents will act. As a result, there is a need to collaboratively design simulators and human studies that more accurately capture the real-world dynamics and mechanisms for which agents will behave, as well as studying more RL-specific properties to help create a more overarching picture of these agents in context.

## 8.3 Bidirectional Opportunities for HCAI and RL

Having outlined key goals and challenges in human-centered RL, we now discuss the significant opportunities to integrate the perspectives of human-centered AI (HCAI) and RL, with benefits flowing in both directions. Although these fields often address overlapping problems, they frequently do so with different assumptions, priorities, and language. Bridging these perspectives opens up space for new questions and solutions. For example, what HCI researchers call user mental models—how users understand a system's behavior—can inform how RL researchers think about policy transparency, while RL's formalism for sequential decision-making offers tools for modeling complex human-AI interaction. We identify a few opportunities here. First, we discuss how taking an RL perspective can be useful for making progress on human-centered challenges, especially with the proliferation of AI agents for real-world applications. Second, we discuss how human-centered questions can benefit and support algorithm design in RL.

### 8.3.1 RL for HCAI

Many core problems in human-centered AI (HCAI) can benefit from an RL perspective. While much of HCAI has focused on static systems or one-shot interactions, RL focuses on agents that make decisions over time, adapt based on feedback, and influence their environment through extended interaction. This view introduces opportunities to expand the HCAI literature to address challenges that arise from how learning agents will behave in real-world settings.

One central challenge is understanding the bidirectional nature of human–AI interaction. This dynamic setting raises new research questions: how do interaction patterns evolve, when do users trust or intervene, and how does agent behavior shape user expectations? However, current studies often focus on narrow domains—such as coding assistants—and lack tools for investigating these questions across diverse tasks and modalities. Moreover, existing RLHF methods typically rely on simplistic and constrained models of user behavior and feedback. But, to improve alignment, we need mechanisms that help users give feedback that is both more expressive and better aligned with their true preferences, which we can support through interface design. In what follows, we discuss the aforementioned opportunities where RL can inform and extend HCAI.

**Bidirectional human-AI agent interaction.**   Because RL agents change over time, we must move away from the assumption of static systems and, instead, incorporate and disentangle user learning effects from system learning effects. More broadly, an interesting direction is studying how people form expectations and intuitions about the *learning process* of such agents. What mental models do they form about feedback and its ability to shift behavior? Exploration further complicates this dynamic: RL agents may explore suboptimal actions to obtain better performance in the future [171], which users could misinterpret as incompetence or misalignment. Studying these interactions in RL environments can help HCI researchers better model, support, and design for agents that evolve alongside their users.

**Scalable tools for analyzing human-agent interaction.** People increasingly are interacting with agents for various tasks [96]. But we currently lack tools to study these interactions beyond narrow applications, preventing identification of generalizable interaction principles and creating methodological fragmentation. Without cross-domain research infrastructure, each study must develop ad-hoc evaluation methods. As agents become multi-modal and multi-task capable, single-domain studies become less representative of real deployment scenarios. Such platforms should also enable the study of emergent interaction patterns that arise when humans and agents compete, collaborate, and interact across broad domains [3]. Multi-agent RL can contribute toward the goal of studying these interactions by providing a cross-domain means of simulating interaction patterns. If the human modeling component is sufficiently accurate, then these findings can generalize to real human-AI interactions while alleviating the human burden of interacting with an agent over a long period of time.

**Design considerations for feedback elicitation.** In RLHF, we typically make some assumption about user noise with respect to preferences. For example, we may assume a BTL model, as we mention in Section 2.1.4. However, relatively little work has investigated the method by which users give feedback, and how designing better interfaces and scaffolding users may lead to less noisy feedback overall [126]. For example, introducing a decoy option—one that is clearly inferior to a specific alternative but not to others—can bias users toward selecting that alternative [325], indicating that preferences are context-dependent and susceptible to noise introduced by presentation. Designing better interfaces can help elicit more accurate feedback to enable better integration with learning algorithms to train AI agents. Mechanism design may also offer tools to help elicit and aggregate potentially conflicting preferences.

## 8.3.2 HCAI for RL

We now discuss how the proliferation of AI agents in human environments creates opportunities for RL researchers to consider and integrate human-centered questions into the design of algorithms. We start by identifying key assumptions in the preference learning paradigm that break down when real users are involved. These assumptions include aggregating diverse and conflicting user preferences and incorporating multiple modalities of feedback. We then comment on how integrating human norms can help AI alignment generalize in the absence of specific task instructions.

**Alignment at scale.** RLHF pipelines typically aggregate feedback across users [221, 262]. But simply averaging feedback can obscure disagreement, reinforce majority biases, or even result in decisions that are suboptimal from the perspective of all users. As AI agents are increasingly deployed in multi-user environments, they must navigate conflicting preferences and trade-offs across stakeholders. As a result, we need RL algorithms that can model distributions over preferences, detect when conflicts arise, and reason about how to resolve them. Conflict resolution may involve integrating voting-based mechanisms or through facilitating iterative interactions and negotiations with stakeholders. Studying how alignment breaks down at scale can also inform when human oversight, delegation, or appeals should be built into the system, especially when no single action is universally acceptable.

**RLHF with diverse and underrepresented preferences.**    Another challenge is most RLHF pipelines are inflexible in how users can express feedback. Users may express preferences in hybrid forms: through corrections [270], pairwise comparisons [221, 291], or the state of the world [289]. But current methods for preference modeling assume a single fixed feedback modality, unnecessarily restricting important data sources and limiting the types of preferences that can be learned. If an agent only accepts one form of input, users may struggle to express certain preferences at all. Going forward, a key open question is: how can we design learning algorithms that enable agents to make effective use of diverse feedback modalities, particularly to model rare but important preferences in the long tail? But even with flexible and multi-modal feedback, agents must still generalize to novel states and contexts in which there is no explicit guidance on how to act.

**Norm-guided decision making.**    Norms enable people to generalize behavior in the absence of specific instructions. They encode shared expectations across people, roles, and contexts to formalize acceptable behavior. This shift from individual to collective expectations becomes critical as agents interact with diverse user populations—such as different cultures [13, 182]—or operate in multi-agent environments [279]. Although individual preferences may vary widely, norms serve as a common-sense behavioral prior to help agents act reasonably even when no ground-truth preference is available. However, it is not clear how to model and integrate norms into the decision-making process of RL agents [218]. For example, norms can span multiple levels of abstraction, from low-level action constraints to high-level interaction patterns. Additional key challenges ahead include developing scalable algorithms for norm induction, designing agents that can detect and explain violations in real time, and building mechanisms for adapting behavior to context-sensitive and potentially conflicting norms.

# 9 Conclusion

To close, in this dissertation we showcase a variety of new developments in AI that have been designed for human-centered challenges that emerge from the current or future deployment of AI agents. As we have shown, there are ample opportunities for future research. These challenges require novel methodological developments to connect and build upon the current work in HCI, RL, and application. To address these pressing challenges in such an unprecedented era, I believe we must rethink assumptions at scale to bring forward a future that leads to human flourishing.

# Part V

# APPENDICES

# Appendix: Building and Understanding Human-Like Agents

## 1 Hyperparameters

Table 1 details the hyperparameters used for training all three of the agent types with PPO.

## 2 Behavioral Study Details

In this section, we include additional details about our MTurk behavioral study.

**Data collected.** In addition to the consent, familiarity, and HNTT questions, we collected the following data: timing data broken down by page and the order in which the trials were presented to each participant.

**Full instructions.** We detail the full instructions included in the MTurk study here.

| Hyperparameter | Value |
|---|---:|
| Batch size | 2048 |
| Dropout rate | 0.1 |
| Learning rate | 2.5e-4 |
| Optimizer | Adam [164] |
| Gamma | 0.99 |
| $\lambda$ | 0.95 |
| Clip range | 0.2 |
| Gradient norm clipping coefficient | 0.5 |
| Entropy coefficient $c_2$ | 0.0 |
| Value function coefficient $c_1$ | 0.5 |
| Minibatches per update | 4 |
| Training epochs per update | 4 |
| Replay buffer size | 5 x batch size |

Table 1: **Hyperparameters for training the *symbolic*, *hybrid*, and *reward-shaping* agents.** We train all of the agents with the PPO algorithm [284]. For additional detail on what these hyperparameters correspond to, we encourage an interested reader to refer to the original PPO paper. We provide these hyperparameter values for reproducibility.

We are conducting a survey on navigation in video games for a research project. Please read the **Description** and **Requirements**, and then select the link below to complete the survey. At the end of the survey, you will receive a code to paste into the box below to receive credit for taking our survey.
**Description**:
- **Overview**: The survey is anonymous and includes a required consent form, comprehension check, some background info, and 6 video sections with 3 questions each. All questions are marked *required.
- **Time required**: about **30 minutes**.
- **Compensation**: you will receive a fixed compensation of **$6.50** for completing the task, with potential for a **$1 bonus** for a high-quality response. For example, copy/pasting answers, or responses that are not specific to the videos on each page, will not get the bonus.
- The MTurk HIT has a 1-hour duration. It will **not** allow you to submit after 1-hour has passed (*remember to submit or return HITs within 1-hour so you don't time out!*)
- If you start the task but change your mind, you may terminate your participation at any time and **return the HIT** within 1-hour, but you will **not** be paid for returned HITs or partial completions.

**Requirements**:
- You must complete all the questions.
- You must not have previously completed a HIT called "Navigation Turing Test (NTT)". Repeat participants are ineligible and will not be paid.
- You cannot participate from tablets or mobile phones.

**Make sure to leave this window open as you complete the survey.** When you are finished, you will return to this page to paste the code into the box.

# 3  Coding Details for Human-Like Characteristics

We include the coding guide agreed upon and used by both annotators when annotating their responses. Figure 1 shows a screenshot of the guide.

Goals of Codes

- Independence
- Separation

General Protocol

- Ideally, we want nice separation between codes. We also don't want the codes to be highly dependent. For example, we don't want every instance of goal directed to also be coded as collision avoidance. Therefore, if something can be coded in two ways, it is better to code it as the more *specific* instantiation.
- If we think collision avoidance, only code as collision avoidance and not environment receptivity.
- If we think smoothness, only code as smoothness and not mechanistic. This is true even if the response calls out stereotyped bot behavior.

Codes

1. **Smoothness of movement:** refers to the quality of the agent's navigation or camera movement. This code considers both more immediate jerky actions and more temporally-extended zig-zagging behavior.
   a. Key words:
      i. Smooth (+)
      ii. Janky (-)
      iii. Jerky (-)
      iv. Straight (+)
      v. Swerve (-)
      vi. Moving camera (-)
      vii. Steady camera (+)
      viii. Fluid (+)
      ix. Rigid (-)
      x. Zig-zag (-)
      xi. Clunky (-)
2. **Goal directed** looks at how intentional the agent's behavior seems to be. Participants describing behavior that pertains to a perceived goal, even if that goal is not the main one in the video, is included in this code. For example, if they refer to intentional exploration, we code it as goal directed.
   a. Key words:
      i. "Knew where they were going" (+)
      ii. Moving with intention (+)
      iii. Navigating with purpose (+)
      iv. Moving with focus (+)
      v. Making decisions (+)
3. **Collision avoidance** can be considered an instantiation of goal-directed behavior, but we set up a separate code because we believe it is an important feature in its own right. To preserve the independence of codes, behaviors that are coded as collision avoidance should not be coded as any other code --- except when e.g., the participant calls out another goal-directed behavior in the same response. An example of collision avoidance would be an agent trying not to "crash" into other objects. An example of a response that would not be considered collision avoidance is if an agent tries to collect a power-up.
   a. Key words:
      i. Collide (-)
      ii. Avoid (+)
      iii. Crash (-)
      iv. Run into obstacle (-)
4. **Environment receptivity** aims to capture the agent's relationship with the game environment. In a real-world setting, this might look like a person walking on a path instead of the grass or being responsive to the environment (such as a walk sign). Dynamic. The focus here is on the *intention* of the behavior. This code also aims to capture abiding by *norms* or (potentially unspoken) conventions that people may only have a sense for.
   a. Keywords:
      i. Explore (+)
      ii. Stay on the path (+)
      iii. Collecting power-ups (+)
      iv. Already knows everything about the environment (-)
      v. Take shortcut (+)
      vi. Seeing through walls (-)
      vii. Took shortcut (+)
5. **Non-mechanistic** is a more nebulous code that tries to capture pre-conceived notions of human imperfection. An example of mechanistic behavior is one identified as being too perfect. This code also captures any references to mistakes or error correction.
   a. Keywords:
      i. Too perfect (-)
      ii. Makes mistakes (+)
      iii. Precise (-)
      iv. Micro-corrections (-)
      v. Random (+)
      vi. Overcorrect (-)
6. **Intuition** refers to feelings that a behavior was more human-like, without sufficient specific justification. We include this code to capture instances where participants can identify what they believe is more human-like behavior but struggle to express it.
   a. Keywords:
      i. Natural
      ii. Feeling
      iii. Seems to be
      iv. Normal
      v. Realistic
7. **Self-reference** is a code to capture when judges relate the agent behavior to their own play. This could look like participants mentioning that they would feel ill if they played this way or mentioning that they typically collect all power-ups when playing.
   a. Keywords:
      i. "Like how I play"

Discard

1. If the participant seems to randomly assign human-like or not. This may look like saying, "I don't know" or "I can't tell". In contrast, if they say that the response is based off of feeling, we categorize it as intuition.
2. If the free-form response is in conflict with the actual response.

Figure 1: **Coding guide used by the annotators**. The guide includes a description of the codes, the general protocol, and the discard guidelines.

# Appendix: Reducing Human Annotation Burden for Concept-Based Policies

| Concept Name | Value Range | GPT-4o Error |
|---|:---:|:---:|
| Cart Position | $(-2.4, 2.4)$ | $0.01 \pm 0.00$ |
| Cart Velocity | $\mathbb{R}$ | $0.31 \pm 0.12$ |
| Pole Angle | $(-0.2095, 0.2095)$ | $0.01 \pm 0.00$ |
| Pole Angular Velocity | $\mathbb{R}$ | $0.63 \pm 0.14$ |

Table 1: **Concepts and GPT-4o labeling errors in PixelCartPole.** All concepts are continuous. Error is measured using MSELoss. The $\pm$ [value] part shows the standard deviation. GPT-4o struggles the most with the Pole Angular Velocity. This may be because angular velocity is a second-order property, which requires a comparison of changes in angle over time, and the stacked frames may not suffice for reliable estimation.

| Concept Name | Value Range | GPT-4o Error |
|---|:---:|:---:|
| Agent Position x | 3 | $0.03 \pm 0.06$ |
| Agent Position y | 3 | $0.04 \pm 0.06$ |
| Agent Direction | 4 | $0.20 \pm 0.05$ |
| Obstacle 1 Position x | 3 | $0.08 \pm 0.02$ |
| Obstacle 1 Position y | 3 | $0.19 \pm 0.04$ |
| Obstacle 2 Position x | 3 | $0.22 \pm 0.09$ |
| Obstacle 2 Position y | 3 | $0.12 \pm 0.02$ |
| Direction Movable Right | 2 | $0.19 \pm 0.06$ |
| Direction Movable Down | 2 | $0.14 \pm 0.05$ |
| Direction Movable Left | 2 | $0.13 \pm 0.02$ |
| Direction Movable Up | 2 | $0.07 \pm 0.08$ |

Table 2: **Concepts and GPT-4o labeling errors in DynamicObstacles.** All concepts are discrete. Error is 1 minus classification accuracy. The $\pm$ [value] part shows the standard deviation. GPT-4o is most successful at estimating the agent's x and y positions, likely because they are relatively easy to compute given clear instructions.

# 1 Concept Details

In this section, we provide more details about the concepts that we use in our experiments. We first specify the concepts more precisely, then provide a brief investigation of the concept errors produced by GPT-4o. The goal of this investigation is to summarize where state-of-the-art VLMs currently perform well and where they could improve. We then emphasize the importance of well-specified concepts for concept-based RL by visualizing some example start configurations and the associated concepts for one of the environments.

| Concept Name | Value Range | GPT-4o Error |
|---|---|---|
| Agent Position x at Frame 1 | 160 | $0.81 \pm 0.14$ |
| Agent Position y at Frame 1 | 210 | $0.93 \pm 0.05$ |
| Enemy Position x at Frame 1 | 160 | $0.78 \pm 0.13$ |
| Enemy Position y at Frame 1 | 210 | $0.90 \pm 0.07$ |
| Agent Position x at Frame 2 | 160 | $0.67 \pm 0.06$ |
| Agent Position y at Frame 2 | 210 | $0.88 \pm 0.04$ |
| Enemy Position x at Frame 2 | 160 | $0.71 \pm 0.05$ |
| Enemy Position y at Frame 2 | 210 | $0.83 \pm 0.08$ |

Table 3: **Concepts and GPT-4o labeling errors in Boxing.** All concepts are discrete. Error is 1 minus classification accuracy. Predictions are considered correct if within distance 5 of the ground truth. The $\pm$ [value] part shows the standard deviation.

| Concept Name | Value Range | GPT-4o Error |
|---|---|---|
| Ball Position x | 176 | $0.99 \pm 0.00$ |
| Ball Position y | 88 | $0.98 \pm 0.00$ |
| Ball Velocity x | 176 | $0.95 \pm 0.00$ |
| Ball Velocity y | 176 | $0.80 \pm 0.00$ |
| Paddle Velocity | 176 | $0.78 \pm 0.05$ |
| Paddle Acceleration | 352 | $0.79 \pm 0.04$ |
| Paddle Jerk | 704 | $0.83 \pm 0.04$ |

Table 4: **Concepts and GPT-4o labeling errors in Pong.** All concepts are discrete. Error is 1 minus classification accuracy.

**Different environments have various concepts.** In addition to Table 4.3 from Chapter 4, Tables 1 to 4 provide more details regarding the concepts used in each environment, categorizing them by their names, types, and value ranges. For the PixelCartPole environment, all concepts such as Cart Position, Cart Velocity, Pole Angle, and Pole Angular Velocity are continuous. In contrast, the DynamicObstacles environment lists discrete concepts, including Agent and Obstacle positions, with corresponding value ranges.

| Agent position: (2, 3) | Agent position: (2, 2) | Agent position: (0, 1) |
| Agent direction: down | Agent direction: up | Agent direction: up |
| Key position: (3, 1) | Key position: (2, 1) | Key position: (1, 5) |
| Door position: (4, 2) | Door position: (3, 4) | Door position: (2, 2) |
| Door open: no | Door open: no | Door open: no |
| Can move right: yes | Can move right: no | Can move right: no |
| Can move down: yes | Can move down: yes | Can move down: no |
| Can move left: yes | Can move left: yes | Can move left: no |
| Can move up: yes | Can move up: yes | Can move up: yes |

Figure 1: **Example start configurations of the DoorKey environment, along with the associated concepts for each configuration.** This figure shows that concepts must be general enough to apply to many different configurations. As a result, the agent cannot memorize the exact position of the door and key.

**Precise concept definitions are crucial for enabling correct agent behavior.** We visualize the start configurations for DoorKey in Figure 1 to illustrate the importance of well-defined concepts in concept-based interpretable RL. DoorKey includes a variety of initial states, each with different positions of the agent, key, and door. This diversity highlights a challenge in concept-based RL: how to define concepts that are both specific enough to be meaningful and general enough to apply across all possible scenarios. An agent trained with poorly defined concepts might perform well in some configurations but fail to generalize to others, leading to suboptimal performance in new or unseen environments. As a result, concept engineering is its own separate but important problem.

# 2  LICORICE Details

In this section, we provide additional implementation details for LICORICE. The model architecture has been mostly described in the main text and we state all additional details here. The number of neurons in the concept layer is exactly the number of concepts. For continuous concept values, we directly use a linear layer to map from features to concept values. For discrete concept values, since different concepts have different numbers of categories, we create one linear classification head for each single concept, and to predict the final action, we calculate the class with the largest predicted probability for each concept.

**Architecture.** The network begins with a CNN-based feature extractor $f_\theta : \mathcal{X} \to \mathbb{R}^d$, comprising three convolutional layers that map input state $x \in \mathcal{X}$ to a $d$-dimensional feature vector $h = f_\theta(x)$. We then introduce a linear layer $g_\phi : \mathbb{R}^d \to \mathbb{R}^k$ for concept prediction ($g$), mapping extracted features to $k$ concept values. For regression tasks, each of the $k$ predicted concept is represented as a real value; for classification tasks, each concept maps to one categorical value derived from a classification head. The action prediction component ($f$) is implemented as an MLP with two fully connected layers, each containing 64 neurons with Tanh activation, taking only $c$ as input to produce action logits $a = \mathrm{MLP}(c)$. Notably, we share the CNN feature extractor $f_\theta$ between policy and value functions, a decision informed by improved performance observed in preliminary experiments.

**Value-based methods.** If we were to use a value-based method as the RL backbone, we would need to make the following changes. First, we would need to modify $V(s, a)$ or $Q(s, a)$ to include a concept bottleneck, such that $Q(s, a) = f(g(s))$. Then, we can conduct interleaved training in a similar way to LICORICE.



Figure 2: Architecture of our concept-bottleneck actor-critic method.

