
Learning to Assist Humans without Inferring Rewards

Vivek Myers¹ Evan Ellis¹

Sergey Levine¹ Benjamin Eysenbach² Anca Dragan¹

¹UC Berkeley ²Princeton University

Abstract

Assistive agents should make humans’ lives easier. Classically, such assistance is studied through the lens of inverse reinforcement learning, where an assistive agent (e.g., a chatbot, a robot) infers a human’s intention and then selects actions to help the human reach that goal. This approach requires inferring intentions, which can be difficult in high-dimensional settings. We build upon prior work that studies assistance through the lens of empowerment: an assistive agent aims to maximize the influence of the human’s actions such that they exert a greater control over the environmental outcomes and can solve tasks in fewer steps. We lift the major limitation of prior work in this area — scalability to high-dimensional settings — with contrastive successor representations. We formally prove that these representations estimate a similar notion of empowerment to that studied by prior work and provide a ready-made mechanism for optimizing it. Empirically, our proposed method outperforms prior methods on synthetic benchmarks, and scales to *Overcooked*, a cooperative game setting. Theoretically, our work connects ideas from information theory, neuroscience, and reinforcement learning, and charts a path for representations to play a critical role in solving assistive problems.¹

1 Introduction

AI agents deployed in the real world should be helpful to humans. When we know the utility function of the humans an agent could interact with, we can directly train assistive agents through reinforcement learning with the known human objective as the agent’s reward. In practice, agents rarely have direct access to a scalar reward corresponding to human preferences (if such a consistent model even exists) [1], and must infer them from human behavior [2, 3]. This inference can be challenging, as humans may act suboptimally with respect to their stated goals, not know their goals, or have changing preferences [4]. Optimizing a misspecified reward function can have poor consequences [5].

An alternative paradigm for assistance is to train agents that are *intrinsically* motivated to assist humans, rather than directly optimizing a model of their preferences. An analogy can be drawn to a parent raising a child. A good parent will empower the child to make impactful decisions and flourish, rather than proscribing an “optimal” outcome for the child. Likewise, AI agents might seek to *empower* the human agents they interact with, maximizing their capacity to change the environment [6]. In practice, concrete notions of empowerment can be difficult to optimize as an objective, requiring extensive modeling assumptions that don’t scale well to the high-dimensional settings deep reinforcement learning agents are deployed in.

What is a good intrinsic objective for assisting humans that doesn’t require these assumptions? We propose a notion of assistance based on maximizing the influence of the human’s actions on the

¹Code: https://github.com/vivekmyers/empowerment_successor_representations
Website: <https://empowering-humans.github.io>

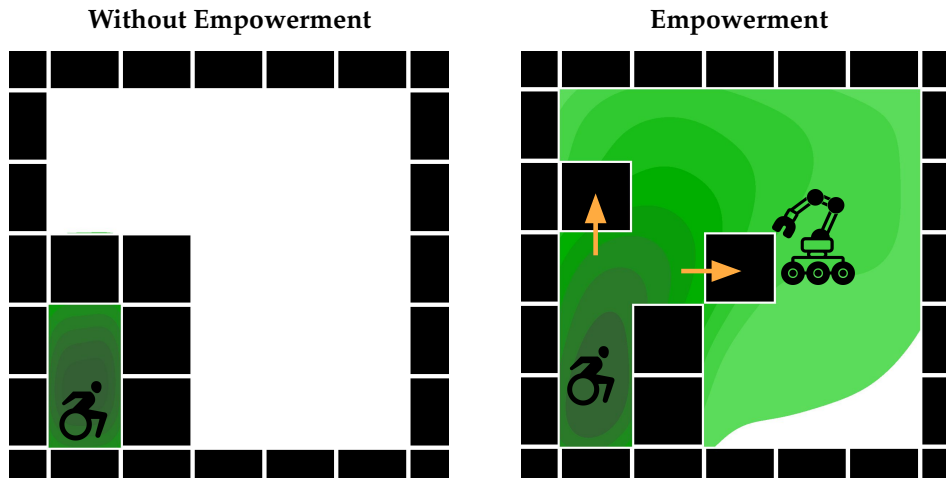


Figure 1: We propose an algorithm training assistive agents to empower human users—the assistant should take actions that enable human users to visit a wide range of future states, and the human’s actions should exert a high degree of influence over the future outcomes. Our algorithm scales to high-dimensional settings, opening the door to building assistive agents that need not directly reason about human intentions.

environment. This approach only requires one structural assumption: the AI agent is interacting with an environment where there is a notion of actions taken by the human agent—a more general setting than the case where we model the human actions as the outcome of some optimization procedure, as in inverse RL [7, 8] or preference-based RL [9].

Prior work has studied many effective objectives for empowerment. For instance, Du et al. [6] approximates human empowerment as the variance in the final states of random rollouts. Despite excellent results in certain settings, this approach can be challenging to scale to higher dimensional settings, and does not necessarily enable human users to achieve the goals they want to achieve. By contrast, our approach exclusively empowers the human with respect to the distribution of (useful) behaviors induced by their current policy, and can be implemented through a simple objective derived from contrastive successor features, which can then be optimized with scalable deep reinforcement learning (Fig. 1). We provide a theoretical framework connecting our objective to prior work on empowerment and goal inference, and empirically show that agents trained with this objective can assist humans in the Overcooked environment [10] as well as the obstacle gridworld assistance benchmark proposed by Du et al. [6].

Our core contribution is a novel objective for training agents that are intrinsically motivated to assist humans without requiring a model of the human’s reward function. Our objective, Empowerment via Successor Representations (ESR), maximizes the influence of the human’s actions on the environment, and, unlike past approaches for assistance without reward inference, is based on a scalable model-free objective that can be derived from learned successor features that encode which states the human is likely to want to reach given their current action. Our objective empowers the human to reach the desired states, not all states, without assuming a human model. We analyze this objective in terms of empowerment and goal inference, drawing novel mathematical connections between time-series representations, decision-making, and assistance. We empirically show that agents trained with our objective can assist humans in two benchmarks proposed by past work: the Overcooked environment [10] and an obstacle-avoidance gridworld [6].

2 Related Work

Our approach broadly connects ideas from contrastive representation learning and intrinsic motivation to the problem of assisting humans.

Assistive Agents. There are two lines of past work on assistive agents that are most relevant.

The first line of work focuses on the setting of an assistance game [2], where a robot (AI) agent tries to optimize a human reward of which it is initially unaware. Practically, inverse reinforcement learning (IRL) can be used in such a setting to infer the human’s reward function and assist the human in achieving their goals [3]. The key challenge with this approach is that it requires modeling the human’s reward function. This can be difficult in practice, especially if the human’s behavior is not well-modeled by the reward architecture. Slightly misspecified reward functions can lead to catastrophic outcomes (i.e., directly harmful behavior in the assistance context) [11–13]. By contrast, our approach does not require modeling the human’s reward function.

The second line of work focuses on empowerment-like objectives for assistance and shared autonomy. Empowerment generally refers to a measure of an agent’s ability to influence the environment [14, 15]. In the context of assistance, Du et al. [6] show one such approximation of empowerment (AvE) can be approximated in simple environments through random rollouts to assist humans. Meanwhile, empowerment-like objectives have been used in shared autonomy settings to assist humans with teleoperation [16] and general assistive interfaces [17]. A key limitation of these approaches for general assistance is they only model empowerment over one time step. Our approach enables a more scalable notion of empowerment that can be computed over multiple time steps.

Intrinsic Motivation. Intrinsic motivation broadly refers to agents that accomplish behaviors in the absence of an externally-specified reward or task [18]. Common applications of intrinsic motivation in single-agent reinforcement learning include exploration and skill discovery [19–21], empowerment [15, 14], and surprise minimization [22, 23, 15]. When applied to settings with humans, these objectives may lead to antisocial behavior [5]. Our approach applies intrinsic motivation to the setting of assisting humans, where the agent’s goal is an empowerment objective—to maximize the human’s ability to change the environment.

Information-theoretic Decision Making. Information-theoretic approaches have seen broad applicability across unsupervised reinforcement learning [24, 15, 19]. These methods have been applied to goal-reaching [25], skill discovery [26, 27, 20, 28, 29], and exploration [21, 30, 31]. In the context of assisting humans, information-theoretic methods have primarily been used to reason about the human’s goals or rewards [32–34].

Our approach is made possible by advances in contrastive representation learning for efficient estimation of the mutual information of sequence data [35]. While these methods have been widely used for representation learning [36, 37] and reinforcement learning [38–41], to the best of our knowledge prior work has not used these contrastive techniques for learning assistive agents.

3 The Information Geometry of Empowerment

We will first state a general notion of an assistive setting, then show how an empowerment objective based on learned successor representations can be used to assist humans without making assumptions about the human following an underlying reward function. In Section 5, we provide empirical evidence supporting these claims.

3.1 Preliminaries

Formally, we adapt the notation of Hadfield-Menell et al. [2], and assume a “robot” (**R**) and “human” (**H**) policy are training together in an MDP $M = (\mathcal{S}, \mathcal{A}_H, \mathcal{A}_R, R, P, \gamma)$. The states s consist of the joint states of the robot and the human; we do not have separate observations for the human and robot. At any state $s \in \mathcal{S}$, the robot policy selects actions distributed according to $\pi_R(a^R | s)$ for $a^R \in \mathcal{A}_R$ and the human selects actions from $\pi_H(a^H | s)$ for $a^H \in \mathcal{A}_H$. The transition dynamics are defined by a distribution $P(s' | s, a^H, a^R)$ over the next state $s' \in \mathcal{S}$ given the current state $s \in \mathcal{S}$ and actions $a^H \in \mathcal{A}_H$ and $a^R \in \mathcal{A}_R$, as well as an initial state distribution $P(s_0)$. For notational convenience, we will additionally define random variables \mathfrak{s}_t to represent the state at time t , and $\mathfrak{a}_t^R \sim \pi_R(\cdot | \mathfrak{s}_t)$ and $\mathfrak{a}_t^H \sim \pi_H(\cdot | \mathfrak{s}_t)$ to represent the human and robot actions at time t , respectively.