**Feature extractor.** If we use the actor-critic paradigm, we propose to share a feature extractor between the policy and value networks, shown in Figure 2. Intuitively, this choice can offer several advantages compared with using image or predicted concepts as input for both networks. Sharing a feature extractor enables both networks to benefit from a common, rich representation of the input data, reducing the number of parameters to be trained. More importantly, it balances the updates of the policy and value networks. In experiments, we observe that directly using the raw image as input for both networks complicated policy learning. Conversely, relying solely on predicted concepts for the value network may limit its accuracy in value estimation, particularly if the concepts do not capture all the nuances relevant to the value predictions.

# 3 VLM Details

We detail our prompts for each environment here.

**PixelCartPole**

> **Prompt:** Here are the past 4 rendered frames from the CartPole environment. Please use these images to estimate the following values in the latest frame (the last one):
> - Cart Position, within the range (-2.4, 2.4)
> - Cart Velocity
> - Pole Angle, within the range (-0.2095, 0.2095)
> - Pole Angular Velocity
>
> Additionally, please note that the last action taken was [last action].
>
> Please carefully determine the following values and give concise answers one by one. Make sure to return an estimated value for each parameter, even if the task may look challenging.
>
> Follow the reporting format:
> - Cart Position: estimated_value
> - Cart Velocity: estimated_value
> - Pole Angle: estimated_value
> - Pole Angular Velocity: estimated_value

**DoorKey**

**Prompt:** Here is an image of a 4x4 grid composed of black cells, with each cell either empty or containing an object. Each cell is defined by an integer-valued coordinate system starting at (1, 1) for the top-left cell. The coordinates increase rightward along the x-axis and downward along the y-axis. Within this grid, there is a red isosceles triangle representing the agent, a yellow cell representing the door (which may visually disappear if the door is open), a yellow key icon representing the key (which may disappear), and one green square representing the goal. Carefully analyze the grid and report on the following attributes, focusing only on the black cells as the gray cells are excluded from the active black area.

Detailed Instructions:

1. Agent Position: Identify and report the coordinates (x, y) of the red triangle (agent). Ensure the accuracy by double-checking the agent's exact location within the grid.
2. Agent Direction: Specify the direction the red triangle is facing, which is the orientation of the vertex (pointy corner) of the isosceles triangle. Choose from 'right', 'down', 'left', or 'up'. Clarify that this direction is independent of movement options.
3. Key Position: Provide the coordinates (x, y) where the key is located. If the key is absent, report as (0, 0). Verify visually that the key is present or not before reporting.
4. Door Position:
   - Position: Determine and report the coordinates (x, y) of the door.
   - Status: Assess whether the door is open or closed (closed means the door is visible as a whole yellow cell, while open means the door disappears visually). Report as 'true' for open and 'false' for closed. Double-check the door's appearance to confirm if it is open or closed.
5. Direction Movable: Evaluate and report whether the agent can move one cell in each specified direction, namely, the neighboring cell in that direction is active and empty (not key, closed door, or grey inactive cell):
   - Right (x + 1): Check the cell to the right.
   - Down (y + 1): Check the cell below.
   - Left (x - 1): Check the cell to the left.
   - Up (y - 1): Check the cell above.
   Each direction's feasibility should be reported as 'true' if clear and within the grid, and 'false' otherwise.

Reporting Format: Carefully report each piece of information sequentially, following the format 'name: answer'. Ensure each response is precise and reflects careful verification of the grid details as viewed.

**DynamicObstacles**

**Prompt:** Here is an image of a 3x3 grid composed of black cells, with each cell either empty or containing an object. Each cell is defined by an integer-valued coordinate system starting at (1, 1) for the top-left cell. The coordinates increase rightward along the x-axis and downward along the y-axis. Within this grid, there is a red isosceles triangle representing the agent, two blue balls representing obstacles, and one green square representing the goal. Please carefully determine the following values and give concise answers one by one:

1. Agent Position: Identify and report the coordinates (x, y) of the red triangle (agent). Ensure the accuracy by double-checking the agent's exact location within the grid.
2. Agent Direction: Specify the direction the red triangle is facing, which is the orientation of the vertex (pointy corner) of the isosceles triangle. Choose from 'right', 'down', 'left', or 'up'. Clarify that this direction is independent of movement options.
3. Obstacle Position: Identify and report the coordinates of the two obstacles in ascending order. Compare the coordinates by their x-values first. If the x-values are equal, compare by their y-values.
    a) First Obstacle: Provide the coordinates (x, y) of the first blue ball.
    b) Second Obstacle: Provide the coordinates (x, y) of the second blue ball.
4. Direction Movable: Evaluate and report whether the agent can move one cell in each specified direction, namely, the neighboring cell in that direction is active and empty (not obstacle or out of bounds):
    • Right (x + 1): Check the cell to the right.
    • Down (y + 1): Check the cell below.
    • Left (x - 1): Check the cell to the left.
    • Up (y - 1): Check the cell above.
    Each direction's feasibility should be reported as 'true' if clear and within the grid, and 'false' otherwise.

Reporting Format: Carefully report each piece of information sequentially, following the format 'name: answer'. Ensure each response is precise and reflects careful verification of the grid details as viewed.

**Boxing**

> **Prompt:** Here are two consecutive rendered frames from the Atari Boxing environment. The game screen is 160x210 pixels, with (0, 0) at the top-left corner. The x-coordinate increases rightward, and the y-coordinate increases downward. For each frame, estimate the following values as integers:
> - The white player's x and y coordinates
> - The black player's x and y coordinates
>
> Please carefully determine the following values and give concise answers one by one. Make sure to return an estimated value for each one, even if the task may look challenging.
>
> Follow the reporting format:
> - frame 1 white x: estimated_value
> - frame 1 white y: estimated_value
> - frame 1 black x: estimated_value
> - frame 1 black y: estimated_value
> - frame 2 white x: estimated_value
> - frame 2 white y: estimated_value
> - frame 2 black x: estimated_value
> - frame 2 black y: estimated_value

**Pong**

**Prompt:** Here are four consecutive rendered frames from the Atari Pong environment. The game screen is 84x84 pixels, with (0, 0) at the top-left corner. You need to look at the ball and the paddle in the right. The x-coordinate increases downward, and the y-coordinate increases rightward.
Please determine the following values and give concise answers one by one. Make sure to return an estimated value for each one, even if the task may look challenging.
- x coordinate of the ball relative to the paddle in the first frame
- y coordinate of the ball in the first frame
- x velocity of the ball calculated from the first and second frames
- y velocity of the ball calculated from the first and second frames
- velocity of the paddle calculated from the first and second frames
- acceleration of the paddle calculated from the first three frames
- jerk of the paddle calculated from the four frames

Follow the reporting format in your response.
- x coordinate of the ball: estimated_value
- y coordinate of the ball: estimated_value
- x velocity of the ball: estimated_value
- y velocity of the ball: estimated_value
- velocity of the paddle: estimated_value
- acceleration of the paddle: estimated_value
- jerk of the paddle: estimated_value

# 4 Experimental Details

We now provide important hyperparameters and details regarding our experimental results to support reproducibility.

## 4.1 Hyperparameters

**Behavior learning hyperparameters.** For the PPO hyperparameters, we set $4 \cdot 10^6$ total timesteps for PixelCartPole and DoorKey, $10^6$ for DynamicObstacles, $1.5 \cdot 10^7$ for Boxing, and $10^7$ for Pong. For PixelCartPole, DoorKey, and DynamicObstacles, we use 8 vectorized environments, horizon $T = 4096$, 10 epochs for training, batch size of 512, learning rate $3 \cdot 10^{-4}$, entropy coefficient 0.01, and value function coefficient 0.5. For Boxing and Pong, we use 8 vectorized environments, horizon $T = 1024$, 4 epochs for training, batch size of 256, learning rate $3 \cdot 10^{-4}$, entropy coefficient 0.01, and value function coefficient 0.5. For all other hyperparameters, we use the default values from Stable Baselines 3 [263].

**Concept learning hyperparameters.** For the concept training, we set 100 epochs with Adam optimizer with the learning rate linearly decaying from $3 \cdot 10^{-4}$ to 0 for each iteration in PixelCartPole, Boxing, and Pong. In DoorKey and DynamicObstacles, we use the same optimizer and initial learning rate, yet set 50 epochs instead and set early stopping with threshold linearly increasing from 10 to 20, to incentivize the concept network not to overfit in earlier iterations. The batch size is 32. We model concept learning for PixelCartPole as a regression problem (minimizing mean squared error). We model concept learning for DoorKey, DynamicObstacles, Boxing, and Pong as classification problems.

**LICORICE-specific hyperparameters.** For LICORICE, we set the ratio for active learning $\tau = 10$, batch size to query labels in the active learning module $b = 20$, and the number of ensemble models $N = 5$ for the first three environments. For the complex environments Boxing and Pong, we choose the ratio for active learning $\tau = 4$, batch size to query labels in the active learning module $b = B_m/5$ and number of ensemble models $N = 5$ to make a balance between performance and speed. The default number of iterations chosen in our algorithm is $M = 4$ for PixelCartPole, $M = 2$ for DoorKey, DynamicObstacles, and $M = 5$ for Boxing and Pong. The sample acceptance rate $p = 0.02$ for PixelCartPoleand $p = 0.1$ for the other four environments. The Random-Q baseline uses the same sample acceptance rate $p$ as LICORICE. In the complex environments Boxing and Pong, we use both the KL-divergence penalty and PPO loss at the end of the algorithm to improve optimization with $\beta = 0.01$, as mentioned in Section 4.3.

## 4.2 Other Experimental Details

**Computational resources.**  We use NVIDIA A6000 and NVIDIA RTX 6000 Ada Generation. Each of our training program uses less than 3GB GPU memory. For Boxing and Pong, each run takes less than 20 hours to finish. For PixelCartPole and DoorKey, each run takes less than 9 hours to finish. For DynamicObstacles, each run takes less than 2 hours to finish.

**Experimental replicability.**  For all experiments, we use 5 seeds [123, 456, 789, 1011, 1213] to train the models. We then evaluate on the environment with seed 42 with 100 episodes and take the average of the reward across the episodes and the concept error across the observations within each episode.

# 5  Additional Results

In this section, we present both the numerical results as tables from the figures in the main paper, in addition to experimental results not reported in the main paper. As a reminder, the reward is reported as the fraction of the reward upper bound. For PixelCartPole, the $c$ error is the MSE. For the other environments, the $c$ error is 1 - accuracy. Sequential-Q, Disagreement-Q, Random-Q, and the LICORICE variants are all budget-constrained, whereas CPM is given an unlimited budget.

## 5.1 Balancing Concept Performance and Environment Reward

For completeness, we present all of the numerical results from Figure 4.2, Chapter 7 in Table 5. Here, the budget-constrained algorithms are given budgets of $B = [500, 300, 300, 3000, 5000]$ for each environment, from top to bottom; CPM is given an unlimited budget (in practice, it uses 4M, 4M, 1M, 15M, 10M concept labels respectively). Our method enjoys low variance across all environments in terms of both concept error and reward.

## 5.2 Budget Allocation Effectiveness

In Table 6, we present an extension of the results in Figure 4.3, including the standard deviation. As expected, as the budget increases, the standard deviation for both the reward and concept error tends to decrease. The one exception is the reward for PixelCartPole. Interestingly, the standard deviation is highest for $B = 400$. We suspect this is because the concept errors may be more critical here, leading to higher variance in the reward performance.

Figure 3: **LICORICE+GPT-4o waits until the coast is clear to move to the goal.** It appears to make a small mistake in the bottom left, requiring it to wait slightly longer than necessary to navigate to the goal.

## 5.3 Integration with Vision-Language Models

In Table 7, we present an extension of the results in Figure 4.3 in the main paper, including the standard deviation. Interestingly, the standard deviation for the reward obtained by using LICORICE with GPT-4o as the annotator does not always follow the same trend as shown in Table 6 (when we assume access to a more accurate human annotator). Instead, the standard deviation is relatively consistent for PixelCartPole, regardless of the budget. It steadily decreases for DoorKey, as expected. However, in DynamicObstacles, we see an increase when $B = 300$. We are not sure of the cause of this. Perhaps at this point the algorithm begins overfitting to the errors in the labels from GPT-4o (the concept error rate is the same for 200 and 300 labels). Further investigation is required to understand the underlying factors contributing to this anomaly.

**Different tasks and concepts have various difficulties for GPT-4o.** In addition to Table 4.3 from Chapter 4, Tables 1 to 4 list detailed concept errors for concepts in all environments. In PixelCartPole, cart position and pole angle have smaller errors, while velocities require understanding multiple frames and thus are harder to predict accurately. For DoorKey and DynamicObstacles, different concepts have slightly varying concept errors, indicating visual tasks have different difficulties for GPT-4o. Agent direction has as high as around 0.3 prediction error for both DoorKey and DynamicObstacles. The concept accuracies in DoorKey are generally higher than DynamicObstacles, suggesting GPT-4o struggles more with a larger grid.

Figure 4: **Learning curves for ablations in three environments: PixelCartPole, DoorKey, and DynamicObstacles.** Shaded region shows 95% CI, calculated using 1000 bootstrap samples.

**A GPT-4o-trained RL agent recovers from a mistake.** In Figure 3, we show an example of a GPT-4o-trained RL agent on DynamicObstacles, in which the agent appears to wait until it is safe to move towards the goal: the green square. The image sequence shows the agent (red triangle) starting from its initial position and moving to the right. It then is cornered by an obstacle (blue circle), then both obstacles. In the bottom left corner frame, it appears to make a mistake by turning to the right, meaning it missed a window to escape. Finally, it moves to the goal when the path is clear in the second-to-last and last frames (bottom right). This behavior highlights that the agent may still learn reasonable behavior even if the concept labels may be incorrect (causing it to make a mistake, as in the bottom left frame).

## 5.4 Ablation Study

In Figure 4 shows all learning curves for ablations in all environments, extending the results from Table 4.4 in the main dissertation. In PixelCartPole, we clearly see the benefit of iterative strategies on reward: LICORICE, LICORICE-DE and LICORICE-AC consistently increase in reward and converge at high levels, but LICORICE-IT converges at less than $50\%$ of the optimal reward. We also see that LICORICE-IT also achieves higher concept error, indicating that it struggles to learn the larger distribution of concepts induced by a non-optimal policy. In DoorKey, all algorithms steadily increase in reward. However, we can see a dip at around $10^6$ environment steps where we begin the second iteration for LICORICE, LICORICE-AC, and LICORICE-DE. Although LICORICE-IT achieves similar reward to LICORICE, LICORICE achieves lower concept error, which is beneficial for the goal of interpretability. Finally, in DynamicObstacles, LICORICE-AC lags behind the most in terms of both reward and concept error. This result indicates that the active learning component is most important for this environment.

157

| | Algorithm | R ↑ | c Error ↓ |
|---|---|---|---|
| | Sequential-Q | $0.23 \pm 0.09$ | $0.10 \pm 0.04$ |
| | Disagreement-Q | $0.36 \pm 0.20$ | $0.10 \pm 0.04$ |
| | Random-Q | $0.34 \pm 0.14$ | $0.07 \pm 0.02$ |
| PixelCartPole | LICORICE | $1.00 \pm 0.00$ | $0.02 \pm 0.00$ |
| | LICORICE+GPT-4o | $0.06 \pm 0.01$ | $0.25 \pm 0.08$ |
| | CPM | $1.00 \pm 0.00$ | $0.01 \pm 0.00$ |
| | Sequential-Q | $0.46 \pm 0.09$ | $0.46 \pm 0.05$ |
| | Disagreement-Q | $0.76 \pm 0.12$ | $0.39 \pm 0.07$ |
| | Random-Q | $0.89 \pm 0.03$ | $0.27 \pm 0.08$ |
| DoorKey | LICORICE | $0.99 \pm 0.01$ | $0.05 \pm 0.01$ |
| | LICORICE+GPT-4o | $0.83 \pm 0.06$ | $0.31 \pm 0.01$ |
| | CPM | $1.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | Sequential-Q | $0.79 \pm 0.44$ | $0.01 \pm 0.01$ |
| | Disagreement-Q | $0.99 \pm 0.01$ | $0.00 \pm 0.00$ |
| | Random-Q | $0.98 \pm 0.01$ | $0.02 \pm 0.01$ |
| DynamicObstacles | LICORICE | $0.99 \pm 0.00$ | $0.00 \pm 0.00$ |
| | LICORICE+GPT-4o | $0.88 \pm 0.09$ | $0.13 \pm 0.03$ |
| | CPM | $0.97 \pm 0.02$ | $0.00 \pm 0.00$ |
| | Sequential-Q | $0.87 \pm 0.20$ | $0.32 \pm 0.26$ |
| | Disagreement-Q | $0.83 \pm 0.14$ | $0.20 \pm 0.21$ |
| | Random-Q | $0.96 \pm 0.04$ | $0.14 \pm 0.10$ |
| Boxing | LICORICE | $1.00 \pm 0.02$ | $0.05 \pm 0.01$ |
| | LICORICE+GPT-4o | $0.22 \pm 0.01$ | $0.98 \pm 0.01$ |
| | CPM | $1.00 \pm 0.05$ | $0.00 \pm 0.00$ |
| | Sequential-Q | $0.98 \pm 0.02$ | $0.54 \pm 0.03$ |
| | Disagreement-Q | $0.94 \pm 0.09$ | $0.57 \pm 0.05$ |
| | Random-Q | $0.94 \pm 0.05$ | $0.60 \pm 0.04$ |
| Pong | LICORICE | $1.00 \pm 0.02$ | $0.24 \pm 0.16$ |
| | LICORICE+GPT-4o | $0.13 \pm 0.17$ | $0.71 \pm 0.01$ |
| | CPM | $1.00 \pm 0.00$ | $0.00 \pm 0.01$ |

Table 5: **Evaluation of the reward $R$ and concept error achieved by all methods in all environments.** This is a extended table from Figure 4.2. The $\pm$ [value] part shows the standard deviation.

| Environment | B | M | R $\uparrow$ | c Error $\downarrow$ |
|---|---|---|---|---|
| PixelCartPole | 300 | 1 | $0.28 \pm 0.10$ | $0.15 \pm 0.07$ |
| | 400 | 3 | $0.77 \pm 0.22$ | $0.05 \pm 0.01$ |
| | 500 | 4 | $1.00 \pm 0.00$ | $0.02 \pm 0.00$ |
| DoorKey | 100 | 1 | $0.77 \pm 0.04$ | $0.26 \pm 0.04$ |
| | 200 | 2 | $0.93 \pm 0.02$ | $0.09 \pm 0.01$ |
| | 300 | 2 | $0.99 \pm 0.01$ | $0.05 \pm 0.01$ |
| DynamicObstacles | 100 | 1 | $0.96 \pm 0.02$ | $0.04 \pm 0.01$ |
| | 200 | 1 | $0.98 \pm 0.01$ | $0.01 \pm 0.00$ |
| | 300 | 2 | $0.99 \pm 0.00$ | $0.00 \pm 0.00$ |
| Boxing | 1000 | 5 | $0.94 \pm 0.05$ | $0.11 \pm 0.01$ |
| | 2000 | 5 | $0.98 \pm 0.05$ | $0.10 \pm 0.06$ |
| | 3000 | 5 | $1.00 \pm 0.02$ | $0.05 \pm 0.01$ |
| Pong | 1000 | 5 | $0.99 \pm 0.00$ | $0.65 \pm 0.04$ |
| | 3000 | 5 | $1.00 \pm 0.00$ | $0.43 \pm 0.05$ |
| | 5000 | 5 | $1.00 \pm 0.02$ | $0.24 \pm 0.16$ |

Table 6: **Performance of LICORICE on all environments for varying budgets.** We select $M$ to optimize $R$ - $c$ error without additional weighting since they are observed to have the same magnitude. The reward is reported as the fraction of the reward upper bound. The $\pm$ [value] part shows the standard deviation. This shows a more complete version of the results in Figure 4.3.

| Environment | B | M | LICORICE+GPT-4o | | GPT-4o |
| | | | $R\uparrow$ | $c$ Error $\downarrow$ | $c$ Error |
| --- | --- | --- | --- | --- | --- |
| PixelCartPole | 300 | 1 | $0.06 \pm 0.02$ | $0.17 \pm 0.19$ | $0.19 \pm 0.25$ |
| | 400 | 2 | $0.07 \pm 0.02$ | $0.25 \pm 0.15$ | $0.22 \pm 0.15$ |
| | 500 | 2 | $0.06 \pm 0.01$ | $0.25 \pm 0.08$ | $0.24 \pm 0.07$ |
| DoorKey | 100 | 1 | $0.71 \pm 0.04$ | $0.45 \pm 0.04$ | $0.31 \pm 0.03$ |
| | 200 | 2 | $0.74 \pm 0.04$ | $0.33 \pm 0.01$ | $0.31 \pm 0.03$ |
| | 300 | 2 | $0.83 \pm 0.06$ | $0.31 \pm 0.01$ | $0.30 \pm 0.03$ |
| DynamicObstacles | 100 | 1 | $0.27 \pm 0.37$ | $0.20 \pm 0.03$ | $0.16 \pm 0.03$ |
| | 200 | 1 | $0.93 \pm 0.05$ | $0.13 \pm 0.03$ | $0.12 \pm 0.05$ |
| | 300 | 1 | $0.88 \pm 0.09$ | $0.13 \pm 0.03$ | $0.13 \pm 0.03$ |
| Boxing | 1000 | 5 | $0.34 \pm 0.09$ | $0.98 \pm 0.01$ | $0.78 \pm 0.01$ |
| | 2000 | 5 | $0.23 \pm 0.11$ | $0.98 \pm 0.01$ | $0.79 \pm 0.01$ |
| | 3000 | 5 | $0.22 \pm 0.01$ | $0.98 \pm 0.01$ | $0.79 \pm 0.02$ |
| Pong | 1000 | 5 | $0.00 \pm 0.00$ | $0.79 \pm 0.08$ | $0.88 \pm 0.01$ |
| | 3000 | 5 | $0.00 \pm 0.00$ | $0.74 \pm 0.07$ | $0.86 \pm 0.03$ |
| | 5000 | 5 | $0.13 \pm 0.17$ | $0.71 \pm 0.01$ | $0.87 \pm 0.02$ |