Empowerment. Our work builds on a long line of prior methods that use information theoretic objectives for RL. Specifically, we adopt *empowerment* as an objective for training an assistive agent [6, 42, 43]. This section provides the mathematical foundations for empowerment, as developed

in prior work. Our work will build on the prior work by (1) providing an information geometric interpretation of what empowerment does (Section 3.3) and (2) providing a scalable algorithm for estimating and optimizing empowerment, going well beyond the gridworlds studied in prior work.

The idea behind empowerment is to think about the changes that an agent can effect on a world; an agent is more empowered if it can effect a larger degree of change over future outcomes. Following prior work [25, 43, 42], we measure empowerment by looking at how much the actions taken *now* affect outcomes *in the future*. An agent with a high degree of empowerment exerts a high degree of control of the future states by simply changing the actions taken now. Like prior work, we measure this degree of control through the mutual information $I(\mathfrak{s}^+; a^{\mathbf{H}})$ between the current action $a^{\mathbf{H}}$ and the future states \mathfrak{s}^+ . Note that these future states might occur many time steps into the future.

Empowerment depends on several factors: the environment dynamics, the choice of future actions, the current state, and other agents in the environment. Different problem settings involve maximizing empowerment using these different factors. In this work, we study the setting where a “human” agent and a “robot” agent collaborate in an environment; the robot will aim to maximize the empowerment of the human. This problem setting was introduced in prior work [6]. Compared with other mathematical frameworks for learning assistive agents [44], framing the problem in terms of empowerment means that the assistive agent need not infer the human’s underlying intention, an inference problem that is typically challenging [45, 46].

We now define our objective. To do this, we introduce random variable \mathfrak{s}^+ , which corresponds to a state sampled $K \sim \text{Geom}(1 - \gamma)$ steps into the future under the behavior policies $\pi_{\mathbf{H}}$ and $\pi_{\mathbf{R}}$. We will use $\rho(\mathfrak{s}^+ | s_t)$ to denote the density of this random variable; this density is sometimes referred to as the discounted state occupancy measure. We will use mutual information to measure how much the action a_t at time t changes this distribution:

$$I(a_t^{\mathbf{H}}; \mathfrak{s}^+ | s_t) \triangleq \mathbb{E}_{s_t, s_{t+k}, a_t^{\mathbf{H}}, a_t^{\mathbf{R}}} \left[\log \frac{P(\mathfrak{s}_{t+K} = s_{t+k} | \mathfrak{s}_t = s_t, a_t^{\mathbf{H}} = a_t)}{P(\mathfrak{s}_{t+K} = s_{t+k} | \mathfrak{s}_t = s_t)} \right]. \quad (1)$$

Our overall objective is *empowerment*, $\mathcal{E}(\pi_{\mathbf{H}}, \pi_{\mathbf{R}})$: the mutual information between the human’s actions and the future states \mathfrak{s}^+ while interacting with the robot:

$$\mathcal{E}(\pi_{\mathbf{H}}, \pi_{\mathbf{R}}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t I(a_t^{\mathbf{H}}; \mathfrak{s}^+ | s_t) \right]. \quad (2)$$

Note that this objective resembles an RL objective: we do not just want to maximize this objective greedily at each time step, but rather want the assistive agents to take actions now that help the human agent reach states where it will have high empowerment in the future.

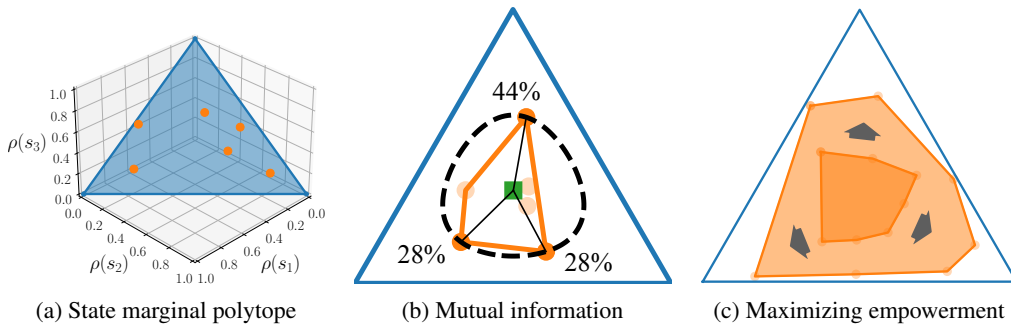


Figure 2: **The Information Geometry of Empowerment**, illustrating the analysis in Section 3.3. (Left) For a given state s_t and assistant policy $\pi_{\mathbf{R}}$, we plot the distribution over future states for 6 choices of the human policy $\pi_{\mathbf{H}}$. In a 3-state MDP, we can represent each policy as a vector lying on the 2-dimensional probability simplex. We refer to the set of all possible state distributions as the *state marginal polytope*. (Center) Mutual information corresponds to the distance between the center of the polytope and the vertices that are maximally far away. (Right) Empowerment corresponds to maximizing the size of this polytope. For example, when an assistive agent moves an obstacle out of a human user’s way, the human user can spend more time at desired state.