Table 7: **Performance of LICORICE with GPT-4o integrated into the loop for all environments across different budgets, along with the concept labeling error of GPT-4o as a reference.** We select $M$ to optimize $R$ - $c$ error without additional weighting since they are observed to have the same magnitude. GPT-4o $c$ error is evaluated on a random sample of 50 observations from the same rollout set used for LICORICE+GPT-4o, due to API cost constraints. This shows a more complete version of the results in Figure 4.3.

| Environment | $B$ | Sequential-Q | | Disagreement-Q | | Random-Q | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | $R\uparrow$ | $c$ Error $\downarrow$ | $R\uparrow$ | $c$ Error $\downarrow$ | $R\uparrow$ | $c$ Error $\downarrow$ |
| PixelCartPole | 300 | $0.17 \pm 0.03$ | $0.11 \pm 0.03$ | $0.22 \pm 0.04$ | $0.20 \pm 0.04$ | $0.24 \pm 0.07$ | $0.14 \pm 0.04$ |
| | 400 | $0.23 \pm 0.05$ | $0.10 \pm 0.03$ | $0.26 \pm 0.06$ | $0.13 \pm 0.03$ | $0.28 \pm 0.07$ | $0.09 \pm 0.05$ |
| | 500 | $0.23 \pm 0.07$ | $0.10 \pm 0.03$ | $0.39 \pm 0.28$ | $0.09 \pm 0.03$ | $0.35 \pm 0.16$ | $0.07 \pm 0.02$ |
| DoorKey | 100 | $0.33 \pm 0.07$ | $0.51 \pm 0.08$ | $0.61 \pm 0.14$ | $0.33 \pm 0.07$ | $0.70 \pm 0.07$ | $0.30 \pm 0.02$ |
| | 200 | $0.43 \pm 0.09$ | $0.49 \pm 0.03$ | $0.72 \pm 0.09$ | $0.28 \pm 0.07$ | $0.85 \pm 0.04$ | $0.24 \pm 0.04$ |
| | 300 | $0.51 \pm 0.11$ | $0.45 \pm 0.08$ | $0.76 \pm 0.12$ | $0.41 \pm 0.06$ | $0.90 \pm 0.02$ | $0.27 \pm 0.08$ |
| DynamicObstacles | 100 | $0.43 \pm 0.47$ | $0.04 \pm 0.04$ | $0.82 \pm 0.18$ | $0.04 \pm 0.00$ | $0.34 \pm 0.46$ | $0.05 \pm 0.05$ |
| | 200 | $0.76 \pm 0.43$ | $0.03 \pm 0.03$ | $0.99 \pm 0.02$ | $0.01 \pm 0.00$ | $0.76 \pm 0.42$ | $0.03 \pm 0.02$ |
| | 300 | $0.78 \pm 0.44$ | $0.01 \pm 0.01$ | $0.99 \pm 0.01$ | $0.00 \pm 0.00$ | $0.97 \pm 0.01$ | $0.02 \pm 0.02$ |
| Boxing | 1000 | $0.84 \pm 0.20$ | $0.35 \pm 0.12$ | $0.93 \pm 0.05$ | $0.25 \pm 0.11$ | $0.83 \pm 0.24$ | $0.25 \pm 0.05$ |
| | 2000 | $0.88 \pm 0.12$ | $0.19 \pm 0.11$ | $0.93 \pm 0.06$ | $0.14 \pm 0.04$ | $0.83 \pm 0.11$ | $0.30 \pm 0.18$ |
| | 3000 | $0.87 \pm 0.20$ | $0.32 \pm 0.26$ | $0.83 \pm 0.14$ | $0.20 \pm 0.21$ | $0.96 \pm 0.04$ | $0.14 \pm 0.10$ |
| Pong | 1000 | $0.91 \pm 0.12$ | $0.74 \pm 0.07$ | $0.84 \pm 0.21$ | $0.68 \pm 0.06$ | $0.74 \pm 0.21$ | $0.73 \pm 0.05$ |
| | 3000 | $0.98 \pm 0.02$ | $0.58 \pm 0.08$ | $0.98 \pm 0.01$ | $0.64 \pm 0.03$ | $0.96 \pm 0.04$ | $0.60 \pm 0.04$ |
| | 5000 | $0.98 \pm 0.02$ | $0.54 \pm 0.03$ | $0.94 \pm 0.09$ | $0.57 \pm 0.05$ | $0.94 \pm 0.05$ | $0.60 \pm 0.04$ |

Table 8: **Performance of Sequential-Q, Disagreement-Q, and Random-Q for all environments across different budgets.** The $\pm$ [value] part shows the standard deviation. This shows a more complete version of the results in Figure 4.3.

# Appendix: Extracting Interpretable Policies for Multi-Agent Coordination

# 1 Experimental Details

## 1.1 Environments

For all environments, we utilize the initialization and reward scheme as described in the original MAD-DPG paper [193] and Pytorch implementation. The only change we make is to the predator-prey environment, which we describe below.

**Predator-prey.** We follow the definition of the original environment proposed in the multi-agent particle environment [193], with only changes in the partial observation provided to each agent. The observations of the adversary and the agents consist of the concatenation of the following vectors:

1. $[self\_pos, self\_vel]$

2. binarized relative positions and relative velocity (if applicable) of the landmarks and other agents using $sgn(x)$ as the binarizing function

3. binarized relative distance between all pairs of agents on the other team. If the opponent team has agent $a_1, a_2$, then it will be $[sgn(x_1 - x_2), sgn(y1 - y2)]$.

4. binarized relative distance between all pairs of agents on the same team.

For an environment with $K = M = 2$, the observation size will be 22 and 24 respectively for the adversary and the agents.

## 1.2 Hyperparameters

We vary the hyperparameters that would impact training performance of these algorithms 2-3 times and choose the hyperparameters that yield the agents with the best performance. For all environments, we vary the number of rollouts to be $[50, 100]$ and the number of iterations to be $[30, 100]$, while the threshold is fixed at $N - 1$. For the baselines, we also vary the maximum number of samples used for training each agent between $[10000, 30000, 100000]$. For Imitation DT, we did not see much of a performance increase between 3000 and 10000 samples, so we pick the maximum value for fairness of comparison. For Fitted Q-Iteration, we also did not see much performance increase after 30000 samples, so we chose 30000 samples due to time constraints. Table 1 shows the values of the hyperparameters that are utilized by all algorithms. Although we set a maximum number of training iterations for MAVIPER, we stop training early when there is no noticeable performance gain to further improve runtime.

| Algorithm | Environment | Max Training Iterations | Number of Rollouts | Threshold | Max Samples |
|-----------|-------------|-------------------------|--------------------|-----------|--------------|
| IVIPER | Physical Deception | 50 | 50 | | |
| | Cooperative Navigation | 100 | 50 | N/A | 300,000 |
| | Predator-prey | 100 | 100 | | |
| MAVIPER | All | 100 | 50 | $N-1$ | 300,000 |
| Imitation DT | All | 20 | N/A | N/A | 100,000 |
| Fitted Q-Iteration | All | 10 | N/A | N/A | 30,000 |

Table 1: **Hyperparameter values used for all algorithms.**

## 1.3 Algorithm Implementation Details

To optimize running speed, MAVIPER adopts a caching mechanism to avoid training a new decision tree for each data point being predicted. It also does parallelization for the `Predict` function by precomputing all the prediction information upfront, where the each prediction is delayed until all data points are looped over. In this way, MAVIPER can gather all the predictions that a particular tree needs to make and therefore do it in a parallel manner.

Since IVIPER is fully decentralized, training of each decision-tree can be performed in parallel.

For Fitted Q-Iteration, we bin the states into 10 (mostly) evenly-spaced bins:

$$(-\infty, -1.), [-1., -.75), [-.75, -.5), [-.5, -.25), [-.25, 0.),$$
$$[0., .25), [.25, .5), [.5, .75), [.75, 1.), (1., \infty).$$

We note that Fitted Q-Iteration may perform better with a better choice of bins; however, choosing the correct bin values requires either domain knowledge or extensive manual tuning to find the right balance between granularity, number of timesteps for training, and performance.

## 2 Additional Results

In Appendices 2.1 and 2.2, we further present results using the defined environment reward. This reward is not as intuitive as the primary metric for many of these environments, but we present the results here for the sake of completeness. The individual and joint performance ratios are defined in the same way as in Section 5.7 in the main body of the paper, with one caveat. Since we are now measuring reward, which may be negative or positive, we take $\frac{||A|-|B||}{A}$ to report how much more or less $B$ is than $A$.

(a) Physical Deception



(b) Cooperative Navigation



(c) Predator-prey

Figure 1: **Individual performance ratio measured by reward.** Individual performance ratio of agents that use decision-treepolicies compared to expert agents for different maximum depths. Higher is better. Error bars represent the 95% confidence interval.

## 2.1 Individual Performance

Figure 1 shows the individual performance ratio measured by reward on all three environments.

**Physical deception.**  Results for physical deception are shown in Fig. 1a. Interestingly, MAVIPER and IVIPER defenders only significantly outperform both baselines when the maximum depth is 2. When the maximum depth is 4, MAVIPER, IVIPER, and Imitation DT perform similarly, with Fitted Q-Iteration barely reaching above .50 for all depths. When the maximum depth is 6, Imitation DT actually achieves the highest defender reward, significantly outperforming all other algorithms. However, as shown in Figure 5.3 in the main body, MAVIPER significantly outperforms all algorithms for all maximum depths when reporting the success ratio. This is because the reward is in part dependent on the performance of the adversary. In other words, a poorly-performing defender can achieve similar performance to a high-performing defender if the poor-performing defender is paired with a high-performing adversary. We see a similar pattern for the adversary performance: IVIPER, MAVIPER, and Fitted Q-Iteration all significantly outperform the Fitted Q-Iteration baseline for different maximum depths. We note that MAVIPER tends to perform better for lower maximum depths, which is desirable for interpretability.

**Cooperative navigation.**  Results for cooperative navigation are shown in Figure 1b. IVIPER and MAVIPER significantly outperform the Fitted Q-Iteration baseline for all maximum depths. However, MAVIPER only significantly outperforms Imitation DT when the maximum depth is 2. Otherwise, IVIPER, MAVIPER, and Imitation DT all perform similarly.

**Predator-prey**  Results for predator-prey are shown in Fig. 1c. MAVIPER prey significantly outperform all other algorithms for maximum depths of 2 and 6. For a maximum depth of 4, MAVIPER and IVIPER algorithms both significantly outperform the baselines. In contrast, MAVIPER predators only significantly outperform all other algorithms for a maximum depth of 4. For a maximum depth of 2, MAVIPER and IVIPER significantly outperform the baselines, For a maximum of depth of 6, MAVIPER, IVIPER, and Imitation DT significantly outperform Fitted Q-Iteration. Again, we note that this performance is not necessarily reflected in the results using the collision metric in Figure 5.3c.

(a) Physical Deception

(b) Cooperative Navigation

(c) Predator-prey

Figure 2: **Joint performance ratio, measured by reward, of agents with decision-tree policies compared to expert agents for different maximum depths.** Here, we evaluate the agents jointly to assess the coordination ability of the algorithms that train decision-tree policies. Error bars represent the 95% confidence interval. Higher is better.

## 2.2 Joint Performance

Fig. 2 shows the joint performance ratio measured by reward in all three environments.

**Physical deception.** Figure 2a shows the results on physical deception. MAVIPER defenders significantly outperform all other algorithms on this environment for maximum depths of 2 and 4. For a maximum depth of 6, IVIPER, MAVIPER, and Imitation DTall perform similarly. Note that MAVIPER again achieves good performance for lower maximum depths, demonstrating its promise as an algorithm for producing interpretable policies. We also note that, again, the reward metric is somewhat deceptive: when measuring the success conditions in the environment (as in Figure 5.4 in the main body of the paper), MAVIPER significantly outperforms all other algorithms for all maximum depths.

**Cooperative navigation.** Figure 2b depicts the joint agent performance on the cooperative navigation environment. Interestingly, we see that MAVIPER significantly outperforms all other algorithms for all maximum depths, despite obtaining similar individual performance. Consequently, this means that MAVIPER better captures the desired coordinated behavior than all of the other algorithms.

**Predator-prey.**    Figure 2c shows the results for the predator-prey environment. MAVIPER prey significantly outperform other algorithms for a maximum depth of 2. For maximum depths of 6 and 8, it achieves slightly better (but not statistically significant) performance than IVIPER, and significantly outperforms the two baselines. MAVIPER and IVIPER predators enjoy similar performance for maximum depths of 2 and 6. For a maximum depth of 4, MAVIPER significantly outperforms all other algorithms. Note that the correct behavior in this environment is challenging to capture with a small decision tree, as the number of features is either 22 or 24, depending on the agent type.

# Appendix: Aligning Reinforcement Learning Policies with Human Feedback

# 1 Additional Evolutionary Algorithm Details

Here, we specify the details of the evolutionary algorithm (EA) that are excluded from the main paper.

**Tournament selection mechanism.** For selecting individuals during reproduction, we use tournament selection with a size of $k = 3$. Specifically, $k$ random individuals from the population compete based on their fitness, and the best among them is chosen as a parent. This process ensures that better-performing individuals are more likely to be selected while maintaining diversity by allowing weaker individuals a chance to reproduce.

**Individual encoding.** Instead of directly operating on the tree policies, we represent each policy as a real-valued feature vector in continuous space. Each policy is vectorized in depth-first order, with all vectors having a consistent length determined by a pre-specified maximum depth $D$. Specifically, each feature vector has a dimension of $\mathbb{R}^{2^{D+1}-1}$, encoding the structure and information of the policy. This allows for more efficient manipulation and comparison within the continuous feature space, avoiding the need to store and operate on the full decision trees directly.

**Crossover.** We choose to preserve individuals that are scored highly according to $\hat{\boldsymbol{\theta}}$, then perform subtree swapping as a crossover operator. Subtree swapping first randomly selects a node within the tree structure of the two parents in the crossover operation. Then, the subtrees rooted at this node are exchanged, resulting in two new individuals.

**Mutation.** To more broadly explore the space of possible trees, we perform mutations on randomly selected nodes of a tree in one of three ways, determined by tunable weights on the three mutation strategies: random, subtree rebuilding, and subtree removal. Random mutation ($\mathcal{M}_{\text{rand}}$) selects a random feature and a corresponding split value (or selects a new random action) to create a new node $v'$. Subtree rebuilding ($\mathcal{M}_{\text{subtree}}$) first retrieves the subset of data $\mathcal{D}_{\text{path}} \in \mathcal{D}$ corresponding to a chosen node. The subtree is then reconstructed using imitation learning with $\mathcal{D}_{\text{path}}$. Subtree removal replaces the subtree rooted at the selected node with an action.

**Fitness function calculation speedup.** Because we envision that this algorithm could actually interact with end users, we prioritize efficiency. One challenge is that, in each EA loop, we need to estimate the quality, or fitness, of each new decision-tree policy using $\hat{\boldsymbol{\theta}}$. When the feature vector comprises of simple policy attributes, this estimation is fast. However, when environment return is included, we often need to perform rollouts for a fixed number of episodes $\rho$ to estimate the expected return. This process could be computationally expensive, especially with a large number of candidates. To address this challenge, we propose an adaptive policy evaluation algorithm that dynamically adjusts its evaluation process based on the policy's performance profile, allowing for fewer episodes to be evaluated for less promising candidates. We use a UCB-inspired bound [301] to estimate the expected return for each tree as

$$\tilde{R} = \hat{R} + c\sqrt{\frac{\ln(n)}{n}}, \tag{1}$$

where $n$ is the number of episodes performed. As we gather more samples (increasing $n$), our confidence interval narrows, allowing us to be more certain estimate on the expected return of each arm. To allow for fewer episodes to be rolled out for less promising arms, we design a stopping condition

$$c\sqrt{\frac{\ln(n)}{n}} \leq \frac{\epsilon_{\text{base}}}{1 + \epsilon_{\text{scale}} \cdot \tilde{R}}, \tag{2}$$

which is updated after each rollout. This formulation embodies a key insight: as the estimate $\tilde{R}$ increases (suggesting a potentially better policy), we decrease $\epsilon$, demanding higher precision. This allows us to differentiate more accurately between high-performing policies.

| Hyperparameter | PotholeWorld-v0 | CartPole-v1 |
| --- | --- | --- |
| Population Size | 75 | 75 |
| Mutation Rate $\mu$ | 0.9 | 0.9 |
| Crossover Rate $\rho$ | 0.8 | 0.8 |
| Tournament Size $k$ | 3 | 3 |
| Max Generations | 125 | 100 |

Table 1: **Hyperparameters of the evolutionary algorithm for generating the upper and lower bounds.**

## 2 Experimental Details

In this section, we first describe the initialization process used to generate the set of initial policies. We then provide detailed descriptions of the hyperparameter settings used in our experiments.

### 2.1 Initialization Details

In real-world applications, policies typically must not only satisfy user preferences but also perform competently in their intended domain. We leverage this insight to develop an initialization method for candidate policies that are both performant and diverse with the goal of providing a strong starting point for preference elicitation. Our approach builds on VIPER [17], a DAgger-based interactive imitation learning algorithm [270] in which an RL policy $\pi^*$ acts as an expert to guide the training of decision-tree policies. Following this framework, we train a high-quality RL policy $\pi^*$ for the domain of interest. Then, we generate a pool of decision trees, which serve as candidate tree policies.

**Use randomization to promote diversity.** To ensure that there is sufficient diversity in the initial set of DT policies, we make two key modifications to VIPER. First, instead of only outputting the final best tree $\hat{\pi}$ according to return or alignment with the expert $\pi^*$, we instead use *all* tree policies generated during the $V$ VIPER iterations. Second, we introduce *depth randomization*, in which we choose a set of maximum depths to constrain the resulting depths of the trees. From this process, we obtain an initial population of $|\hat{\Pi}_0| = V \times D$ tree policies, which we then convert into their corresponding $d$-dimensional feature vectors $\mathbf{f}_{\hat{\pi}} \in \mathbb{R}^d \leftarrow \phi(\hat{\pi}, \cdot), \forall \hat{\pi} \in \hat{\Pi}_0$.

### 2.2 Hyperparameters

**Evolutionary algorithm for upper and lower bounds.** We use the EA to arrive at the upper and lower bounds on the preference score for each preference for each environment. We present the specific configurations for the hyperparameters in Table 1. This table includes essential parameters such as the mutation rate, crossover rate, population size, and the number of generations, among others.

| Hyperparameter | PotholeWorld-v0 | CartPole-v1 |
|---|---|---|
| Expert Algorithm | DQN | DQN |
| Number of Batch Rollouts | 10 | 10 |
| Maximum Samples | 200000 | 200000 |
| Maximum Iterations | 20 | 20 |
| Training Fraction | 0.8 | 0.8 |
| Reweighting | True | True |
| Number of Test Rollouts | 50 | 50 |

Table 2: **VIPER algorithm hyperparameters.**

| Hyperparameter | PotholeWorld-v0 | CartPole-v1 |
|---|---|---|
| Population Size | 75 | 75 |
| Number of Trees in Initial Population | 675 | 150 |
| Mutation Rate $\mu$ | 0.9 | 0.9 |
| Crossover Rate $\rho$ | 0.8 | 0.8 |
| Tournament Size $k$ | 3 | 3 |
| Max Generations | 75 | 75 |
| Expert for Imitation Learning | DQN | DQN |
| EA Start (Query Number) | 5 | 5 |
| EA Spacing (every $k$ steps) | 5 | 5 |

Table 3: **Hyperparameters of PASTEL.**

**VIPER.** The hyperparameters specific to the VIPER algorithm are detailed in Table 2, including the number of batch rollouts, maximum samples, and the fraction of data used for training. These number of batch rollout controls how many expert trajectories are generated during training, providing the data needed to distill the policy. The maximum number of samples limits the total state-action pairs collected. The fraction of data used for training determines how much of the collected dataset is used in each iteration, balancing the need for immediate policy updates with the potential benefit of retaining data for future use. The maximum iterations specifies how many rounds of VIPER are completed. Reweighting means that the algorithm reweights the samples used to train the decision-tree policy according to the associated Q-values.