3.2 Intuition and Geometry of Empowerment

Intuitively, the assistive agent should aim to maximize the number of future outcomes. We will mathematically quantify this in terms of the discounted state occupancy measure, $\rho^\pi(\mathfrak{s}^+ | s)$. Intuitively, an agent has a large empowerment if the future states for one action are very different from the future actions after taking a different action; i.e., when $\rho(a_t = a_1; \mathfrak{s}^+ | s_t)$ is quite different from $\rho(a_t = a_2; \mathfrak{s}^+ | s_t)$ for actions $a_1 \neq a_2$. The mutual information (Eq. (1)) quantifies this degree of control: $I(a_t; \mathfrak{s}^+ | s_t)$.

One way of understanding this mutual information is through *information geometry* [47, 48, 48, 49]. For a fixed current state s_t , assistant policy π_R and human policy π_H , each potential action a_t that the human takes induces a different distribution over future states: $\rho^{\pi_R, \pi_H}(\mathfrak{s}^+ | s_t, a_t)$. We can think about the set of these possible distributions: $\{\rho^{\pi_R, \pi_H}(\mathfrak{s}^+ | s_t, a_t) | a_t \in \mathcal{A}\}$. Figure 2 (*Left*) visualizes this distribution on a probability simplex for 6 choices of action a_t . If we look at any possible distribution over actions, then this set of possible future distributions becomes a polytope (see orange polygon in Fig. 2 (*Center*)).

Intuitively, the mutual information $I(a_t; \mathfrak{s}^+ | s_t)$ used to define our empowerment objective corresponds to the *size* or *volume* of this state marginal polytope. This intuition can be formalized by using results from information geometry [50–52]. The human policy $\pi_H(a_t | s_t)$ places probability mass on the different points in Figure 2 (*Center*). Maximizing the mutual information corresponds to “picking out” the state distributions that are maximally spread apart (see probabilities in Fig. 2 (*Center*)). To make this formal, define

$$\rho(\mathfrak{s}^+ | s_t) \triangleq \mathbb{E}_{\pi(a_t | s_t)}[\rho(\mathfrak{s}^+ | s_t, a_t)] \quad (3)$$

as the *average* state distribution from taking the human’s actions (see green square in Fig. 2 (*Center*)).

Remark 3.1. *Mutual information corresponds to the distance between the average state distribution (Eq. 3) and the furthest achievable state distributions:*

$$I(a_t; \mathfrak{s}^+ | s_t) = \max_{a_t} D_{KL}(\rho(a_t; \mathfrak{s}^+ | s_t) || \rho(\mathfrak{s}^+ | s_t)) \triangleq d_{max}. \quad (4)$$

This distance is visualized as the black lines in Fig. 2. When we talk about the “size” of the state marginal polytope, we are specifically referring to the length of these black lines (as measured with a KL divergence).

This sort of mutual information is a way for measuring the degree of control that an agent exerts on an environment. This measure is well defined for any agent/policy; that agent need not be maximizing mutual information, and could instead be maximizing some arbitrary reward function. This point is important in our setting: this means that the assistive agent can estimate and maximize the empowerment of the human user *without having to infer what reward function the human is trying to maximize*.

Finally, we come back to our empowerment objective, which is a discounted sum of the mutual information terms that we have been analyzing above. This empowerment objective says that the human is more empowered when this set has a larger size — i.e., the human can visit a wider range of future state (distributions). The empowerment objective says that the assistive agent should act to try to maximize the size of this polytope. Importantly, this maximization problem is done sequentially: the assistive agent wants the size of this polytope to be large both at the current state and at future states; the human’s actions should exert a high degree of influence over the future outcomes both now and in the future. Thus, our overall objective looks at a sum of these mutual informations.

Not only does this analysis provides a geometric picture for what empowerment is doing, it also lays the groundwork for formally relating empowerment to reward.

3.3 Relating Empowerment to Reward

In this section we take aim at the question: when humans are well-modeled as optimizing a reward function, when does maximizing empowerment help humans maximize their rewards? Answering this question is important because for empowerment to be a safe and effective assistive objective, it should enable the human to better achieve their goals. We show that under certain assumptions,

empowerment yields a provable lower bound on the average-case reward achieved by the human for sufficiently long-horizon empowerment (i.e., $\gamma \rightarrow 1$).