**PASTEL.** In Table 3, we provide the hyperparameters for PASTEL. Both environments, PotholeWorld-v0 and CartPole-v1, share similar configurations, with the only difference that PotholeWorld-v0 uses more trees in the initial (VIPER-generated) population.

| $\widehat{\pi}$ | Reward | Depth | Num. Leaves | Takes Action | Score $(\widehat{\theta}_1)$ | Score $(\widehat{\theta}_2)$ | Score $(\widehat{\theta}_3)$ |
|---|---|---|---|---|---|---|---|
| 1 | 24 | 4 | 15 | 1 | 0.55 | 0.57 | 0.55 |
| 2 | 23 | 4 | 11 | 1 | 0.56 | 0.58 | 0.56 |
| 3 | 19 | 2 | 4 | 1 | 0.99 | 0.87 | 0.99 |
| 4 | 8 | 2 | 3 | 1 | 0.99 | 0.83 | 0.99 |
| 5 | 5 | 0 | 1 | 0 | 2.48 | 2.57 | 2.48 |
| 6 | 24 | 5 | 7 | 1 | 0.35 | 0.44 | 0.35 |
| 7 | 5 | 1 | 2 | 1 | 1.2 | 0.97 | 1.20 |



Figure 1: **Example illustration of how one can leverage and inspect the learned preferences and features.** This figure shows an example from a PASTEL run in the PotholeWorld-v0 environment. The top table shows the different values of the features (Reward, Depth, etc.) that correspond to each policy $\hat{\pi}$ in the Pareto frontier. It also shows the overall score according to each $\hat{\theta}$, as well as the overall $\hat{\pi}$ that is selected as the best (in green) and the ones that were not selected (in pink). The three plots at the bottom show the directions of the $\hat{\theta}$ for different combinations of feature comparisons. $\hat{\theta}_1$ and $\hat{\theta}_3$ are very similar, so there are only 2 visible vectors on the plots.

# 3 Additional Results

## Additional Example of User Interaction

We provide another visualization showing how a user can interact with the learned preference models and policy features in Figure 1. This figure shows an example from a PASTEL run in PotholeWorld-v0. The top part of the figure showcases that the user can investigate which of the $\hat{\theta}$ voted for each of the $\hat{\pi}$, in addition to the overall winner. The bottom part highlights that the user can further decompose the assessment into pairwise comparisons over features to see the direction that the $\hat{\theta}$ point in for each combination. The Pareto frontier filtering means that users would only have to inspect a small subset of alternatives, rather than hundreds or thousands.

# Appendix: Learning from Human Feedback for Fuzzy Tasks

|  | TrueSkill | Elo | Individual Analysis |
|---|---|---|---|
| Easily compare algorithms | ✓ | ✓ | ✗ |
| Handle multi-algorithm comparisons | ✓ | ✗ | ✗ |
| Provide uncertainty estimates | ✓ | ✗ | ✗ |

Table 1: **Comparison of different potential evaluation systems for the BASALT benchmark.** We use TrueSkill for the BASALT evaluation protocol.

# 1 BASALT Benchmark

## 1.1 TrueSkill

We provide an extensive discussion on the use of TrueSkill over other procedures and evaluation metrics. For a summary of these comparisons, see Table 1.

**TrueSkill vs. Elo**

The Elo rating system [88] is a popular zero-sum system initially created for chess. The gain in rating points of a player, or algorithm, after a win is equivalent to the loss in points for their competitor. The outcome of the game and the difference in ratings between algorithms influence the specific number of points transferred. Elo was designed for two-player games, whereas TrueSkill can handle multi-player games. In the case of BASALT, this means that human evaluations can be performed over more than two algorithms, which may be more cost-effective for researchers to deploy.

Furthermore, Elo does not incorporate uncertainty estimates for an algorithm's performance, whereas TrueSkill does. In other words, it only provides point estimates without uncertainty quantification. When comparing two algorithms with similar Elo ratings, we cannot determine whether they are genuinely close in performance or whether we simply lack sufficient data. TrueSkill addresses this by maintaining explicit uncertainty estimates that reflect our confidence in each algorithm's skill level.

Finally, Elo's fixed learning rate means that all agents adjust their ratings at the same pace regardless of how much data supports their current rating. TrueSkill uses uncertainty-weighted updates: algorithms with high uncertainty (typically newer or less-tested methods) experience larger rating changes, while algorithms with low uncertainty (well-established through many evaluations) show more stable ratings.

**TrueSkill vs. Independent Algorithm Performance Assessment**

Due to the difficulty of the BASALT tasks, independently assessing the performance of an algorithm is challenging. For example, in the absence of any comparisons, how would one assess the performance of an algorithm that only partially completes a task? How could we ensure that the score is calibrated for the assessment of other algorithms? These standalone assessments do not offer ways to directly compare multiple algorithms or adjust the rating of an algorithm based on its performance over time, as TrueSkill and Elo do.

## 1.2 Best Practices

To fairly compare agents on a benchmark, the parameters of that benchmark must be concretely defined. Previously, many techniques claimed to solve MineRL ObtainDiamond. Often, these successes involved giving the agent access to the full game state instead of pixel observations, building in human priors through extensive action shaping, or simplifying the underlying environment dynamics. Changes like these should be clearly stated and applied equally to compared methods.

To avoid such discrepancies in the future, we specify some methodological best practices when using the BASALT benchmark: i) only pixel observations provided by the four environments should be used, ii) the environments should not be modified in any way, iii) action shaping [156] is permitted, as long as it applied to all methods, iv) algorithms should be evaluated with TrueSkill [130] using the specific hold-out test seeds provided below, and v) the final evaluations must be conducted with human evaluators. For further guidance on how to fairly compare algorithms, we refer the reader to recent work [251].

The evaluation environment seeds for each task are:

- `MineRLBasaltFindCave-v0`: 14169, 65101, 78472, 76379, 39802, 95099, 63686, 49077, 77533, 31703, 73365.

- `MineRLBasaltMakeWaterfall-v0`: 95674, 39036, 70373, 84685, 91255, 56595, 53737, 12095, 86455, 19570, 40250.

- `MineRLBasaltCreateVillageAnimalPen-v0`: 21212, 85236, 14975, 57764, 56029, 65215, 83805, 35884, 27406, 5681265, 20848.

- `MineRLBasaltBuildVillageHouse-v0`: 52216, 29342, 67640, 73169, 86898, 70333, 12658, 99066, 92974, 32150, 78702.

# 2 Datasheet for BEDD

We present a datasheet [101] for BEDD, which includes details on our 651GB human demonstrations dataset (Demonstrations Dataset) and our TrueSkill evaluations dataset (Evaluation Dataset).

## Motivation

### For what purpose was the dataset created?

The Demonstrations Dataset was provided as a training dataset for competitors to use for the BASALT 2022 competition [213], while the Evaluation Dataset was generated at the end of the competition when evaluating submissions. Below, we describe further purposes for the datasets.

**Demonstrations Dataset.** We believe that this set of successful demonstrations could be useful for training machine learning algorithms to perform complex tasks in Minecraft. All task completions can be thought of as having a positive label denoting successful completion. The dataset could serve as a basis for generating pretrained embeddings, similar to the VPT method. These embeddings could be further finetuned with additionally collected demonstrations, corrections, or other feedback modalities to further refine the behavior. This dataset could also be leveraged for a wealth of insights. Since these tasks are complex and require sequential dependencies, one could analyze and extract the crucial subtasks required to complete the full task. This analysis would support the development of various hierarchical methods. This task decomposition would support further subtask-specific analysis. Understanding the required time to complete each subtask would help identify which tasks take longer or are completed more quickly. This information could provide meaningful insights into task difficulty.

**Evaluation Dataset.** We envision this dataset to be used for both analysis and development. The detailed task-specific questions may be useful for understanding which factors most strongly influence overall assessments of behavior. We envision this dataset to be used to develop reward models that better align with human preferences. This aim aligns with the BASALT benchmark of both understanding and utilizing human preferences to solve fuzzy tasks. Training reward models on pairwise comparisons over algorithms may reveal the latent structure of human preferences over the task space. A more nuanced understanding would emerge as more of the task-specific questions and long-form justification of responses are processed and understood. This training regime would promote the development of more intuitive and human-aligned algorithms.[1] Furthermore, we envision that this dataset could provide a wealth of insights into more general human preferences, which could prove invaluable when transferred or used to warm-start reward models in different domains. By identifying these underlying general preferences, we could inform the design of more universally acceptable and effective AI systems (bearing in mind individual, cultural, and other important differences).

---

[1]We again stress that the goal is not to *completely replace* people. Instead, automated evaluations serve as a way to sidestep the issues with existing evaluation methods, which involve laborious manual assessments.

### Who created the dataset?

The dataset was created by Anssi Kanervisto (Microsoft Research), Stephanie Milani (Carnegie Mellon University), Karolis Ramanauskas (University of Bath), Byron V. Galbraith (Seva Inc.), Steven H. Wang (ETH Zürich), Sander Schulhoff (University of Maryland), Brandon Houghton (OpenAI), Sharada Mohanty (AIcrowd), and Rohin Shah. The dataset was not created on the behalf of any entity.

### Who funded the creation of the dataset?

The FTX Future Fund Regranting Program, Encultured.ai, and Microsoft funded the creation of the Evaluation Dataset, prizes, and compute. OpenAI sponsored the creation of the Demonstrations Dataset.

## Composition

### What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?

The Demonstrations Dataset contains videos of agent trajectories and associated actions taken. The Evaluation Dataset contains videos of trajectories and human evaluations on multiple factors such as "Found Cave" or "Least Village Harm".

Please refer to Section 7.4 for more information. Table 7.1 contains a breakdown of information on the Demonstrations Dataset and Table 7.2 contains a breakdown of information on the Evaluation Dataset. Further information is in Appendices 3.2 and 4.

### Is there a label or target associated with each instance?

The Evaluation Dataset contains labeled data. Agents are evaluated across a range of criteria (Section 7.6), which are labels. Additionally, the Demonstrations Dataset contains videos of trajectories paired with 'labels' of the action(s) being executed at each frame. Our codebase contains files that describe how each trajectory ends (e.g., by ESC key or death.).

### Is any information missing from individual instances?

Not to our knowledge.

### Are there recommended data splits (e.g., training, development/validation, testing)?

No.

**Are there any errors, sources of noise, or redundancies in the dataset?**

There are instances of contractors doing a task incorrectly in the Demonstrations Dataset. For example, in some instances, contractors found a cave and did not stop the episode or made a waterfall and stared at it instead of ending the episode. Please see Appendix 3.2 for more possible issues. We did not manually check the entire dataset, so it may contain additional anomalous activities.

**Do/did we do any data cleaning on the dataset?**

We did not. All data is presented exactly as collected. We provide information on which demonstrations may contain human errors in the repository. This information may be then used to conduct data cleaning by the end user of the dataset.

## Collection Process

### How was the data associated with each instance acquired?

The Demonstrations Dataset was collected using a modded, closed source Minecraft version that recorded users' game frames and actions. Contractors watched agent gameplay side by side in order to evaluate them and create the Evaluation Dataset.

### Who was involved in the data collection process and how were they compensated?

For the Demonstrations Dataset, 20 contractors were hired and paid 20 USD an hour. Additionally, some members of the BASALT team collected data for this dataset. For the Evaluation Dataset, 65 high quality crowdsource workers from Amazon Mechanical Turk were paid to collect the data at approximately 15 USD an hour.

### Over what timeframe was the data collected?

The Demonstrations Dataset was collected between April 2022 and July 2022. The Evaluation Dataset was collected between November 2022 and January 2023.

## Uses

### Has the dataset been used for any tasks already?

Both datasets were used for the BASALT 2022 competition. Please see more details in the competition retrospective [213] and this competitor white paper [201].

**Is there a repository that links to any or all papers or systems that use the dataset?**

This page contains a list of papers and projects that use MineRL. Only some projects from 2022 onward use BEDD data. Competitor papers from the 2022 competition [201] only used the Demonstrations Dataset (the Evaluation Dataset was created when evaluating these submissions).

**Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**

We do not believe so because the data for both datasets was collected from paid contractors and high-quality paid crowdsourcers.

## Distribution

### Will the dataset be distributed to third parties?

Yes, it is free and available online.

### How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?

The Evaluation Dataset exists on Zenodo (DOI: 10.5281/zenodo.8021960) as a zip file.

The Demonstrations Dataset exists on an OpenAI server as JSONL and MP4 files and does not have a DOI.

All data is under the MIT license.

### Have any third parties imposed IP-based or other restrictions on the data associated with the instances?

No.

### Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?

No.

## Maintenance

### Who will be supporting/hosting/maintaining the dataset?

The authors on this paper will provide needed maintenance to the datasets (Demonstrations Dataset, Evaluation Dataset). We do not expect much maintenance to be needed as we will not be adding data to the dataset. However, we will accept PRs from users who have improvements to make to the supporting codebase.

### How can the owner/curator/manager of the dataset be contacted (e.g., email address)?

Please email us at basalt@minerl.io

### Is there an erratum?

There is not, but 1) we mention potential issues with the data in this datasheet, and 2) we provide a list of data within the dataset that we believe to be invalid due to issues such as not properly completing the given task (Appendix 3.2).

### Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?

Yes, but we expect minimal updates to be required, as we do not intend to add more data to the dataset.

# 3  Demonstrations Dataset Details

## 3.1  Human Data Collection Details

Contractors were originally hired through Upwork by responding to the following job posting:

> We are looking for people who want to get paid to play Minecraft. We will want you to describe some things about your experience, so our ideal candidate would possess Native English fluency and have a microphone. You'll need to install java, download a modified version of Minecraft (that collects and uploads your play data and voice), and play Minecraft survival mode! Paid per hour of gameplay. Prior experience in Minecraft is not necessary. We do not collect any data that is unrelated to Minecraft from your computer.

Prior to recording this dataset, all contractors had been working on other datasets in Minecraft[2], and were proficient in taking English language instructions via UpWork and using the provided recording scripts for generating demonstration data. The contractors were initially hired in the context of collecting narrated Minecraft play. However, narrations were not requested for the BASALT dataset, and any incidentally collected narrations will not be released. In total, 20 contractors were recruited and paid 20 USD per hour. We allocated a total of 10k USD for this. To supplement this dataset, three members of BASALT contributed data. Below are the exact transcripts we used to instruct the contractors for the four tasks.

### FindCave

> Task 1 - Recorder version find-cave Look around for a cave. When you are inside one, quit the game by opening main menu and pressing "Save and Quit To Title". You are not allowed to dig down from the surface to find a cave.
> Timelimit: 3 minutes.
> Example recordings: `https://www.youtube.com/watch?v=TclP_ozH-eg`

### MakeWaterfall

> After spawning in a mountainous area with a water bucket and various tools, build a beautiful waterfall and then reposition yourself to "take a scenic picture" of the same waterfall, and then quit the game by opening the menu and selecting "Save and Quit to Title"
> Timelimit: 5 minutes.
> Example recordings: `https://youtu.be/NONcbS85NLA`

---

[2]Previous Minecraft experience was not a hard requirement in the job posting. However, in practice, those without experience did not continue their work.

**CreateVillageAnimalPen**

> After spawning in a village, build an animal pen next to one of the houses in a village. Use your fence posts to build one animal pen that contains at least two of the same animal. (You are only allowed to pen chickens, cows, pigs, sheep or rabbits.) There should be at least one gate that allows players to enter and exit easily. The animal pen should not contain more than one type of animal. (You may kill any extra types of animals that accidentally got into the pen.) Don't harm the village. After you are done, quit the game by opening the menu and pressing "Save and Quit to Title".
>
> You may need to terraform the area around a house to build a pen. When we say not to harm the village, examples include taking animals from existing pens, damaging existing houses or farms, and attacking villagers. Animal pens must have a single type of animal: pigs, cows, sheep, chicken or rabbits.
>
> The food items can be used to lure in the animals: if you hold seeds in your hand, this attracts nearby chickens to you, for example.
>
> Timelimit: 5 minutes. Example recordings: `https://youtu.be/SL07sep7B08`

**BuildVillageHouse**

> Taking advantage of the items in your inventory, build a new house in the style of the village (random biome), in an appropriate location (e.g. next to the path through the village), without harming the village in the process. Then give a brief tour of the house (i.e. spin around slowly such that all of the walls and the roof are visible).
>
> You start with a stone pickaxe and a stone axe, and various building blocks. It's okay to break items that you misplaced (e.g. use the stone pickaxe to break cobblestone blocks). You are allowed to craft new blocks.
>
> Please spend less than ten minutes constructing your house.
>
> You don't need to copy another house in the village exactly (in fact, we're more interested in having slight deviations, while keeping the same "style"). You may need to terraform the area to make space for a new house. When we say not to harm the village, examples include taking animals from existing pens, damaging existing houses or farms, and attacking villagers.
>
> After you are done, quit the game by opening the menu and pressing "Save and Quit to Title".
>
> Timelimit: 12 minutes.
>
> Example recordings: `https://youtu.be/WeVqQN96V_g`

## 3.2 Documentation of Challenges

Like most other datasets, this one contains some issues. Instead of waiting for users to discover them, we preemptively investigate the dataset ourselves. We document the issues below. This is done in the interest of transparency and to ensure the dataset is maximally useful.

## Episode Boundaries

The recording software employed for our dataset splits video and action label files into 5-minute segments. This has no bearing on the FindCave task, as its time limit is 3 minutes. However, the time limits for MakeWaterfall and CreateVillageAnimalPen tasks are 5 minutes, and occasionally episodes exceed this limit by up to 3 seconds, triggering the video splitter. For the BuildVillageHouse task, the time limit is 12 minutes, causing some episodes to be divided into three parts.

Certain training or analysis methods necessitate knowledge of episode boundaries, requiring a reliable method for identifying which videos belong to the same episode. While file labels being unique to an episode would simplify this process, this is not the case. We investigated various systematic approaches to detect where an episode ends and a new one begins, but none proved 100% reliable. Some challenges we faced include: different episodes sharing the same file ID; video splitter not generating filenames with timestamps exactly 5 minutes apart; ESC keys not consistently triggering episode ends, because ESC is also used to close inventory; some episodes concluding in exactly 5 minutes; and a few episodes missing their first 5 minutes, rendering the loading screen in the initial video frame an unreliable boundary indicator. Our most reliable and straightforward solution was to consider two files as part of the same episode if the first file is exactly 5 minutes long, no ESC key is pressed with the GUI closed, and the filename ID is identical. This approach resulted in an error rate of less than 1%, based on manual inspection of the first frames of all videos in the final split.

Leveraging this heuristic, we produced a file for each task containing a list of episodes, with each episode having an associated file list. We also included the step count per episode and two tags specifying (1) whether the episode ended with an ESC key press, (2) whether it was incomplete due to a saving error. Episode ending with an ESC key press indicates a successfully completed task. If an episode does not conclude with an ESC key press, the step count can be utilized to determine the type of episode end — a step count at or slightly above the time limit implies a timeout, while a lower step count indicates player death. The four files are located in the same repository mentioned above. This is the recommended way of using the dataset if episode boundaries are important for the training algorithm.

## Idiosyncratic Episodes

We also noticed some episodes, which pass our filters for being valid and complete episodes but have some unique characteristics. This might be useful to know for data cleaning purposes. We provide some examples below, including the associated file names.

- `gloppy-persimmon-ferret-3e42e8e14be0-20220716-190015` - FindCave episode finishes in under 2 seconds with an ESC press, likely a misclick.

- `squeaky-ultramarine-chihuahua-bbf328311fb8-20220726-132144` - FindCave episode where the player spawns, then falls into a ravine and dies in less than 5 seconds.

- `gloppy-persimmon-ferret-1d8dcc4e2446-20220716-144806` - successful FindCave episode, where the player finds a cave and hits ESC in under 3 seconds.

- `pokey-cyan-spitz-62e7b7415aaf-20220714-093544` - FindCave episode where the video lasts longer than the 3 minute time limit, but the action labels only cover 20 seconds.

- `whiny-ecru-cougar-f153ac423f61-20220712-192050` - FindCave episode, where the player seems to just play Minecraft, making stone tools and such, instead of finding a cave.

- `shabby-pink-molly-*` - a total of 67 FindCave episodes, where the player finds a cave, but does not press ESC to finish the episode. Instead the player proceeds to explore the cave, then exits it and goes looking for other caves.

- `thirsty-lavender-koala-479e09882ca6-20220717-203846` - MakeWaterfall episode, one of several, where the player finishes the waterfall, and stares at it for 3 minutes to timeout instead of pressing ESC.

These were the outliers we found by sorting the data based on various metrics and checking the extremes. They are rare, on the order of tens of episodes total in a dataset with over 13,000 episodes. It would add up to a total of less than 1%. Although we believe that these rare outliers would not influence training outcomes, we recommend that users of this dataset either remove this data from the training set or more carefully check these episodes before including them.

## Video Encoding Differences

We also noticed that the codecs used to encode the videos were not always the same, likely due to subtle differences in the systems contractors used to play the game to generate the data. Codecs are how the data in the videos gets compressed and decompressed. Roughly 88% use H.264 (Constrained Baseline Profile), while most of the remaining ones use H.264 (High Profile). Also, the bitrate of the videos varies. Most bitrates are roughly 4 Mbps, but there are a few outliers on both ends of the distribution. While these differences are imperceptible to the human eye, the training algorithms might pick up on them. Having different encodings has both benefits and issues for training algorithms. The benefits are robustness, adaptability and real-world applicability. The issues include biased training data and increased computational complexity.

# 4 Evaluation Dataset Details

This section contains details about the Evaluation Dataset. In this section, we hope to provide enough detail such that our evaluation pipeline can be easily reproduced.