For constructing the formal bound, we suppose the human is Boltzmann-rational [53, 54] with respect to some reward function $R \sim \mathcal{R}$, where \mathcal{R} is some distribution that could be interpreted as a prior over the human’s objective, a set of skills the human may try and carry out, or a population of humans with different objectives that the agent could be interacting with. Our quantity of interest, the average-case reward achieved by the human with our empowerment objective, is given by

$$\mathcal{J}_{\pi_{\mathbf{R}}}^{\gamma}(\pi_{\mathbf{H}}) = \mathbb{E}_{R \sim \mathcal{R}} \left[\mathbb{E}_{s_0 \sim p_0} [V_{R, \gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_0)] \right] \quad (5)$$

where $V_{R, \gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_0)$ is the value function of the human policy $\pi_{\mathbf{H}}$ under the reward function R when interacting with $\pi_{\mathbf{R}}$. Recalling Eq. (2), we will express the overall empowerment objective we are trying to relate to Eq. (5) as

$$\mathcal{E}_{\gamma}(\pi_{\mathbf{H}}, \pi_{\mathbf{R}}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t I(\mathbf{s}^+; \mathbf{a}_t^{\mathbf{H}} | \tilde{\mathbf{s}}_t) \right]. \quad (6)$$

This notation is formalized in Appendix B.

The two key assumptions used in our analysis are Assumption 3.1, which states that the human will optimize for behaviors that uniformly cover the state space, and Assumption 3.2, which simply states that with infinite time, the human will be able to reach any state in the state space.

Assumption 3.1 (Skill Coverage). *The rewards $R \sim \mathcal{R}$ are uniformly distributed over the scaled $|\mathcal{S}|$ -simplex $\Delta^{|\mathcal{S}|}$ such that:*

$$\left(R + \frac{1}{|\mathcal{S}|} \right) \left(\frac{1}{1-\gamma} \right) \sim \text{Unif}(\Delta^{|\mathcal{S}|}) = \text{Dirichlet}(\underbrace{1, 1, \dots, 1}_{|\mathcal{S}| \text{ times}}). \quad (7)$$

Assumption 3.2 (Ergodicity). *For some $\pi_{\mathbf{H}}, \pi_{\mathbf{R}}$, we have*

$$\mathbb{P}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(\mathbf{s}^+ = s | s_0) > 0 \quad \text{for all } s \in \mathcal{S}, \gamma \in (0, 1). \quad (8)$$

Our main theoretical result is Theorem 3.1, which shows that under these assumptions, maximizing empowerment yields a lower bound on the (squared) average-case reward achieved by the human for sufficiently large γ . In other words, for a sufficiently long empowerment horizon, the empowerment objective Eq. (2) is a meaningful proxy for reward maximization.

Theorem 3.1. *Under Assumption 3.1 and Assumption 3.2, for sufficiently large γ and any $\beta > 0$,*

$$\mathcal{E}_{\gamma}(\pi_{\mathbf{H}}, \pi_{\mathbf{R}})^{1/2} \leq (\beta/e) \mathcal{J}_{\pi_{\mathbf{R}}}^{\gamma}(\pi_{\mathbf{H}}). \quad (9)$$

The proof is in Appendix B.1 To the best of our knowledge, this result provides the first formal link between empowerment maximization and reward maximization. This motivates us to develop a scalable algorithm for empowerment maximization, which we introduce in the following section.

4 Estimating and Maximizing Empowerment with Contrastive Representations

Directly computing Eq. (2) would require access to the human policy, which we don’t have. Therefore, we want a tractable estimation that still performs well in large environments which are more difficult to model due to the exponentially increasing set of possible future states. To better-estimate empowerment, we learn contrastive representations that encode information about which future states are likely to be reached from the current state. These contrastive representations learn to model mutual information between the current state, action, and future state, which we then use to compute the empowerment objective.

4.1 Estimating Empowerment

To estimate this empowerment objective, we need a way of learning the probability ratio inside the expectation. Prior methods such as Du et al. [6] and Salge et al. [42] rollout possible future states

and compute a measure of their variance as a proxy for empowerment, however this doesn't scale when the environment becomes complex. Other methods learn a dynamics model, which also doesn't scale when dynamics become challenging to model [27]. Modeling these probabilities directly is challenging in settings with high-dimensional states, so we opt for an indirect approach. Specifically, we will learn representations that encode two probability ratios. Then, we will be able to compute the desired probability ratio by combining these other probability ratios.

Our method learns three representations:

1. $\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})$ — This representation can be understood as a sort of latent-space model, predicting the future representation given the current state s and the human's current action $a^{\mathbf{H}}$ as well as the robot's current action $a^{\mathbf{R}}$.
2. $\phi'(s, a^{\mathbf{R}})$ — This representation can be understood as an uncontrolled model, predicting the representation of a future state without reference to the current human action $a^{\mathbf{H}}$. This representation is analogous to a value function.
3. $\psi(s^+)$ — This is a representation of a future state.

We will learn these three representations with two contrastive losses, one that aligns $\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}}) \leftrightarrow \psi(s^+)$ and one that aligns $\phi'(s, a^{\mathbf{R}}) \leftrightarrow \psi(s^+)$

$$\max_{\phi, \phi', \psi} \mathbb{E}_{\{(s_i, a_i, s'_i) \sim P(\mathfrak{s}_t, \mathbf{a}_t^{\mathbf{H}}, \mathfrak{s}_{t+k})\}_{i=1}^N} [\mathcal{L}_c(\{\phi(s_i, a_i)\}, \{\psi(s'_i)\}) + \mathcal{L}_c(\{\phi'(s_i)\}, \{\psi(s'_i)\})],$$

where the contrastive loss \mathcal{L}_c is the symmetrized infoNCE objective [35]:

$$\mathcal{L}_c(\{x_i\}, \{y_j\}) \triangleq \sum_{i=1}^N \left[\log \left(\frac{e^{x_i^T y_i}}{\sum_{j=1}^N e^{x_i^T y_j}} \right) + \log \left(\frac{e^{x_i^T y_i}}{\sum_{j=1}^N e^{x_j^T y_i}} \right) \right]. \quad (10)$$

We have colored the index j for clarity. At convergence, these representations encode two probability ratios [24], which we will ultimately be able to use to estimate empowerment (Eq. 2):

$$\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})^T \psi(g) = \log \left[\frac{P(\mathfrak{s}_{t+K} = g \mid \mathfrak{s}_t = s, \mathbf{a}_t^{\mathbf{H}} = a^{\mathbf{H}}, \mathbf{a}_t^{\mathbf{R}} = a^{\mathbf{R}})}{C_1 P(\mathfrak{s}_{t+K} = g)} \right] \quad (11)$$

$$\phi'(s, a^{\mathbf{R}})^T \psi(g) = \log \left[\frac{P(\mathfrak{s}_{t+K} = s_{t+k} \mid \mathfrak{s}_t = s_t, \mathbf{a}_t^{\mathbf{R}} = a^{\mathbf{R}})}{C_2 P(\mathfrak{s}_{t+K} = g)} \right]. \quad (12)$$

Note that our definition of empowerment (Eq. 2) is defined in terms of similar probability ratios. The constants C_1 and C_2 will mean that our estimate of empowerment may be off by an additive constant, but that constant will not affect the solution to the empowerment maximization problem.

4.2 Estimating Empowerment with the Learned Representations

To estimate empowerment, we will look at the difference between these two inner products:

$$\begin{aligned} & \phi(s_{t+K}, a^{\mathbf{R}}, a^{\mathbf{H}})^T \psi(g) - \phi(s_{t+K}, a^{\mathbf{R}})^T \psi(g) \\ &= \log P(s_{t+K} \mid s, a^{\mathbf{H}}) - \log C_1 - \log P(s_{t+K}) - \log P(s_{t+K} \mid s) + \log C_2 + \log P(s_{t+K}) \\ &= \log \frac{P(s_{t+K} \mid s, a^{\mathbf{H}})}{P(s_{t+K} \mid s)} + \log \frac{C_2}{C_1}. \end{aligned}$$

Note that the expected value of the first term is the *conditional* mutual information $I(s_{t+K}; a^{\mathbf{H}} \mid s)$. Our empowerment objective corresponds to averaging this mutual information across all the visited states. In other words, our objective corresponds to an RL problem, where empowerment corresponds to the expected discounted sum of these log ratios:

$$\begin{aligned} \mathcal{E}(\pi_{\mathbf{H}}, \pi_{\mathbf{R}}) &= \mathbb{E}_{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}} \left[\sum_{t=0}^{\infty} \gamma^t I(s^+; a_t^{\mathbf{H}} \mid s_t) \right] \\ &\approx \mathbb{E}_{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}} \left[\sum_{t=0}^{\infty} \gamma^t (\phi(s_t, a^{\mathbf{R}}, a^{\mathbf{H}}) - \phi(s_t, a^{\mathbf{R}}))^T \psi(g) - \log \frac{C_2}{C_1} \right]. \end{aligned}$$

The approximation above comes from function approximation in learning the Bayes optimal representations. Again, note that the constants C_1 and C_2 do not change the optimization problem. Thus, to maximize empowerment we will apply RL to the assistive agent $\pi_{\mathbf{R}}(a \mid s)$ using a reward function

$$r(s, a^{\mathbf{R}}) = (\phi(s_t, a^{\mathbf{R}}, a^{\mathbf{H}}) - \phi(s_t, a^{\mathbf{R}}))^T \psi(g). \quad (13)$$

Algorithm 1: Empowerment via Successor Representations (ESR)

Input: Human policy $\pi_{\text{H}}(a | s)$
Randomly initialize assistive agent policy $\pi_{\text{R}}(a | s)$, and representations $\phi(s, a^{\text{R}}, a^{\text{H}})$, $\psi(s, a^T)$, and $\psi(g)$.
Initialize replay buffer \mathcal{B} .
while not converged **do**
 Collect a trajectory of experience with human policy and assistive agent policy, store in replay buffer \mathcal{B} .
 Update representations $\phi(s, a^{\text{R}}, a^{\text{H}})$, $\psi(s, a^T)$, and $\psi(g)$ with the contrastive losses in Eq. (10).
 Update $\pi_{\text{R}}(a | s)$ with RL using reward function $r(s, a^{\text{R}}, a^{\text{H}}) = (\phi(s, a^{\text{R}}, a^{\text{H}}) - \phi'(s, a^{\text{R}}))^T \psi(g)$.
Return: Assistive policy $\pi_{\text{R}}(a | s)$.