## 4.1 Human Evaluation Details

We estimated that one answer would take 15 minutes. We paid 3.75 USD per HIT for a total of 15 USD per hour. In actuality, workers took 5.12 minutes on average to complete each evaluation, meaning the pay was closer to 43.95 USD per hour. We also provided bonuses to MTurk workers that helped us debug issues with the form during the evaluation. In total, we spent 14,849.16 USD on collecting this data.

The human judges view the following task description on MTurk.

> In this questionnaire, you will watch videos of different players completing tasks in Minecraft, and your task is to judge which of the players is more successful at completing the task. This will take roughly 15 minutes of your time.

### Qualification Criteria

We set the following criteria for selecting human judges. To preview the task, the judges must have had a $99\%$ or greater HIT accept rate and a minimum of $10,000$ completed HITs on MTurk. If the judges had these qualifications, they then took an 8-question Minecraft validity test to confirm that they had at least basic knowledge of Minecraft. We define passing this test as reaching $65\%$ or greater. This means that, for all checkboxes, the user correctly checked or unchecked at least $65\%$ of the boxes.

**The Minecraft qualification test.** For the purpose of completeness, we include this test. In this section, we detail the questions that we asked the human judges. To enable this questionnaire to be deployed, we have included deployable and more user-friendly versions of the quiz in our Github repository. Specifically, we provide a text file and an XML file. The former enables the questions to be more readily copied and pasted into new formats, while the latter slots easily into the MTurk UI. Before answering the quiz questions, the human judges viewed the following prompt:

> This is a Qualification Test to demonstrate your familiarity with Minecraft. This Qualification will enable you to accept HITs released by the BASALT team in relation to the NeurIPS 2022 BASALT Competition.

After reading this prompt, the human judges proceeded to the Minecraft Qualification Test. Figure 1 shows the quiz. These questions were generated by a member of our team who has extensive experience playing Minecraft.

Q1: (select all that apply) A bed in Minecraft...



[ ] can be used to speed up the passage of the night **(correct)**
[ ] can be used to change the respawn location **(correct)**
[ ] Is built using an item dropped by an in-game animal **(correct)**
[ ] Can be built at a furnace, but not at a crafting table.

Q2: (select all that apply) The following events can drop experience orbs:



[ ] Killing monsters **(correct)**
[ ] Mining trees
[ ] Jumping on a coal ore block
[ ] Mining coal **(correct)**

Q3: (select all that apply) The following blocks drop an item when mined by hand:
[ ] Wood log **(correct)**
[ ] Wood plank **(correct)**
[ ] Coal ore
[ ] Iron ore
[ ] Stone

(a) Page 1 of Minecraft Qualification Test

Q4: (select all that apply) Which of the following mobs can damage players:



[ ] Skeletons **(correct)**
[ ] Zombies **(correct)**
[ ] Sheep
[ ] Pigs
[ ] Creepers **(correct)**

Q5: (select all that apply) Which of the following items can be eaten by the player:
[ ] Apples **(correct)**
[ ] Dirt
[ ] Beef **(correct)**
[ ] Wheat
[ ] Bread **(correct)**

Q6: (select all that apply) Which of the following are ways that the player can die in vanilla Minecraft:
[ ] Fall damage **(correct)**
[ ] Drowning **(correct)**
[ ] Burning **(correct)**
[ ] Thirst
[ ] Getting hit repeatedly with a stick **(correct)**

(b) Page 2 of Minecraft Qualification Test

Q7: (select all that apply) What is this mob called?



[ ] Creeper **(correct)**
[ ] Beeper
[ ] Leaper
[ ] Reaper
[ ] Peeper

Q8: (select all that apply) What types of items can you make on a crafting table using only this block?



[ ] Planks **(correct)**
[ ] Sticks **(correct)**
[ ] Bone meal
[ ] Apples
[ ] Shield

(c) Page 3 of Minecraft Qualification Test

Q9: (select all that apply) Jumping onto the following block type does what?



[ ] Makes you drown immediately
[ ] Refreshes your thirst bar
[ ] Heals you
[ ] Removes burning effects **(correct)**
[ ] Stops fall damage **(correct)**

(d) Page 4 of Minecraft Qualification Test

Figure 1: **Minecraft Qualification Test**. We ask participants to complete this test to ensure that they have sufficient knowledge about Minecraft to participate in the study.

**Protocol**

Each time a human judge accepts our HIT, they are redirected to the AICrowd site that collects their answer. The task is randomly chosen from the four available ones, then it is assigned to the human judge. To choose the two videos to show the judge, we aim to maximally reduce uncertainty using the TrueSkill ranking of agents. We ensure that the compared videos are always in the same Minecraft seed (i.e., the same starting location and surroundings).

**MTurk task instructions.**   Below is the full text of the MTurk task instructions given to participants.

> Welcome to the MineRL BASALT 2022 evaluation questionnaire! This will take roughly 15 minutes of your time. Requirements: Knowledge of Minecraft (at least of 5 hours of gameplay time with Minecraft).
> In this questionnaire, you will watch videos of different players completing tasks in Minecraft, and your task is to judge which of the players is more successful at completing the task.
> Videos are shown in pairs, and your task is to select which one of the two is better at solving the task. The page shows the task description. You are also given a set of more specific questions which may help you decide which of the two players completes the task better.
> You may refer to the Example Videos section, for some examples of what are considered as good executions of the said task, and why.
> Your answers will be used in the following ways:
> To rank the solutions in the MineRL BASALT 2022 competition. The answers will be included in the final report of the competition. The answers may be shared publicly to support the research. No personal information is collected.
> You may complete the same questionnaire multiple times, but may be asked to judge players completing a different task. So please ensure that you take note of the Task you are submitting the responses for.

**Task details.**   Each task contained some common questions and some different questions. For all tasks, the human judges were asked to evaluate which player was better overall and which player appeared more human-like. Some tasks contained more task-specific questions than others. For example, BuildVillage-House contained a single question about whether the players harmed the villagers; in contrast, CreateVil-lageAnimalPen contained six questions to determine proper task completion, such as whether the pen contained at least two animals of the same type. Figures 2 to 5 show screenshots of the form for each task that the human judges saw.

## 4.2  Additional Details of the Evaluation Dataset Analysis

We now present the mapping of per-task questions to factors, which are in Figure 7.4 from the main text and elaborated in Tables 2 to 5.

*Appendix: Learning from Human Feedback for Fuzzy Tasks*

## NeurIPS BASALT 2022 competition evaluation

Welcome to the MineRL BASALT 2022 evaluation questionnaire! This will take roughly **15 minutes of your time**.

**Requirements**: Knowledge of Minecraft (at least of 5 hours of gameplay time with Minecraft).

In this questionnaire, you will watch videos of different players completing tasks in Minecraft, and your task is to judge which of the players is more successful at completing the task.

Videos are shown in pairs, and your task is to select which one of the two is better at solving the task. The page shows the task description. You are also given a set of more specific questions which may help you decide which of the two players completes the task better.

You may refer to the **Example Videos** section, for some examples of what are considered as good executions of the said task, and why.

Your answers will be used in the following ways:

- To rank the solutions in the MineRL BASALT 2022 competition.
- The answers will be included in the final report of the competition.
- The answers may be shared publicly to support the research.

No personal information is collected.

You may complete the same questionnaire multiple times, but may be asked to judge players completing a different task. So please ensure that you take note of the Task you are submitting the responses for.

> Task description the players received
>
> ## Find Cave
>
> - The player should search for a cave, and terminate the episode when it is inside one.
>
> **Example Videos**

**Left Player**



685affc4-787c-43a0-a9f6-22a4de3bcdb3

**Right Player**



8a10233d-377a-49f6-847c-554cacc1478b

## Question Set #1

| | Left Player | Right Player |
|---|---|---|
| Direct questions | | |
| **Q1.** Did this player find and enter a cave? | ☐ | ☐ |

## Question Set #2

| | | | | |
|---|---|---|---|---|
| **Q1.** Which player found a cave the fastest? (If neither found a cave, that is a draw.) | Left Player | Draw | Right Player | N/A |
| **Q2.** Which player moved more quickly and efficiently? | Left Player | Draw | Right Player | N/A |
| **Q3.** Which player was better at looking for caves in areas they hadn't already explored? | Left Player | Draw | Right Player | N/A |
| **Q4.** Which player was better at going to areas where it is more likely to find caves? | Left Player | Draw | Right Player | N/A |
| **Q5.** Which player was better at noticing potential caves that entered its field of vision? | Left Player | Draw | Right Player | N/A |
| **Q6.** Which player was better at realizing when it has successfully found a cave? (In other words, which player was better at properly ending the minigame once it had entered a cave?) | Left Player | Draw | Right Player | N/A |
| **Q7.** Which player seemed more human-like (rather than a bot or computer player)? | Left Player | Draw | Right Player | N/A |

## Which player is better overall ?*

| Left Player | Right Player | Draw |
|---|---|---|

Justify your answer with a minimum of 100 characters.*

Submit

**Figure 2: Page that the judges view when assessing agent behavior for the FindCave task.**

190

Figure 3: **Page that the judges view when assessing agent behavior for the MakeWaterfall task.**

Figure 4: **Page that the judges view when assessing agent behavior for the CreateVillageAnimalPen task.**

## NeurIPS BASALT 2022 competition evaluation

Welcome to the MineRL BASALT 2022 evaluation questionnaire! This will take roughly **15 minutes of your time**.

**Requirements**: Knowledge of Minecraft (at least of 5 hours of gameplay time with Minecraft).

In this questionnaire, you will watch videos of different players completing tasks in Minecraft, and your task is to judge which of the players is more successful at completing the task.

Videos are shown in pairs, and your task is to select which one of the two is better at solving the task. The page shows the task description. You are also given a set of more specific questions which may help you decide which of the two players completes the task better.

You may refer to the **Example Videos** section, for some examples of what are considered as good executions of the said task, and why.

Your answers will be used in the following ways:

- To rank the solutions in the MineRL BASALT 2022 competition.
- The answers will be included in the final report of the competition.
- The answers may be shared publicly to support the research.

No personal information is collected.

You may complete the same questionnaire multiple times, but may be asked to judge players completing a different task. So please ensure that you take note of the Task you are submitting the responses for.
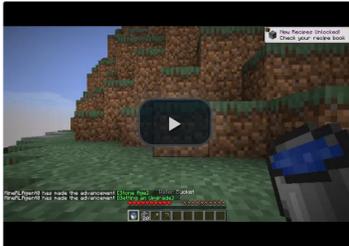
Task description the players received

### Build Village House

- The player should build a new house in the style of the village (random biome), in an appropriate location (e.g. next to the path through the village), without harming the village in the process.
- Then the player should give a brief tour of the house (i.e. spin around slowly such that all of the walls and the roof are visible).

Example Videos

**Left Player**

e5fa2436-44af-4f0a-8b23-54ce6fec9239

**Right Player**

0993b32e-59fa-4406-8748-4872ccaab2a3

## Question Set #1

| | Left Player | Right Player |
|---|---|---|
| Direct questions | | |
| **Q1.** Did the player harm the village? Examples include taking animals from existing pens, damaging existing houses or farms, and attacking villagers. | ☐ | ☐ |

## Question Set #2

| | | | | |
|---|---|---|---|---|
| **Q1.** Which player chose a better location for their house? | Left Player | Draw | Right Player | N/A |
| **Q2.** Which player's structure seemed most like a house? | Left Player | Draw | Right Player | N/A |
| **Q3.** Which player was better at removing unnecessary blocks (or never placing unnecessary blocks)? | Left Player | Draw | Right Player | N/A |
| **Q4.** Which player was better at using the appropriate type of blocks (i.e. the ones that are used in other houses in the village)? | Left Player | Draw | Right Player | N/A |
| **Q5.** Which player's house better matched the "style" of the village? | Left Player | Draw | Right Player | N/A |
| **Q6.** Which player built the better-looking house? | Left Player | Draw | Right Player | N/A |
| **Q7.** Which player seemed more human-like (rather than a bot or computer player)? | Left Player | Draw | Right Player | N/A |

## Which player is better overall ?*

Left Player | Right Player | Draw

Justify your answer with a minimum of 100 characters.*

Submit

Figure 5: **Page that the judges view when assessing agent behavior for the BuildVillageHouse task.**

| Question | Factor |
|---|---|
| **Quantitative** | |
| Did this player find and enter a cave? | Found Cave |
| **Qualitative** | |
| Which player found a cave the fastest? (If neither found a cave, that is a draw) | Found Cave Faster |
| Which player moved more quickly and efficiently? | More Quick and Efficient Movement |
| Which player was better at looking for caves in areas they hadn't already explored? | Better Cave Search in Unknown Areas |
| Which player was better at going to areas where it is more likely to find caves? | Better Navigation to Cave Areas |
| Which player was better at noticing potential caves that entered its field of vision? | Better Cave Detection |
| Which player was better at realizing when it had successfully found a cave? (In other words, which player was better at properly ending the minigame once it had entered a cave?) | Better Cave Perception |
| Which player seemed more human-like (rather than a bot or computer player)? | More Human Like |

Table 2: **FindCave: Mapping from evaluation questions to attributes.**

# 5 Extended Discussion and Future Work

Here we provide a longer discussion about the limitations of our work and suggest future work that stems from these limitations.

**Low-level data discrepancies.** As with any dataset, ours is not without issues. We discussed some limitations of the Demonstrations Dataset in a previous section of the appendix. We found that it is often challenging to delineate episode boundaries with data recording tools and splitting methods. In contrast, with games like Atari, episode boundaries are easily provided by the simulator. Furthermore, since the data is temporally-extended and human-generated, there are idiosyncrasies that will necessarily arise. We emphasize that rather than deploying the dataset and making users find these issues, we dedicated time to investigating these issues and proposing solutions.

| Question | Factor |
|---|---|
| **Quantitative** | |
| Did this player create a waterfall? | Created Waterfall |
| Did this player end the video while looking at a player-constructed waterfall? | Looked at Waterfall |
| | |
| **Qualitative** | |
| Which player moved more efficiently? | More Quick and Efficient Movement |
| Which player chose a better location for their waterfall? (If neither player created a waterfall, select "Draw".) | Better Location for Waterfall |
| Which player took a better "picture" of the waterfall? (If neither player took a picture of a player-constructed waterfall, select "Draw") | Better Picture of Waterfall |
| Which player seemed more human-like (rather than a bot or computer player)? | More Human Like |

Table 3: **MakeWaterfall : Mapping from evaluation questions to attributes.**

**Demographic information.** In both studies, we did not collect any additional demographic information about participants. On one hand, this lack of personally-identifiable information decreases the likelihood of deanonymization of the participants, which is important even in a study with relatively low stakes such as this one. On the other hand, there is a missed opportunity to critically analyze the influence of demographic or other factors on either the produced demonstrations or the assessments. Future work might include some of these demographic details to produce a more in-depth analysis of the assessments.

**English-language focus.** One aspect of our study that offers both a specific focus and an avenue for exploration is the use of English-language descriptions for tasks. Incorporating multiple languages into the dataset would make the benchmark more globally applicable and uncover new insights about how language and cultural context influence human feedback.

| Question | Factor |
|---|---|
| **Quantitative** | |
| Did the player harm the village? Examples include taking animals from existing pens, damaging existing houses or farms, and attacking villagers. | Least Village Harm |
| | |
| **Qualitative** | |
| Which player chose a better location for their house? | Better House Location |
| Which player's structure seemed most like a house? | More House-Like Structure |
| Which player was better at removing unnecessary blocks (or never placing unnecessary blocks)? | Intentional Block Placing |
| Which player was better at using the appropriate type of blocks (i.e., the ones that are used in other houses in the village)? | Appropriate Block Placing |
| Which player's house better matched the "style" of the village? | Better Style Matching |
| Which player built the better-looking house? | More Attractive House |
| Which player seemed more human-like (rather than a bot or computer player)? | More Human Like |

Table 4: **BuildVillageHouse : Mapping from evaluation questions to attributes.**

| Question | Factor |
|---|---|
| **Quantitative** | |
| Did the player harm the village? Examples include taking animals from existing pens, damaging existing houses or farms, and attacking villagers. | Least Village Harm |
| Did the player build an enclosed space (pen) from which animals could not escape? (For this question, it is okay if the pen did not contain animals). | Better Enclosed Space |
| Did the player's pen contain at least two animals of the same type? | Correct Number of Animals |
| Did the player's pen contain only one type of animal (if pen contained no animals, answer is no. Monsters and villagers are not counted)? | Correct Type of Animals |
| Did the player's pen have at least one gate that would allow players to enter and exit (using a block to jump over fence does not count)? | Proper Exit |
| Was the player's pen built next to a house? | Next to House |
| | |
| **Qualitative** | |
| Which player chose a better location for their pen in the village? | Better Location |
| Which player searched more effectively for animals to pen? | More Effective Search |
| After finding the animals, which player penned them more effectively? | More Effective Penning |
| Which player seemed more human-like (rather than a bot or computer player)? | More Human Like |

Table 5: **CreateVillageAnimalPen: Mapping from evaluation questions to attributes.**

# Bibliography

[1] Hello gpt-4o. `https://openai.com/index/hello-gpt-4o/`.

[2] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.

[3] Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. Cooperation, competition, and maliciousness: Llm-stakeholders interactive negotiation. *Advances in Neural Information Processing Systems*, 37:83548–83599, 2024.

[4] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[5] EU Artificial Intelligence Act. The eu artificial intelligence act, 2024.

[6] Peter Aigner and Brenan McCarragher. Human integration into robot control utilising potential fields. In *Proceedings of International Conference on Robotics and Automation*, volume 1, pages 291–296. IEEE, 1997.

[7] Eloi Alonso, Ubisoft La Forge, Maxim Peter, David Goumard, and Joshua Romoff. Deep reinforcement learning for navigation in AAA video games. In *Challenges of Real-World Reinforcement Learning NeurIPS Workshop*, pages 1–13, Montreal, Canada, 2020. NeurIPS.

[8] Ofra Amir, Finale Doshi-Velez, and David Sarne. Agent strategy summarization. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 1203–1207. International Foundation for Autonomous Agents and Multiagent Systems, 2018.

[9] Andrew Anderson, Jonathan Dodge, Amrita Sadarangani, Zoe Juozapaitis, Evan Newman, Jed Irvine, Souti Chattopadhyay, Alan Fern, and Margaret Burnett. Explaining reinforcement learning to mere mortals: An empirical study. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 2019.

[10] Ron Artstein and Massimo Poesio. Inter-coder agreement for computational linguistics. *Computational linguistics*, 34(4):555–596, 2008.

[11] Dilip Arumugam, Jun Ki Lee, Sophie Saskin, and Michael L Littman. Deep reinforcement learning from policy-dependent human feedback. *arXiv preprint arXiv:1902.04257*, 2019.

[12] Christian Arzate Cruz and Jorge Adolfo Ramirez Uresti. Hrlb2: A reinforcement learning based framework for believable bots. *Applied Sciences*, 8(12):2453, 2018.

[13] Edmond Awad, Sohan Dsouza, Richard Kim, Jonathan Schulz, Joseph Henrich, Azim Shariff, Jean-François Bonnefon, and Iyad Rahwan. The moral machine experiment. *Nature*, 563(7729):59–64, 2018.

[14] Thomas Bäck and Hans-Paul Schwefel. An overview of evolutionary algorithms for parameter optimization. *Evolutionary Computation*, 1(1):1–23, 1993.

[15] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

[16] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022.

[17] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy extraction. In *Advances in Neural Information Processing Systems*, pages 2494–2504, 2018.

[18] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.

[19] William H. Beluch, Tim Genewein, Andreas Nürnberger, and Jan M. Köhler. The power of ensembles for active learning in image classification. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9368–9377, 2018.

[20] Viktor Bengs, Róbert Busa-Fekete, Adil El Mesaoudi-Paul, and Eyke Hüllermeier. Preference-based online learning with dueling bandits: A survey. *Journal of Machine Learning Research*, 22(7):1–108, 2021.

[21] Jon Louis Bentley, Hsiang-Tsung Kung, Mario Schkolnick, and Clark D Thompson. On the average number of maxima in a set of vectors and applications. *Journal of the ACM (JACM)*, 25(4):536–543, 1978.

[22] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.

[23] Michael S Bernstein, Joel Brandt, Robert C Miller, and David R Karger. Crowds in two seconds: Enabling realtime crowd-powered interfaces. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 33–42, 2011.

[24] Tom Bewley and Jonathan Lawry. Tripletree: A versatile interpretable representation of black box agents and their environments. *CoRR, abs/2009.04743*, 2020.

[25] Sushrut Bhalla, Sriram Ganapathi Subramanian, and Mark Crowley. Deep multi agent reinforcement learning for autonomous driving. In *Proceedings of the Canadian Conference on Artificial Intelligence*, pages 67–78. Springer, 2020.

[26] Ioana Bica, Daniel Jarrett, Alihan Hüyük, and Mihaela van der Schaar. Learning "what-if" explanations for sequential decision-making. In *Proceedings of the International Conference on Learning Representations*, 2021.

[27] Simone Borsci, Alessio Malizia, Martin Schmettow, Frank Van Der Velde, Gunay Tariverdiyeva, Divyaa Balaji, and Alan Chamberlain. The chatbot usability scale: the design and pilot of a usability scale for interaction with ai-based conversational agents. *Personal and Ubiquitous Computing*, 26:95–119, 2022.

[28] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

[29] Juergen Branke, Salvatore Corrente, Salvatore Greco, and Walter Gutjahr. Efficient pairwise preference elicitation allowing for indifference. *Computers & Operations Research*, 88:175–186, 2017.