4.3 Algorithm Summary

We propose an actor-critic method for learning the assistive agent. Our method will alternate between updating these contrastive representations and using them to estimate a reward function (Eq. (13)) that is optimized via RL. We summarize the algorithm in Algorithm 1. In practice, we use SAC [55] as our RL algorithm. In our experiments, we will also study the setting where the human user updates their policy alongside the assistive agent.

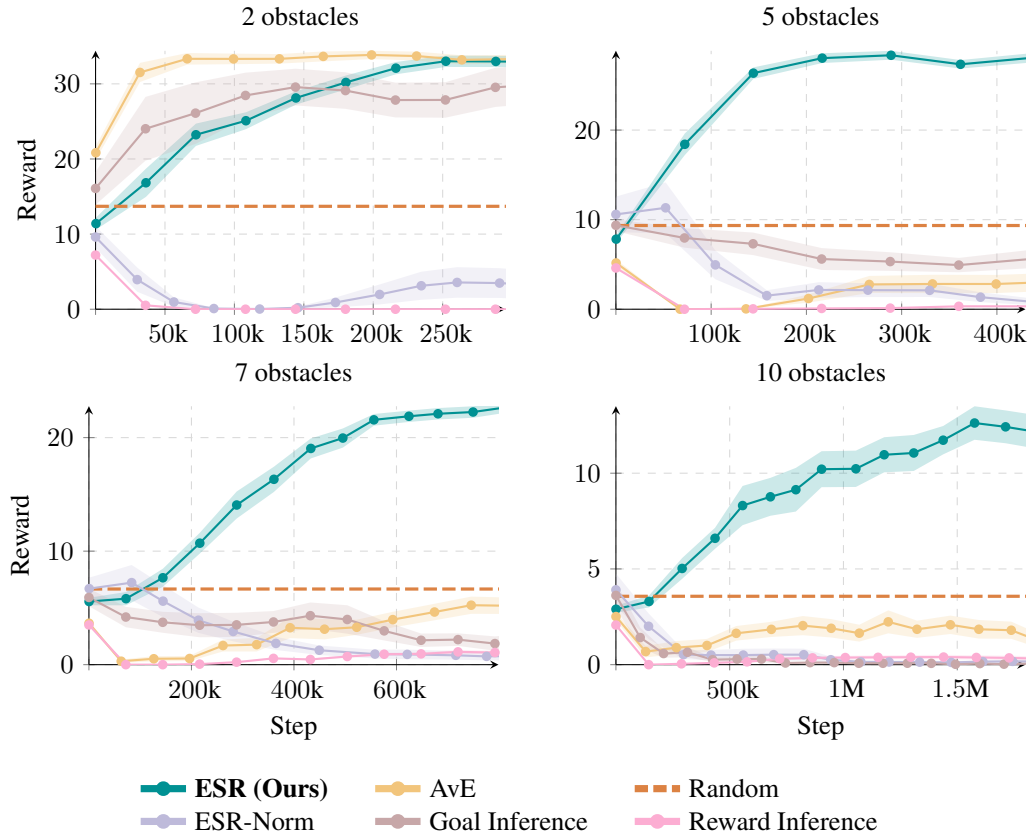


Figure 3: We apply our method to the benchmark proposed in prior work [6], visualized in Fig. 4a. The four subplots show variant tasks of increasing complexity (more blocks), (± 1 SE). We compare against AvE [6], the Goal Inference baseline from [6] which assumes access to a world model, and Reward Inference [56] where we recover the reward from a learned q-value. These prior approaches fail on all except the easiest task, highlighting the importance of scalability.

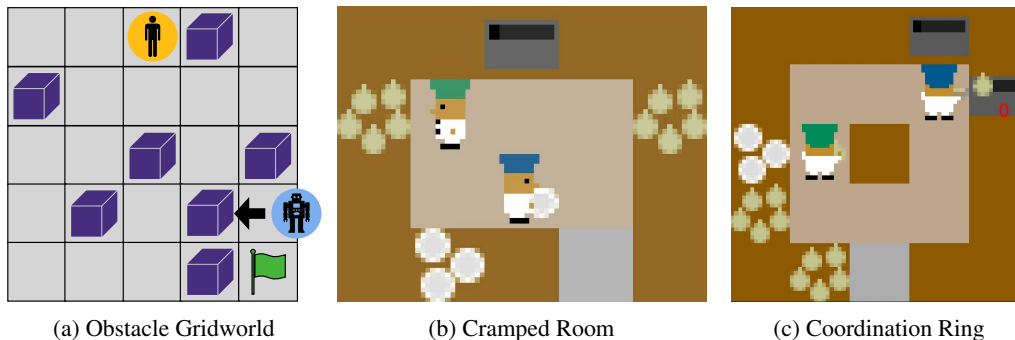


Figure 4: (a) The modified environment from Du et al. [6] scaled to $N = 7$ blocks, and (b, c) the two layouts of the Overcooked environment [10].