[30] Leo Breiman, Jerome Friedman, Richard A Olshen, and Charles J Stone. *Classification and regression trees*. Routledge, 2017.

[31] Marc Brittain and Peng Wei. Autonomous air traffic controller: A deep multi-agent reinforcement learning approach. *arXiv preprint arXiv:1905.01303*, 2019.

[32] Erik Brockbank. *Building Mental Models of Others Over Repeated Interactions*. University of California, San Diego, 2023.

[33] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

[34] David A Broniatowski et al. Psychological foundations of explainability and interpretability in artificial intelligence. *NIST, Tech. Rep*, 2021.

[35] Daniel S Brown, Yuchen Cui, and Scott Niekum. Risk-aware active inverse reinforcement learning. In *Conference on Robot Learning*, pages 362–372. PMLR, 2018.

[36] Yuri Burda, Harri Edwards, Deepak Pathak, Amos Storkey, Trevor Darrell, and Alexei A Efros. Large-scale study of curiosity-driven learning. *arXiv preprint arXiv:1808.04355*, 2018.

[37] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.

[38] Micah Carroll, Orr Paradise, Jessy Lin, Raluca Georgescu, Mingfei Sun, David Bignell, Stephanie Milani, Katja Hofmann, Matthew Hausknecht, Anca Dragan, and Sam Devlin. Uni [mask]: Unified inference in sequential decision problems. *Advances in Neural Information Processing Systems*, 35:35365–35378, 2022.

[39] Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. On the utility of learning about humans for human-ai coordination. *Advances in neural information processing systems*, 32, 2019.

[40] Ben Carterette. System effectiveness, user models, and user utility: a conceptual framework for investigation. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in information retrieval*, pages 903–912, 2011.

[41] Mike Casey. The dynamics of discrete-time computation, with application to recurrent neural networks and finite state machine extraction. *Neural computation*, 8(6):1135–1178, 1996.

[42] Center for Security and Emerging Technology. Autonomous cyber defense: A framework for understanding current and future capabilities. Technical report, Georgetown University, 2024.

[43] Urszula Chajewska, Daphne Koller, and Ronald Parr. Making rational decisions using adaptive utility elicitation. In *Aaai/Iaai*, pages 363–369, 2000.

[44] Kathryn Chaloner and Isabella Verdinelli. Bayesian experimental design: A review. *Statistical Science*, pages 273–304, 1995.

[45] Kushal Chauhan, Rishabh Tiwari, Jan Freyberg, Pradeep Shenoy, and Krishnamurthy Dvijotham. Interactive concept bottleneck models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 5948–5955, 2023.

[46] Yevgen Chebotar, Karol Hausman, Yao Lu, Ted Xiao, Dmitry Kalashnikov, Jake Varley, Alex Irpan, Benjamin Eysenbach, Ryan Julian, Chelsea Finn, et al. Actionable models: Unsupervised offline reinforcement learning of robotic skills. In *Proceedings of the International Conference on Machine Learning*, 2021.

[47] Ramnath K Chellappa and Raymond G Sin. Personalization versus privacy: An empirical examination of the online consumer's dilemma. *Information technology and management*, 6:181–202, 2005.

[48] Ja-Shen Chen, Tran-Thien-Y Le, and Devina Florence. Usability and responsiveness of artificial intelligence chatbot on online customer experience in e-retailing. *International Journal of Retail & Distribution Management*, 49(11):1512–1531, 2021.

[49] Jianyu Chen, Shengbo Eben Li, and Masayoshi Tomizuka. Interpretable end-to-end urban autonomous driving with latent deep reinforcement learning. *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[50] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.

[51] Rex Chen, Stephanie Milani, Zhicheng Zhang, Norman Sadeh, and Fei Fang. Making teams and influencing agents: Efficiently coordinating decision trees for interpretable multi-agent reinforcement learning. *arXiv preprint arXiv:2505.19316*, 2025.

[52] Valerie Chen, Nari Johnson, Nicholay Topin, Gregory Plumb, and Ameet Talwalkar. Use-case-grounded simulations for explanation evaluation. *arXiv preprint arXiv:2206.02256*, 2022.

[53] Xi Chen, Paul N Bennett, Kevyn Collins-Thompson, and Eric Horvitz. Pairwise ranking aggregation in a crowdsourced setting. In *Proceedings of the ACM International Conference on Web Search and Data Mining*, pages 193–202, 2013.

[54] Ziheng Chen, Fabrizio Silvestri, Gabriele Tolomei, He Zhu, Jia Wang, and Hongshik Ahn. Relace: Reinforcement learning agent for counterfactual explanations of arbitrary predictive models. *arXiv preprint arXiv:2110.11960*, 2021.

[55] Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.

[56] Andrew Cho, Jason M Woo, Brian Shi, Aishwaryaa Udeshi, and Jonathan SH Woo. The application of matec (multi-ai agent team care) framework in sepsis care. *arXiv preprint arXiv:2503.16433*, 2025.

[57] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30, 2017.

[58] Seong Yeub Chu, Jong Woo Kim, and Mun Yong Yi. Think together and work better: Combining humans' and llms' think-aloud outcomes for effective text evaluation. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–23, 2025.

[59] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pages 208–214, 2011.

[60] Jack Clark and Dario Amodei. Faulty reward functions in the wild, 2016.

[61] Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, pages 2048–2056. PMLR, 2020.

[62] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, 1960.

[63] Katherine Maeve Collins, Matthew Barker, Mateo Espinosa Zarlenga, Naveen Raman, Umang Bhatt, Mateja Jamnik, Ilia Sucholutsky, Adrian Weller, and Krishnamurthy Dvijotham. Human uncertainty in concept-based ai systems. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 869–889, 2023.

[64] Wolfram Conen and Tuomas Sandholm. Preference elicitation in combinatorial auctions. In *Proceedings of the ACM Conference on Electronic Commerce*, pages 256–259, 2001.

[65] Daphne Cornelisse, Aarav Pandya, Kevin Joseph, Joseph Suárez, and Eugene Vinitsky. Building reliable sim driving agents by scaling self-play. *arXiv preprint arXiv:2502.14706*, 2025.

[66] Leonardo Lucio Custode and Giovanni Iacca. Interpretable ai for policy-making in pandemics. *arXiv preprint arXiv:2204.04256*, 2022.

[67] Yinglong Dai, Haibin Ouyang, Hong Zheng, Han Long, and Xiaojun Duan. Interpreting a deep reinforcement learning model with conceptual embedding and performance analysis. *Applied Intelligence*, pages 1–17, 2022.

[68] Yuxin Dai, Qimei Chen, Jun Zhang, Xiaohui Wang, Yilin Chen, Tianlu Gao, Peidong Xu, Siyuan Chen, Siyang Liao, Huaiguang Jiang, et al. Enhanced oblique decision tree enabled policy extraction for deep reinforcement learning in power system emergency control. *Electric Power Systems Research*, 209:107932, 2022.

[69] Ali Darejeh and Dalbir Singh. A review on user interface design principles to increase software usability for users with less computer literacy. *Journal of Computer Science*, 9(11):1443, 2013.

[70] Sylvain Daronnat, Leif Azzopardi, Martin Halvey, and Mateusz Dubiel. Inferring trust from users' behaviours; agents' predictability positively affects trust, task performance and cognitive load in human-agent real-time collaboration. *Frontiers in Robotics and AI*, 8:642201, 2021.

[71] Devleena Das, Sonia Chernova, and Been Kim. State2explanation: Concept-based explanations to benefit agent learning and user understanding. *Advances in Neural Information Processing Systems*, 36, 2023.

[72] Thomas Degris, Olivier Sigaud, and Pierre-Henri Wuillemin. Learning the structure of factored markov decision processes in reinforcement learning problems. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 257–264, 2006.

[73] Frederic Dehais, Vsevolod Peysakhovich, Sébastien Scannella, Jennifer Fongue, and Thibault Gateau. "automation surprise" in aviation: real-time solutions. In *Proceedings of the 33rd annual ACM conference on Human Factors in Computing Systems*, pages 2525–2534, 2015.

[74] Quentin Delfosse, Jannis Blüml, Bjarne Gregori, Sebastian Sztwiertnia, and Kristian Kersting. Ocatari: Object-centric atari 2600 reinforcement learning environments. In *Advances in Neural Information Processing Systems*, 2023.

[75] Quentin Delfosse, Hikaru Shindo, Devendra Dhami, and Kristian Kersting. Interpretable and explainable logical policies via neurally guided symbolic abstraction. *Advances in Neural Information Processing Systems*, 36, 2024.

[76] Zehang Deng, Yongjian Guo, Changzhou Han, Wanlun Ma, Junwu Xiong, Sheng Wen, and Yang Xiang. AI agents under threat: A survey of key security challenges and future pathways. *ACM Computing Surveys*, 57(7):1–36, 2025.

[77] Sam Devlin, Raluca Georgescu, Ida Momennejad, Jaroslaw Rzepecki, Evelyn Zuniga, Gavin Costello, Guy Leroy, Ali Shaw, and Katja Hofmann. Navigation turing test (ntt): Learning to evaluate human-like navigation. In *Proceedings of the International Conference on Machine Learning*, pages 2644–2653. PMLR, 2021.

[78] Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. Openai baselines. `https://github.com/openai/baselines`, 2017.

[79] Yashesh Dhebar, Kalyanmoy Deb, Subramanya Nageshrao, Ling Zhu, and Dimitar Filev. Interpretable-ai policies using evolutionary nonlinear decision trees for discrete action systems. *arXiv preprint arXiv:2009.09521*, 2020.

[80] Marius-Constantin Dinu, Markus Hofmarcher, Vihang P Patil, Matthias Dorfer, Patrick M Blies, Johannes Brandstetter, Jose A Arjona-Medina, and Sepp Hochreiter. Xai and strategy extraction via reward redistribution. In *International Workshop on Extending Explainable AI Beyond Deep Models and Classifiers*, pages 177–205. Springer, 2022.

[81] Shayan Doroudi, Vincent Aleven, and Emma Brunskill. Where's the reward? a review of reinforcement learning for instructional sequencing. *International Journal of Artificial Intelligence in Education*, 29(4):568–620, 2019.

[82] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.

[83] Joanna Drummond and Craig Boutilier. Preference elicitation and interview minimization in stable matchings. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014.

[84] Yinuo Du, Zimeng Song, Stephanie Milani, Cleotilde Gonzales, and Fei Fang. Learning to play an adaptive cyber deception game. In *The 13th Workshop on Optimization and Learning in Multiagent Systems, AAMAS*, 2022.

[85] Yan Duan, Xi Chen, Rein Houthooft, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *Proceedings of the International Conference on Machine Learning*, pages 1329–1338. PMLR, 2016.

[86] Hippolyte Dubois, Patrick Le Callet, and Antoine Coutrot. Visualizing navigation difficulties in video game experiences. In *2021 13th International Conference on Quality of Multimedia Experience (QoMEX)*, pages 77–80. IEEE, 2021.

[87] Matthias Ehrgott. *Multicriteria optimization*, volume 491. Springer Science & Business Media, 2005.

[88] Arpad E Elo. *The rating of chessplayers, past and present*. Arco Pub., 1978.

[89] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. Implementation matters in deep RL: A case study on ppo and trpo. In *International Conference on Learning Representations*, 2019.

[90] Brochu Eric, Nando Freitas, and Abhijeet Ghosh. Active preference learning with discrete choice data. *Advances in Neural Information Processing Systems*, 20, 2007.

[91] Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6(Apr):503–556, 2005.

[92] Mateo Espinosa Zarlenga, Pietro Barbiero, Gabriele Ciravegna, Giuseppe Marra, Francesco Giannini, Michelangelo Diligenti, Zohreh Shams, Frederic Precioso, Stefano Melacci, Adrian Weller, et al. Concept embedding models: Beyond the accuracy-explainability trade-off. *Advances in Neural Information Processing Systems*, 35:21400–21413, 2022.

[93] European Union. General Data Protection Regulation (GDPR). `https://eur-lex.europa.eu/eli/reg/2016/679/oj`, 2016. Regulation (EU) 2016/679.

[94] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362, 2022.

[95] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[96] Adam Fourney, Gagan Bansal, Hussein Mozannar, Cheng Tan, Eduardo Salinas, Friederike Niedtner, Grace Proebsting, Griffin Bassman, Jack Gerrits, Jacob Alber, et al. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*, 2024.

[97] Nicholas Frosst and Geoffrey Hinton. Distilling a neural network into a soft decision tree. *arXiv preprint arXiv:1711.09784*, 2017.

[98] Nobuto Fujii, Yuichi Sato, Hironori Wakama, Koji Kazai, and Haruhiro Katayose. Evaluating human-like behaviors of video-game agents autonomously acquired with biological constraints. In *International Conference on Advances in Computer Entertainment Technology*, pages 61–76. Springer, 2013.

[99] Sandra Garcia-Rivadulla. Personalization vs. privacy: An inevitable trade-off? *IFLA journal*, 42(3):227–238, 2016.

[100] Josh Gardner, Christopher Brooks, and Ryan Baker. Evaluating the fairness of predictive student models through slicing analysis. In *Proceedings of the 9th international conference on learning analytics & knowledge*, pages 225–234, 2019.

[101] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 64(12):86–92, 2021.

[102] Claire Glanois, Paul Weng, Matthieu Zimmer, Dong Li, Tianpei Yang, Jianye Hao, and Wulong Liu. A survey on interpretable reinforcement learning. *Machine Learning*, pages 1–44, 2024.

[103] Adam Gleave, Mohammad Taufeeque, Juan Rocamonde, Erik Jenner, Steven H Wang, Sam Toyer, Maximilian Ernestus, Nora Belrose, Scott Emmons, and Stuart Russell. imitation: Clean imitation learning implementations. *arXiv preprint arXiv:2211.11972*, 2022.

[104] Astrid Glende. Agent design to pass computer games. In *Proceedings of the 42nd Annual Southeast Regional Conference*, pages 414–415, 2004.

[105] Suat Gönül, Tuncay Namlı, Ahmet Coşar, and İsmail Hakkı Toroslu. A reinforcement learning based algorithm for personalization of digital, just-in-time, adaptive interventions. *Artificial Intelligence in Medicine*, 115:102062, 2021.

[106] Omer Gottesman, Joseph Futoma, Yao Liu, Sonali Parbhoo, Leo Anthony Celi, Emma Brunskill, and Finale Doshi-Velez. Interpretable off-policy evaluation in reinforcement learning by highlighting influential transitions. *arXiv preprint, arXiv:2002.03478*, 2020.

[107] Jonathan Gray, Kavya Srinet, Yacine Jernite, Haonan Yu, Zhuoyuan Chen, Demi Guo, Siddharth Goyal, C Lawrence Zitnick, and Arthur Szlam. Craftassist: A framework for dialogue-enabled interactive agents. *arXiv preprint arXiv:1907.08584*, 2019.

[108] Djordje Grbic, Rasmus Berg Palm, Elias Najarro, Claire Glanois, and Sebastian Risi. Evocraft: A new challenge for open-endedness. In *Applications of Evolutionary Computation: 24th International Conference, EvoApplications 2021, Held as Part of EvoStar 2021, Virtual Event, April 7–9, 2021, Proceedings 24*, pages 325–340. Springer, 2021.

[109] Shane Griffith, Kaushik Subramanian, Jonathan Scholz, Charles L Isbell, and Andrea L Thomaz. Policy shaping: Integrating human feedback with reinforcement learning. *Advances in Neural Information Processing Systems*, 26, 2013.

[110] Niko Grupen, Natasha Jaques, Been Kim, and Shayegan Omidshafiei. Concept-based understanding of emergent multi-agent behavior. In *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.

[111] Stefan Freyr Gudmundsson, Philipp Eisen, Erik Poromaa, Alex Nodet, Sami Purmonen, Bartlomiej Kozakowski, Richard Meurling, and Lele Cao. Human-like playtesting with deep learning. In *2018 IEEE Conference on Computational Intelligence and Games (CIG)*, pages 1–8. IEEE, 2018.

[112] David Gunning and David W Aha. Darpa's explainable artificial intelligence program. *AI Magazine*, 40(2):44–58, 2019.

[113] Wei Guo and Peng Wei. Explainable deep reinforcement learning for aircraft separation assurance. *4th Digital Avionics Systems Conference*, 2022.

[114] Ujjwal Das Gupta, Erik Talvitie, and Michael Bowling. Policy tree: Adaptive representation for policy gradient. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2547–2553, 2015.

[115] William H Guss, Mario Ynocente Castro*, Sam Devlin*, Brandon Houghton*, Noboru Sean Kuno*, Crissman Loomis*, Stephanie Milani*, Sharada Mohanty*, Keisuke Nakata*, Ruslan Salakhutdinov*, John Schulman*, et al. The MineRL 2020 competition on sample efficient reinforcement learning using human priors. *The 34th Conference on Neural Information Processing Systems (NeurIPS) Competition Track*, 2020.

[116] William H. Guss, Cayden Codel*, Katja Hofmann*, Brandon Houghton*, Noboru Kuno*, Stephanie Milani*, Sharada Mohanty*, Diego Perez Liebana*, Ruslan Salakhutdinov*, Nicholay Topin*, Manuela Veloso*, and Phillip Wang*. The MineRL competition on sample efficient reinforcement learning using human priors. In *The 23rd Conference on Neural Information Processing Systems Competition Track*, 2019.

[117] William Hebgen Guss, Stephanie Milani, Nicholay Topin, Brandon Houghton, Sharada Mohanty, Andrew Melnik, Augustin Harter, Benoit Buschmaas, Bjarne Jaster, Christoph Berganski, et al. Towards robust and domain agnostic reinforcement learning competitions: Minerl 2020. In *NeurIPS 2020 Competition and Demonstration Track*, pages 233–252. PMLR, 2021.

[118] Rotem D Guttman, Jessica Hammer, Erik Harpstead, and Carol J Smith. Play for real (ism)-using games to predict human-ai interactions in the real world. *Proceedings of the ACM on Human-Computer Interaction*, 5:1–17, 2021.

[119] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the International Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.

[120] Dylan Hadfield-Menell, Anca D Dragan, Pieter Abbeel, and Stuart Russell. The off-switch game. In *AAAI Workshops*, 2017.

[121] Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. Cooperative inverse reinforcement learning. *Advances in Neural Information Processing Systems*, 29, 2016.

[122] Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*, 2021.

[123] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.

[124] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.

[125] Torkel Hafting, Marianne Fyhn, Sturla Molden, May-Britt Moser, and Edvard I Moser. Microstructure of a spatial map in the entorhinal cortex. *Nature*, 436(7052):801–806, 2005.

[126] Stephane Hatgis-Kessell, W Bradley Knox, Serena Booth, Scott Niekum, and Peter Stone. Influencing humans to conform to preference models for rlhf. *arXiv preprint arXiv:2501.06416*, 2025.

[127] Simon Hecker, Dengxin Dai, Alexander Liniger, Martin Hahner, and Luc Van Gool. Learning accurate and human-like driving using semantic maps and attention. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2346–2353, Piscataway, New Jersey, 2020. IEEE.

[128] Daniel Hein, Alexander Hentschel, Thomas Runkler, and Steffen Udluft. Particle swarm optimization for generating interpretable fuzzy reinforcement learning policies. *Engineering Applications of Artificial Intelligence*, 65:87 – 98, 2017.

[129] Daniel Hein, Steffen Udluft, and Thomas A Runkler. Interpretable policies for reinforcement learning by genetic programming. *Engineering Applications of Artificial Intelligence*, 76:158–169, 2018.

[130] Ralf Herbrich, Tom Minka, and Thore Graepel. Trueskill™: a bayesian skill rating system. *Advances in Neural Information Processing Systems*, 19, 2006.

[131] Alexandre Heuillet, Fabien Couthouis, and Natalia Díaz-Rodríguez. Collective explainable ai: Explaining cooperative strategies and agent contribution in multiagent reinforcement learning with shapley values. *IEEE Computational Intelligence Magazine*, 17(1):59–71, 2022.

[132] Philip Hingston. A new design for a turing test for bots. In *Proceedings of the 2010 IEEE Conference on Computational Intelligence and Games*, pages 345–350, Piscataway, New Jersey, 2010. IEEE.

[133] G Hinton, O Vinyals, and J Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[134] Vita Hinze-Hoare. The review and analysis of human computer interaction (hci) principles. *arXiv preprint arXiv:0707.3638*, 2007.

[135] Sean D Holcomb, William K Porter, Shaun V Ault, Guifen Mao, and Jin Wang. Overview on deepmind and its alphago zero ai. In *Proceedings of the 2018 International Conference on Big Data and Education*, pages 67–71, 2018.

[136] Ryan Hoque, Ashwin Balakrishna, Ellen Novoseller, Albert Wilcox, Daniel S Brown, and Ken Goldberg. Thriftydagger: Budget-aware novelty and risk gating for interactive imitation learning. *arXiv preprint arXiv:2109.08273*, 2021.

[137] Brandon Houghton, Stephanie Milani, Nicholay Topin, William Guss, Katja Hofmann, Diego Perez-Liebana, Manuela Veloso, and Ruslan Salakhutdinov. Guaranteeing reproducibility in deep learning competitions. *NeurIPS Challenges in Machine Learning workshop (CiML)*, 2019.

[138] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1314–1324, 2019.

[139] Sandy H Huang, Kush Bhatia, Pieter Abbeel, and Anca D Dragan. Establishing appropriate trust via critical states. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3929–3936. IEEE, 2018.

[140] Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.

[141] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward learning from human preferences and demonstrations in atari. *Advances in Neural Information Processing Systems*, 31, 2018.

[142] Panagiotis G Ipeirotis, Foster Provost, and Jing Wang. Quality management on amazon mechanical turk. In *Proceedings of the ACM SIGKDD Workshop on Human Computation*, pages 64–67, New York City, New York, 2010. ACM.

[143] Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, pages 2961–2970. PMLR, 2019.

[144] Ozan Irsoy, Olcay Taner Yıldız, and Ethem Alpaydın. Soft decision trees. In *Proceedings of the International Conference on Pattern Recognition*, pages 1819–1822. IEEE, 2012.

[145] Athul Paul Jacob, David J Wu, Gabriele Farina, Adam Lerer, Hengyuan Hu, Anton Bakhtin, Jacob Andreas, and Noam Brown. Modeling strong and human-like gameplay with kl-regularized search. In *Proceedings of the International Conference on Machine Learning*, pages 9695–9728. PMLR, 2022.

[146] Mikhail Jacob, Sam Devlin, and Katja Hofmann. "it's unwieldy and it takes a lot of time"—challenges and opportunities for creating agents in commercial games. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 16, pages 88–94, Palo Alto, CA, 2020. AAAI Press.

[147] Hong Jun Jeon, Smitha Milli, and Anca Dragan. Reward-rational (implicit) choice: A unifying formalism for reward learning. *Advances in Neural Information Processing Systems*, 33:4415–4426, 2020.

[148] Aman Jhunjhunwala. Policy extraction via online q-value distillation. *Masters Thesis, University of Waterloo*, 2019.

[149] Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. The malmo platform for artificial intelligence experimentation. In *Proceedings of the International Joint Conference on Artificial Intelligence*, volume 16, pages 4246–4247, 2016.

[150] Karolis Jucys, George Adamopoulos, Mehrab Hamidi, Stephanie Milani, Mohammad Reza Samsami, Artem Zholus, Sonia Joseph, Blake Richards, Irina Rish, and Özgür Şimşek. Interpretability in action: Exploratory analysis of vpt, a minecraft agent. *arXiv preprint arXiv:2407.12161*, 2024.

[151] Arthur Juliani, Vincent-Pierre Berges, Ervin Teng, Andrew Cohen, Jonathan Harper, Chris Elion, Chris Goy, Yuan Gao, Hunter Henry, Marwan Mattar, et al. Unity: A general platform for intelligent agents. *arXiv preprint arXiv:1809.02627*, 2018.

[152] Zoe Juozapaitis, Anurag Koul, Alan Fern, Martin Erwig, and Finale Doshi-Velez. Explainable reinforcement learning via reward decomposition. In *Proceedings of the International Joint Conference on Artificial Intelligence Workshop on Explainable Artificial Intelligence*, 2019.

[153] Daniel Kahneman, Olivier Sibony, and Cass R Sunstein. *Noise: A flaw in human judgment*. Hachette UK, 2021.

[154] Anssi Kanervisto, Stephanie Milani, Karolis Ramanauskas, Byron Galbraith, Steven H. Wang, Brandon Houghton, Sharada Mohanty, and Rohin Shah. The MineRL BASALT competition on learning from human feedback. *The 36th Conference on Neural Information Processing Systems (NeurIPS) Competition Track*, 2022.

[155] Anssi Kanervisto*, Stephanie Milani*, Karolis Ramanauskas, Nicholay Topin, Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, Wei Yang, et al. Minerl diamond 2021 competition: Overview, results, and lessons learned. *NeurIPS Competitions and Demonstrations Track*, pages 13–28, 2022.

[156] Anssi Kanervisto, Christian Scheller, and Ville Hautamäki. Action space shaping in deep reinforcement learning. In *2020 IEEE Conference on Games (CoG)*, pages 479–486. IEEE, 2020.

[157] Igor V Karpov, Jacob Schrum, and Risto Miikkulainen. Believable bot navigation via playback of human traces. In *Believable bots*, pages 151–170. Springer, 2013.

[158] Seth Karten, Jake Grigsby, Stephanie Milani, Kiran Vodrahalli, Amy Zhang, Fei Fang, Yuke Zhu, and Chi Jin. The pokéagent challenge: Competitive and long-context learning at scale. *NeurIPS Competition Track*, 2025.

[159] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. Interpreting interpretability: understanding data scientists' use of interpretability tools for machine learning. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2020.

[160] Dmitry Kazhdan, Zohreh Shams, and Pietro Lio. Marleme: A multi-agent reinforcement learning model extraction library. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020.

[161] Eoin M Kenny, Mycal Tucker, and Julie Shah. Towards interpretable deep reinforcement learning with human-friendly prototypes. In *Proceedings of the International Conference on Learning Representations*, 2022.

[162] Steven Kerr. On the folly of rewarding a, while hoping for b. *Academy of Management journal*, 18(4):769–783, 1975.

[163] Man-Je Kim, Kyung-Joong Kim, Seungjun Kim, and Anind K. Dey. Performance evaluation gaps in a real-time strategy game between human and artificial intelligence players. *IEEE Access*, 6:13575–13586, 2018.

[164] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[165] Julia Kiseleva, Ziming Li, Mohammad Aliannejadi, Shrestha Mohanty, Maartje ter Hoeve, Mikhail Burtsev, Alexey Skrynnik, Artem Zholus, Aleksandr Panov, Kavya Srinet, et al. Neurips 2021 competition iglu: Interactive grounded language understanding in a collaborative environment. *arXiv preprint arXiv:2110.06536*, 2021.

[166] Julia Kiseleva, Alexey Skrynnik, Artem Zholus, Shrestha Mohanty, Negar Arabzadeh, Marc-Alexandre Côté, Mohammad Aliannejadi, Milagro Teruel, Ziming Li, Mikhail Burtsev, et al. Iglu 2022: Interactive grounded language understanding in a collaborative environment at neurips 2022. *arXiv preprint arXiv:2205.13771*, 2022.

[167] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *Proceedings of the International Conference on Machine Learning*, pages 5338–5348. PMLR, 2020.

[168] Anurag Koul, Sam Greydanus, and Alan Fern. Learning finite state representations of recurrent policy networks. *arXiv preprint, arXiv:1811.12530*, 2018.

[169] Victoria Krakovna. Specification gaming examples in AI, 2018.

[170] Victoria Krakovna, Jonathan Uesato, Vladimir Mikulik, Matthew Rahtz, Tom Everitt, Ramana Kumar, Zac Kenton, Jan Leike, and Shane Legg. Specification gaming: the flip side of ai ingenuity. *DeepMind Blog*, 2020.

[171] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement learning: A survey. *Information Fusion*, 85:1–22, 2022.

[172] Isaac Lage and Finale Doshi-Velez. Learning interpretable concept-based models with human feedback. *International Conference on Machine Learning: Workshop on Human Interpretability in Machine Learning*, 2020.

[173] Isaac Lage, Daphna Lifschitz, Finale Doshi-Velez, and Ofra Amir. Exploring computational user models for agent policy summarization. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, 2019.

[174] Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.

[175] Cassidy Laidlaw and Anca Dragan. The boltzmann policy distribution: Accounting for systematic suboptimality in human models. In *Proceedings of the International Conference on Learning Representations*, 2022.

[176] Guillaume Lample and Devendra Singh Chaplot. Playing fps games with deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017.

[177] Mikel Landajuela, Brenden K Petersen, Sookyung Kim, Claudio P Santiago, Ruben Glatt, Nathan Mundhenk, Jacob F Pettit, and Daniel Faissol. Discovering symbolic policies with deep reinforcement learning. In *International Conference on Machine Learning*, pages 5979–5989. PMLR, 2021.

[178] Kuang-Huei Lee, Ofir Nachum, Mengjiao Yang, Lisa Lee, Daniel Freeman, Winnie Xu, Sergio Guadarrama, Ian Fischer, Eric Jang, Henryk Michalewski, et al. Multi-game decision transformers. *arXiv preprint arXiv:2205.15241*, 2022.

[179] Joel Lehman, Jeff Clune, and Dusan Misevic. The surprising creativity of digital evolution. In *Artificial Life Conference Proceedings*, pages 55–56. MIT Press, 2018.

[180] Belinda Z Li, Alex Tamkin, Noah Goodman, and Jacob Andreas. Eliciting human preferences with language models. *arXiv preprint arXiv:2310.11589*, 2023.

[181] Cheng Li, Na Zhao, and Hao Wu. Multiple deception resources deployment strategy based on reinforcement learning for network threat mitigation. *Scientific Reports*, 15:16830, 2025.

[182] Huao Li, Stephanie Milani, Vigneshram Krishnamoorthy, Michael Lewis, and Katia Sycara. Perceptions of domestic robots' normative behavior across cultures. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 345–351, 2019.

[183] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the International Conference on World Wide Web*, pages 661–670, 2010.

[184] Shihui Li, Yi Wu, Xinyue Cui, Honghua Dong, Fei Fang, and Stuart Russell. Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4213–4220, 2019.

[185] Wenhao Li, Bo Jin, and Xiangfeng Wang. Sparsemaac: Sparse attention for multi-agent reinforcement learning. In *Database Systems for Advanced Applications: DASFAA 2019 International Workshops: BDMS, BDQM, and GDMA, Chiang Mai, Thailand, April 22–25, 2019, Proceedings 24*, pages 96–110. Springer, 2019.

[186] Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba, and Sheila McIlraith. Steve-1: A generative model for text-to-behavior in minecraft. *Advances in Neural Information Processing Systems*, 36:69900–69929, 2023.

[187] Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.

[188] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine Learning*. Elsevier, 1994.

[189] Xiao-Yang Liu, Ziyi Xia, Jingyang Rui, Jiechao Gao, Hongyang Yang, Ming Zhu, Christina Wang, Zhaoran Wang, and Jian Guo. Finrl-meta: Market environments and benchmarks for data-driven financial reinforcement learning. *Advances in Neural Information Processing Systems*, 35:1835–1849, 2022.

[190] Daniel Livingstone. Turing's test and believable ai in games. *Computers in Entertainment (CIE)*, 4(1):6–es, 2006.

[191] Benedikt Loepp, Tim Hussein, and Jüergen Ziegler. Choice-based preference elicitation for collaborative filtering recommender systems. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 3085–3094, 2014.

[192] Steven Loria et al. textblob documentation. *Release 0.15*, 2(8):269, 2018.

[193] Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in Neural Information Processing Systems*, 30, 2017.

[194] Tyler Lu and Craig Boutilier. Robust approximation and incremental elicitation in voting protocols. In *Proceedings of the International Joint Conference on Artificial Intelligence*, volume 1, pages 287–293, 2011.

[195] Tyler Lu and Craig Boutilier. Vote elicitation with probabilistic preference models: Empirical estimation and cost tradeoffs. In *Algorithmic Decision Theory: Second International Conference, ADT 2011, Piscataway, NJ, USA, October 26-28, 2011. Proceedings 2*, pages 135–149. Springer, 2011.

[196] R Duncan Luce. *Individual choice behavior: A theoretical analysis*. Courier Corporation, 2012.

[197] Ronny Luss, Amit Dhurandhar, and Miao Liu. Local explanations for reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023.

[198] Zhihao Ma, Yuzheng Zhuang, Paul Weng, Hankz Hankui Zhuo, Dong Li, Wulong Liu, and Jianye Hao. Learning symbolic rules for interpretable deep reinforcement learning. *arXiv preprint arXiv:2103.08228*, 2021.

[199] Brian Mac Namee. Proactive persistent agents-using situational intelligence to create support characters in character-centric computer games. 2004.

[200] Prashan Madumal, Tim Miller, Liz Sonenberg, and Frank Vetere. Explainable reinforcement learning through a causal lens. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.

[201] Federico Malato, Florian Leopold, Amogh Raut, Ville Hautamäki, and Andrew Melnik. Behavioral cloning via search in video pretraining latent space. *arXiv preprint arXiv:2212.13326*, 2022.

[202] Kleanthis Malialis and Daniel Kudenko. Distributed response to network intrusions using multiagent reinforcement learning. *Engineering Applications of Artificial Intelligence*, 41:270–284, 2015.

[203] Jiayuan Mao, Tomás Lozano-Pérez, Josh Tenenbaum, and Leslie Kaelbling. Pdsketch: Integrated domain programming, learning, and planning. *Advances in Neural Information Processing Systems*, 35:36972–36984, 2022.

[204] Catherine C Marshall and Frank M Shipman. Experiences surveying the crowd: Reflections on methods, participation, and reliability. In *Proceedings of the Annual ACM Web Science Conference*, pages 234–243, 2013.

[205] Laetitia Matignon, Guillaume J Laurent, and Nadine Le Fort-Piat. Independent reinforcement learners in cooperative markov games: a survey regarding coordination problems. *The Knowledge Engineering Review*, 27(1):1–31, 2012.

[206] Andrew Kachites McCallum. *Reinforcement learning with selective perception and hidden state*. University of Rochester, 1996.

[207] R Andrew McCallum. Reinforcement learning with selective perception and hidden state. *PhD Thesis, University of Rochester, Department of Computer Science*, 1997.

[208] Matheus RF Mendonça, Heder S Bernardino, and Raul F Neto. Simulating human behavior in fighting games using reinforcement learning and artificial neural networks. In *2015 14th Brazilian symposium on computer games and digital entertainment (SBGames)*, pages 152–159. IEEE, 2015.

[209] Zili Meng, Minhu Wang, Jiasong Bai, Mingwei Xu, Hongzi Mao, and Hongxin Hu. Interpreting deep learning-based networking systems. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication*, pages 154–171, 2020.

[210] Microsoft Corporation. Randomized controlled trial for copilot for security: Productivity findings. Technical report, Microsoft Corporation, January 2024.

[211] Stephanie Milani*, Zhou Fan*, Saurabh Gulati, Thanh Nguyen, Fei Fang, and Amulya Yadav. Intelligent tutoring strategies for students with autism spectrum disorder: A reinforcement learning approach. In *The 34th AAAI Conference on Artificial Intelligence (AAAI) Workshop on Artificial Intelligence for Education*, 2020.

[212] Stephanie Milani, Arthur Juliani, Ida Momennejad, Raluca Georgescu, Jaroslaw Rzepecki, Alison Shaw, Gavin Costello, Fei Fang, Sam Devlin, and Katja Hofmann. Navigates like me: Understanding how people evaluate human-like ai in video games. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–18, 2023.

[213] Stephanie Milani, Anssi Kanervisto, Karolis Ramanauskas, Sander Schulhoff, Brandon Houghton, Sharada Mohanty, Byron Galbraith, Ke Chen, Yan Song, Tianze Zhou, et al. Towards solving fuzzy tasks with human feedback: A retrospective of the minerl basalt 2022 competition. *arXiv preprint arXiv:2303.13512*, 2023.

[214] Stephanie Milani, Anssi Kanervisto, Karolis Ramanauskas, Sander Schulhoff, Brandon Houghton, and Rohin Shah. BEDD: The MineRL BASALT evaluation and demonstrations dataset for training and benchmarking agents that solve fuzzy tasks. *Advances in Neural Information Processing Systems*, 36:32867–32878, 2023.

[215] Stephanie Milani, Weiran Shen, Kevin S Chan, Sridhar Venkatesan, Nandi O Leslie, Charles Kamhoua, and Fei Fang. Harnessing the power of deception in attack graph-based security games. In *Proceedings of the International Conference on Decision and Game Theory for Security*, pages 147–167. Springer, 2020.

[216] Stephanie Milani, Nicholay Topin, Brandon Houghton, William H Guss, Sharada P Mohanty, Keisuke Nakata, Oriol Vinyals, and Noboru Sean Kuno. Retrospective analysis of the 2019 minerl competition on sample efficient reinforcement learning. In *NeurIPS 2019 Competition and Demonstration Track*, pages 203–214. PMLR, 2020.

[217] Stephanie Milani, Nicholay Topin, Zheyuan Ryan Shi, Charles Kamhoua, Evangelos E. Papalexakis, and Fei Fang. Extracting decision tree policies for interpretable multi-agent reinforcement learning. In *The 36th AAAI Conference on Artificial Intelligence (AAAI) Workshop on Explainable Agency in Artificial Intelligence*, 2022.

[218] Stephanie Milani, Nicholay Topin, and Katia Sycara. Penalty-modified markov decision processes: Efficient incorporation of norms into sequential decision making problems. In *Multidisciplinary Conference on Reinforcement Learning and Decision Making*, 2019.

[219] Stephanie Milani, Nicholay Topin, Manuela Veloso, and Fei Fang. Explainable reinforcement learning: A survey and comparative review. *ACM Computing Surveys (CSUR)*, 2023.

[220] Stephanie Milani*, Zhicheng Zhang*, Nicholay Topin, Ryan Shi, Charles Kamhoua, Evangelos E Papalexakis, and Fei Fang. Maviper: Learning decision tree policies for interpretable multi-agent reinforcement learning. In *Proceedings of the European Conference on Machine Learning (ECML)*, 2022.

[221] Stephanie Milani, Zhicheng Zhang, Nicholay Topin, Lirong Xia, and Fei Fang. Aligning agent policies with preferences: Human-centered interpretable reinforcement learning. *Under Review*, 2025.

[222] Silvia Milano, Mariarosaria Taddeo, and Luciano Floridi. Recommender systems and their ethical challenges. *AI & SOCIETY*, 35(4):957–967, 2020.

[223] Smitha Milli, Luca Belli, and Moritz Hardt. From optimizing engagement to measuring value. *arXiv preprint arXiv:2008.12623*, 2020.

[224] Maximiliano Miranda, Antonio A Sánchez-Ruiz, and Federico Peinado. A neuroevolution approach to imitating human-like play in ms. pac-man video game. In *CoSECivi*, pages 113–124, 2016.

[225] Azalia Mirhoseini, Anna Goldie, Mustafa Yazgan, Joe Wenjie Jiang, Ebrahim Songhori, Shen Wang, Young-Joon Lee, Eric Johnson, Omkar Pathak, Azade Nazi, et al. A graph placement methodology for fast chip design. *Nature*, 594(7862):207–212, 2021.

[226] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing Atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.

[227] Sharada Mohanty, Erik Nygren, Florian Laurent, Manuel Schneider, Christian Scheller, Nilabha Bhattacharya, Jeremy Watson, Adrian Egli, Christian Eichenberger, Christian Baumberger, et al. Flatland-rl: Multi-agent reinforcement learning on trains. *arXiv preprint arXiv:2012.05893*, 2020.

[228] Christopher Molnar. *Interpretable Machine Learning*. 2019.

[229] Michael Moritz, Eric Topol, and Pranav Rajpurkar. Coordinated ai agents for advancing healthcare. *Nature Biomedical Engineering*, pages 1–7, 2025.

[230] Yoshinari Motokawa and Toshiharu Sugawara. Mat-dqn: Toward interpretable multi-agent deep reinforcement learning for coordinated activities. In *International Conference on Artificial Neural Networks*, pages 556–567. Springer, 2021.

[231] Hussein Mozannar, Hunter Lang, Dennis Wei, Prasanna Sattigeri, Subhro Das, and David Sontag. Who should predict? exact algorithms for learning to defer to humans. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2023.

[232] Maxim Mozgovoy and Iskander Umarov. Behavior capture: Building believable and effective ai agents for video games. *International Journal of Arts & Sciences*, 4(20):243, 2011.

[233] Kara Murias, Kathy Kwok, Adrian Gil Castillejo, Irene Liu, and Giuseppe Iaria. The effects of video game use on performance in a virtual navigation task. *Computers in Human Behavior*, 58:398–406, 2016.

[234] Mohamed Musbah, Matthew Mitchell Dixon, Geoffrey Jacoby Gordon, Mahmoud Adada, Soroush Mehri, Andrew James McNamara, and Jonathan David Morrison. Providing automated user input to an application during a disruption, U.S. Patent 11213746, Oct. 2019.

[235] Anjali Narayan-Chen, Prashant Jayannavar, and Julia Hockenmaier. Collaborative dialogue in minecraft. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 5405–5415, 2019.

[236] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proceedings of the International Conference on Machine Learning*, volume 99, pages 278–287, 1999.

[237] Tuomas Oikarinen, Subhro Das, Lam M Nguyen, and Tsui-Wei Weng. Label-free concept bottleneck models. *arXiv preprint arXiv:2304.06129*, 2023.

[238] John O'Keefe and Lynn Nadel. *The hippocampus as a cognitive map*. Oxford: Clarendon Press, New York City, New York, 1978.

[239] Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. Optimal and approximate q-value functions for decentralized pomdps. *Journal of Artificial Intelligence Research*, 32:289–353, 2008.

[240] Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.

[241] Matthew L Olson, Roli Khanna, Lawrence Neal, Fuxin Li, and Weng-Keen Wong. Counterfactual state explanations for reinforcement learning agents via generative deep learning. *Artificial Intelligence*, 295:103455, 2021.

[242] Christian W Omlin and C Lee Giles. Constructing deterministic finite-state automata in recurrent neural networks. *Journal of the ACM (JACM)*, 43(6):937–972, 1996.

[243] Amy Orben and Andrew K Przybylski. The association between adolescent well-being and digital technology use. *Nature Human Behaviour*, 3(2):173–182, 2019.

[244] Laurent Orseau and M Armstrong. Safely interruptible agents. In *Conference on Uncertainty in Artificial Intelligence*. Association for Uncertainty in Artificial Intelligence, 2016.

[245] A Ross Otto, Sean Devine, Eric Schulz, Aaron M Bornstein, and Kenway Louie. Context-dependent choice and evaluation in real-world consumer behavior. *Scientific Reports*, 12(1):17744, 2022.

[246] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

[247] Sharon Oviatt. Ten myths of multimodal interaction. *Communications of the ACM*, 42(11):74–81, 1999.

[248] Sharon Oviatt, Björn Schuller, Philip Cohen, Daniel Sonntag, and Gerasimos Potamianos. *The handbook of multimodal-multisensor interfaces, volume 1: Foundations, user modeling, and common modality combinations*. Morgan & Claypool, 2017.

[249] Gabriele Paolacci, Jesse Chandler, and Panagiotis G Ipeirotis. Running experiments on amazon mechanical turk. *Judgment and Decision Making*, 5(5):411–419, 2010.

[250] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.

[251] A. Patterson, S. Neumann, M. White, and A. White. Empirical design in reinforcement learning. *arXiv preprint arXiv:2304.01315*, 2023.

[252] Radek Pelánek. Applications of the elo rating system in adaptive educational systems. *Computers & Education*, 98:169–179, 2016.

[253] Svetlin Penkov and Subramanian Ramamoorthy. Learning programmatically structured representations with perceptor gradients. In *Proceedings of the International Conference on Learning Representations*, 2019.

[254] Diego Perez-Liebana, Katja Hofmann, Sharada Prasanna Mohanty, Noboru Kuno, Andre Kramer, Sam Devlin, Raluca D. Gaina, and Daniel Ionita. The multi-agent reinforcement learning in Malmö (MARLÖ) competition. *arXiv preprint arXiv:1901.08129*, 2019.

[255] My Phan, Kianté Brantley*, Stephanie Milani*, Soroush Mehri*, Gokul Swamy*, and Geoffrey J Gordon. When is transfer learning possible? In *Proceedings of the International Conference on Machine Learning*, 2024.

[256] Eleonora Poeta, Gabriele Ciravegna, Eliana Pastor, Tania Cerquitelli, and Elena Baralis. Concept-based explainable artificial intelligence: A survey. *arXiv preprint arXiv:2312.12936*, 2023.

[257] LA Prashanth, Cheng Jie, Michael Fu, Steve Marcus, and Csaba Szepesvári. Cumulative prospect theory meets reinforcement learning: Prediction and control. In *Proceedings of the International Conference on Machine Learning*, pages 1406–1415. PMLR, 2016.

[258] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.

[259] Larry D Pyeatt. Reinforcement learning with decision trees. In *Applied Informatics*, pages 26–31, 2003.

[260] Larry D Pyeatt and Adele E Howe. Decision tree function approximation in reinforcement learning. In *Proceedings of the International Symposium on Adaptive Systems: Evolutionary Computation and Probabilistic Graphical Models*, volume 2, pages 70–77, 2001.

[261] JR Quinlan. Induction of decision trees. *Machine Learning*, 1(1), 1986.

[262] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.

[263] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021.

[264] Tom Rainforth, Adam Foster, Desi R Ivanova, and Freddie Bickford Smith. Modern bayesian experimental design. *Statistical Science*, 39(1):100–114, 2024.

[265] Inioluwa Deborah Raji, Emily M Bender, Amandalynne Paullada, Emily Denton, and Alex Hanna. Ai and the everything in the whole wide world benchmark. *arXiv preprint arXiv:2111.15366*, 2021.

[266] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.

[267] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.

[268] Anthony E Richardson, Morgan E Powers, and Lauren G Bousquet. Video game experience predicts virtual, but not real navigation performance. *Computers in Human Behavior*, 27(1):552–560, 2011.

[269] Ariel Rosenfeld, Moshe Cohen, Matthew E Taylor, and Sarit Kraus. Leveraging human knowledge in tabular reinforcement learning: A study of human subjects. *The Knowledge Engineering Review*, 33:1–26, 2018.

[270] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pages 627–635, 2011.

[271] Aaron M Roth, Nicholay Topin, Pooyan Jamshidi, and Manuela Veloso. Conservative q-improvement: Reinforcement learning for an interpretable decision-tree policy. *arXiv preprint, arXiv:1907.01180*, 2019.

[272] Cynthia Rudin. Algorithms for interpretable machine learning. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1519–1519, 2014.

[273] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019.

[274] Cynthia Rudin, Chaofan Chen, Zhi Chen, Haiyang Huang, Lesia Semenova, and Chudi Zhong. Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistics Surveys*, 16:1–85, 2022.

[275] Christian Rupprecht, Cyril Ibrahim, and Christopher J Pal. Finding and visualizing weaknesses of deep reinforcement learning agents. In *Proceedings of the International Conference on Learning Representations*, 2020.

[276] Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. pearson, 2016.

[277] Aadirupa Saha. Optimal algorithms for stochastic contextual preference bandits. *Advances in Neural Information Processing Systems*, 34:30050–30062, 2021.

[278] Christoph Salge, Michael Cerny Green, Rodgrigo Canaan, and Julian Togelius. Generative design in minecraft (gdmc) settlement generation competition. In *Proceedings of the International Conference on the Foundations of Digital Games*, pages 1–10, 2018.

[279] Atrisha Sarkar, Andrei Ioan Muresanu, Carter Blair, Aaryam Sharma, Rakshit S Trivedi, and Gillian K Hadfield. Normative modules: A generative agent architecture for learning norms that supports multi-agent cooperation. *arXiv preprint arXiv:2405.19328*, 2024.

[280] William Saunders, Girish Sastry, Andreas Stuhlmueller, and Owain Evans. Trial without error: Towards safe reinforcement learning via human intervention. *arXiv preprint arXiv:1707.05173*, 2017.

[281] Matthias Scheutz, Paul Schermerhorn, James Kramer, and David Anderson. First steps toward natural human-like hri. *Autonomous Robots*, 22(4):411–423, 2007.

[282] Douglas Schuler and Aki Namioka. *Participatory design: Principles and practices*. CRC Press, 1993.

[283] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *Proceedings of the International Conference on Machine Learning*, pages 1889–1897. PMLR, 2015.

[284] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[285] Peter Seiler, Bongsob Song, and J Karl Hedrick. Development of a collision avoidance system. *SAE transactions*, pages 1334–1340, 1998.

[286] Alessandro Sestini, Linus Gisslén, Joakim Bergdahl, Konrad Tollmar, and Andrew D Bagdanov. Automated gameplay testing and validation with curiosity-conditioned proximal trajectories. *IEEE Transactions on Games*, 2022.

[287] H Sebastian Seung, Manfred Opper, and Haim Sompolinsky. Query by committee. In *Proceedings of the Annual Workshop on Computational Learning Theory*, pages 287–294, 1992.

[288] Nihar B Shah, Sivaraman Balakrishnan, Joseph Bradley, Abhay Parekh, Kannan Ramch, and Martin J Wainwright. Estimation from pairwise comparisons: Sharp minimax bounds with topology dependence. *Journal of Machine Learning Research*, 17(58):1–47, 2016.

[289] Rohin Shah, Dmitrii Krasheninnikov, Jordan Alexander, Pieter Abbeel, and Anca Dragan. Preferences implicit in the state of the world. *arXiv preprint arXiv:1902.04198*, 2019.

[290] Rohin Shah, Steven H Wang, Cody Wild, Stephanie Milani, Anssi Kanervisto, Vinicius G Goecks, Nicholas Waytowich, David Watkins-Valls, Bharat Prakash, Edmund Mills, Divyansh Garg, Alexander Fries, Alexandra Souly, Chan Jun Shern, Daniel del Castillo, and Tom Lieberum. Retrospective on the 2021 basalt competition on learning from human feedback. In *Proceedings of the NeurIPS 2021 Competition and Demonstration Track*, 2022.

[291] Rohin Shah, Cody Wild, Steven H Wang, Neel Alex, Brandon Houghton, William Guss, Sharada Mohanty, Anssi Kanervisto, Stephanie Milani, Nicholay Topin, Pieter Abbeel, Stuart Russell, and Anca Dragan. The MineRL BASALT competition on learning from human feedback. *The 35th Conference on Neural Information Processing Systems (NeurIPS) Competition Track*, 2021.

[292] Noor Shaker, Julian Togelius, Georgios N Yannakakis, Likith Poovanna, Vinay S Ethiraj, Stefan J Johansson, Robert G Reynolds, Leonard K Heether, Tom Schumann, and Marcus Gallagher. The turing test track of the 2012 mario ai championship: entries and evaluation. In *2013 IEEE Conference on Computational Inteligence in Games (CIG)*, pages 1–8. IEEE, 2013.

[293] Jie-Jing Shao, Hao-Ran Hao, Xiao-Wen Yang, and Yu-Feng Li. Learning for long-horizon planning via neuro-symbolic abductive imitation, 2024.

[294] Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100, 1953.

[295] Hong Shen, Haojian Jin, Ángel Alexander Cabrera, Adam Perer, Haiyi Zhu, and Jason I Hong. Designing alternative representations of confusion matrices to support non-expert public understanding of algorithm performance. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW2):1–22, 2020.

[296] Ivaxi Sheth and Samira Ebrahimi Kahou. Auxiliary losses for learning generalizable concept-based models. *Advances in Neural Information Processing Systems*, 36, 2023.

[297] Zheyuan Ryan Shi, Leah Lizarondo, and Fei Fang. A recommender system for crowdsourcing food rescue platforms. In *Proceedings of the web conference 2021*, pages 857–865, 2021.

[298] Andrew Silva, Matthew Gombolay, Taylor Killian, Ivan Jimenez, and Sung-Hyun Son. Optimization methods for interpretable differentiable decision trees applied to reinforcement learning. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pages 1855–1865. PMLR, 2020.

[299] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.

[300] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, 2017.

[301] Aleksandrs Slivkins et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286, 2019.

[302] Reinhardt Smit and Hanlie Smuts. Game-based learning-teaching artificial intelligence to play minecraft: a systematic literature review. 2023.

[303] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, pages 5887–5896. PMLR, 2019.

[304] Bhuman Soni and Philip Hingston. Bots trained to play like a human are more fun. 2008.

[305] Hossein Azari Soufiani, David C Parkes, and Lirong Xia. Preference elicitation for general random utility models. *arXiv preprint arXiv:1309.6864*, 2013.

[306] Sarath Sreedharan, Siddharth Srivastava, and Subbarao Kambhampati. Tldr: Policy summarization for factored ssp problems using temporal abstractions. Proceedings of the 30th International Conference on Automated Planning and Scheduling, 2020.

[307] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.

[308] Jonathan Stray. Aligning ai optimization to community well-being. *International Journal of Community Well-Being*, 3(4):443–463, 2020.

[309] Alexander L Strehl, Carlos Diuk, and Michael L Littman. Efficient structure learning in factored-state mdps. In *Proceedings of the 22nd AAAI Conference on Artificial Intelligence (AAAI-07)*, pages 645–650, 2007.

[310] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.

[311] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[312] Jonathan Sykes and Melissa Federoff. Player-centred game design. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*, pages 1731–1734, 2006.

[313] István Szita. Reinforcement learning in games. In *Reinforcement learning*, pages 539–577. Springer, 2012.

[314] Karan Taneja, Pratyusha Maiti, Sandeep Kakar, Pranav Guruprasad, Sanjeev Rao, and Ashok K Goel. Jill watson: A virtual teaching assistant powered by chatgpt. In *International Conference on Artificial Intelligence in Education*, pages 324–337. Springer, 2024.

[315] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.

[316] Fabien Tencé and Cédric Buche. Automatable evaluation method oriented toward behaviour believability for video games. *arXiv preprint arXiv:1009.0501*, 2010.

[317] The Times of India. Ai-driven traffic signals go live at 3 junctions, pilot rollout by june 15, June 2025. Accessed: 2025-06-05.

[318] Christian Thurau, Christian Bauckhage, and Gerhard Sagerer. Learning human-like movement behavior for computer games. In *Proceedings of the International Conference on the Simulation of Adaptive Behavior*, pages 315–323, 2004.

[319] Erico Tjoa and Cuntai Guan. A survey on explainable artificial intelligence (xai): Towards medical xai. *arXiv preprint, arXiv:1907.07374*, 2019.

[320] Julian Togelius, Georgios N Yannakakis, Sergey Karakovskiy, and Noor Shaker. Assessing believability. In *Believable bots*, pages 215–230. Springer, 2013.

[321] Nicholay Topin, Stephanie Milani, Fei Fang, and Manuela Veloso. Iterative bounding MDPs: Learning interpretable policies via non-interpretable methods. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9923–9931, 2021.

[322] Nicholay Topin and Manuela Veloso. Generation of policy-level explanations for reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.

[323] Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 2023.

[324] Karl Tuyls, Sam Maes, and Bernard Manderick. Reinforcement learning in large state spaces. In *Robot Soccer World Cup*, pages 319–326, 2002.

[325] Amos Tversky and Itamar Simonson. Context-dependent preferences. *Management science*, 39(10):1179–1189, 1993.

[326] Ultralytics. YOLOv5: A state-of-the-art real-time object detection system. `https://docs.ultralytics.com`, 2021.

[327] Berk Ustun and Cynthia Rudin. Learning optimized risk scores. *Journal of Machine Learning Research*, 20(150):1–75, 2019.

[328] William TB Uther and Manuela M Veloso. The lumberjack algorithm for learning linked decision forests. In *The International Symposium on Abstraction, Reformulation, and Approximation*, pages 219–232, 2000.

[329] Teija Vainio. A review of the navigation hci research during the 2000's. *international Journal of Interactive Mobile Technologies (iJIM)*, 4(3), 2010.

[330] Varun Ravi Varma. *Interpretable Reinforcement Learning with the Regression Tsetlin Machine*. PhD thesis, 2021.

[331] Marko Vasic, Andrija Petrovic, Kaiyuan Wang, Mladen Nikolic, Rishabh Singh, and Sarfraz Khurshid. Moët: Interpretable and verifiable reinforcement learning via mixture of expert trees. *arXiv preprint, arXiv:1906.06717*, 2019.

[332] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 10, 2017.

[333] Aravind Venugopal, Stephanie Milani, Fei Fang, and Balaraman Ravindran. MABL: Bi-level latent-variable world model for sample-efficient multi-agent reinforcement learning. In *AAMAS*, 2024.

[334] Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri. Programmatically interpretable reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, pages 5045–5054. PMLR, 2018.

[335] Rajeev Verma, Daniel Barrejón, and Eric Nalisnick. Learning to defer to multiple experts: Consistent surrogate losses, confidence calibration, and conformal ensembles. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2023.

[336] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.

[337] Luis Von Ahn and Laura Dabbish. Labeling images with a computer game. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 319–326, 2004.

[338] Marcel Walch, Julian Frommel, Katja Rogers, Felix Schüssel, Philipp Hock, David Dobbelstein, and Michael Weber. Evaluating vr driving simulation from a player experience perspective. In *Proceedings of the 2017 CHI Conference Extended Abstracts on Human Factors in Computing Systems*, pages 2982–2989, 2017.

[339] Dakuo Wang, Elizabeth Churchill, Pattie Maes, Xiangmin Fan, Ben Shneiderman, Yuanchun Shi, and Qianying Wang. From human-human collaboration to human-ai collaboration: Designing ai systems that can work together with people. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*, pages 1–6, 2020.

[340] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.

[341] Renhao Wang, Jiayuan Mao, Joy Hsu, Hang Zhao, Jiajun Wu, and Yang Gao. Programmatically grounded, compositionally generalizable robotic manipulation. *Proceedings of the International Conference on Learning Representations*, 2023.

[342] Ruiyi Wang*, Stephanie Milani*, Jamie C. Chiu, Jiayin Zhi, Shaun M. Eack, Travis Labrum, Samuel M. Murphy, Nev Jones, Kate Hardy, Hong Shen, Fei Fang, and Zhiyu Chen. Patient-{\Psi}: Using large language models to simulate patients for training mental health professionals. *Empirical Methods in Natural Language Processing*, 2024.

[343] Xinzhi Wang, Huao Li, Hui Zhang, Michael Lewis, and Katia Sycara. Explanation of reinforcement learning model in dynamic multi-agent system. *arXiv preprint arXiv:2008.01508*, 2020.

[344] Marieke GM Weernink, Sarah IM Janus, Janine A Van Til, Dennis W Raisch, Jeannette G Van Manen, and Maarten J IJzerman. A systematic review to identify the use of preference elicitation methods in healthcare decision making. *Pharmaceutical medicine*, 28:175–185, 2014.

[345] Hua Wei, Huaxiu Yao, Guanjie Zheng, and Zhenhui Li. Intellilight: A reinforcement learning approach for intelligent traffic light control. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 2496–2505. ACM, 2018.

[346] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv*, abs/2201.11903, 2022.

[347] Eric Wiewiora. *Reward Shaping*, pages 863–865. Springer US, Boston, MA, 2010.

[348] Erik Wijmans, Manolis Savva, Irfan Essa, Stefan Lee, Ari S Morcos, and Dhruv Batra. Emergence of maps in the memories of blind navigation agents. *arXiv preprint arXiv:2301.13261*, 2023.

[349] John Winder, Stephanie Milani, Matthew Landen, Erebus Oh, Shane Parr, Shawn Squire, Cynthia Matuszek, et al. Planning with abstract learned models while learning transferable subtasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 9992–10000, 2020.

[350] Huai-Ning Wu and Mi Wang. Human-in-the-loop behavior modeling via an integral concurrent adaptive inverse reinforcement learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

[351] Bin-Bin Yang, Song-Qing Shen, and Wei Gao. Weighted oblique decision trees. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 5621–5627, 2019.

[352] Chao-Han Huck Yang, I Danny, Te Hung, Yi Ouyang, and Pin-Yu Chen. Causal inference q-network: Toward resilient reinforcement learning. In *Self-Supervision for Reinforcement Learning Workshop-ICLR 2021*, 2021.

[353] Zhuorui Ye*, Stephanie Milani*, Geoffrey J Gordon, and Fei Fang. Licorice: Label-efficient concept-based interpretable reinforcement learning. In *Proceedings of the International Conference on Learning Representations*, 2025.

[354] Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021.

[355] Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and YI WU. The surprising effectiveness of ppo in cooperative multi-agent games. In *Advances in Neural Information Processing Systems*, volume 35, pages 24611–24624, 2022.

[356] Renos Zabounidis, Joseph Campbell, Simon Stepputtis, Dana Hughes, and Katia P Sycara. Concept learning for interpretable multi-agent reinforcement learning. In *Proceedings of the Conference on Robot Learning*, pages 1828–1837. PMLR, 2023.

[357] Lotfi A Zadeh. Fuzzy logic. *Computer*, 21(4):83–93, 1988.

[358] Anthony Zador, Blake Richards, Bence Ölveczky, Sean Escola, Yoshua Bengio, Kwabena Boahen, Matthew Botvinick, Dmitri Chklovskii, Anne Churchland, Claudia Clopath, et al. Toward next-generation artificial intelligence: Catalyzing the neuroai revolution. *arXiv preprint arXiv:2210.08340*, 2022.

[359] Omar Zaidan, Jason Eisner, and Christine Piatko. Using "annotator rationales" to improve machine learning for text categorization. In *Human language technologies 2007: The conference of the North American chapter of the association for computational linguistics; proceedings of the main conference*, pages 260–267, 2007.

[360] Mateo Espinosa Zarlenga, Katherine M Collins, Krishnamurthy Dvijotham, Adrian Weller, Zohreh Shams, and Mateja Jamnik. Learning to receive help: Intervention-aware concept embedding models. *Advances in Neural Information Processing Systems*, 2023.

[361] Hengzhe Zhang, Aimin Zhou, and Xin Lin. Interpretable policy derivation for reinforcement learning based on evolutionary feature synthesis. *Complex & Intelligent Systems*, pages 1–13, 2020.

[362] Li Zhang, Xin Li, Mingzhong Wang, and Andong Tian. Off-policy differentiable logic reinforcement learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 617–632. Springer, 2021.

[363] Yunqi Zhao, Igor Borovikov, Jason Rupert, Caedmon Somers, and Ahmad Beirami. On multi-agent learning in team sports games. *arXiv preprint arXiv:1906.10124*, 2019.

[364] Zhibing Zhao, Haoming Li, Junming Wang, Jeffrey Kephart, Nicholas Mattei, Hui Su, and Lirong Xia. A cost-effective framework for preference elicitation and aggregation. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, 2018.

[365] Jianlong Zhou, Amir H Gandomi, Fang Chen, and Andreas Holzinger. Evaluating the quality of machine learning explanations: A survey on methods and metrics. *Electronics*, 10(5):593, 2021.

[366] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. *arXiv preprint arXiv:2305.17144*, 2023.

[367] Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of maximum causal entropy. *Proceedings of the International Confer- ence on Machine Learning*, 2010.

[368] Evelyn Zuniga*, Stephanie Milani*, Guy Leroy*, Jaroslaw Rzepecki, Raluca Georgescu, Ida Momennejad, Dave Bignell, Mingfei Sun, Alison Shaw, Gavin Costello, Mikhail Jacob, Sam Devlin, and Katja Hofmann. How humans perceive human-like behavior in video game navigation. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*, pages 1–11, 2022.