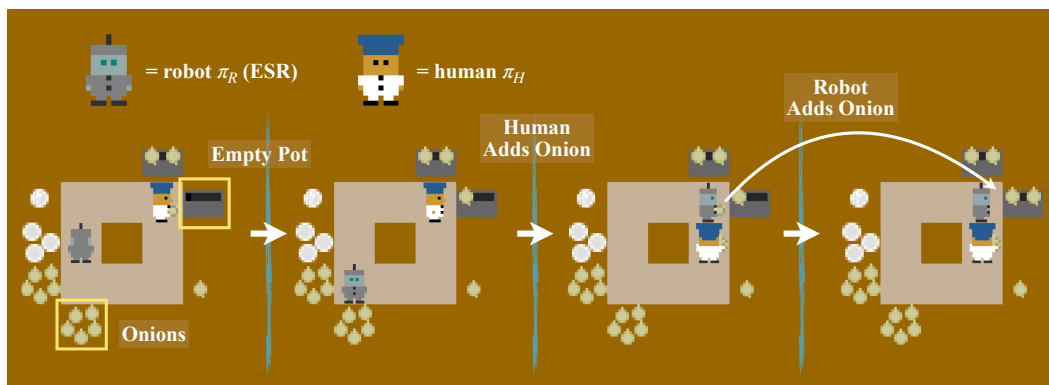


Figure 5: In Coordination Ring, our ESR agent learns to wait for the human to add an onion to the pot, and then adds one itself. There is another pot at the top which is nearly full, but the empowerment agent takes actions to maximize the impact of the human’s actions, and so follows the lead of the human by filling the empty pot.

5 Experiments

We seek to answer two questions with our experiments. *First*, does our approach enable assistance in standard cooperation benchmarks? *Second*, does our approach scale to harder benchmarks where prior methods fail?

Our experiments will use two benchmarks designed by prior work to study assistance: the obstacle gridworld [6] and Overcooked [10]. Our main **baseline** is AvE [6], a prior empowerment-based method. Our conjecture is that both methods will perform well on the lower-dimensional gridworld task, and that our method will scale more gracefully to the higher dimensional Overcooked environment. We will also compare against a naïve baseline where the assistive agent acts randomly.

5.1 Do contrastive successor representations effectively estimate empowerment?

We test our approach in the assistance benchmark suggested in Du et al. [6]. The human (orange) is tasked with reaching a goal state (green) while avoiding the obstacles (purple). The AI assistant can move blocks one step at a time in any direction [6]. While the original benchmark used $N = 2$ obstacles, we will additionally evaluate on harder versions of this task with $N = 5, 7, 10$ obstacles. We show results in Fig. 3. On the easiest task, both our method and AvE achieve similar asymptotic reward, though our method learns more slowly than AvE. However, on the tasks with moderate and high degrees of complexity, our approach (ESR) achieves significantly higher rewards than AvE, which performs worse than a random controller. These experiments support our claim that contrastive successor representations provide an effective means for estimating empowerment, and hint that ESR might be well suited for solving higher dimensional tasks.

5.2 Does our approach scale to tasks with image-based observations?

Our second set of experiments look at scaling ESR to the image-based Overcooked environment. Since contrastive learning is often applied to image domains, we conjectured that ESR would scale gracefully to this setting. We will evaluate our approach in assisting a human policy trained with behavioral cloning taken from Laidlaw and Dragan [57]. The human prepares dishes by picking up ingredients and cooking them on a stove, while the AI assistant moves ingredients and dishes around the kitchen. We focus on two environments within this setting: a cramped room where the human must pass ingredients and dishes through a narrow corridor, and a coordination ring where the human must pass ingredients and dishes around a ring-shaped kitchen (Figs. 4b and 4c). As before, we compare with AvE as well as a naïve random controller. We report results in Table 1. On both tasks, we observe that our approach achieves higher rewards than AvE baseline, which performs no better than a random controller. In Fig. 5, we show an example of one of the collaborative behaviors learned by ESR. Taken together with the results in the previous setting, these results highlight the scalability of ESR to higher dimensional problems.

Table 1: Overcooked Results

Layout	ESR (Ours)	Reward Inference	AvE	Random
Asymmetric Advantages	72.00 ± 5.37	60.33 ± 0.26	36.71 ± 1.71	59.36
Coordination Ring	8.40 ± 0.69	5.96 ± 0.20	5.69 ± 0.93	6.02
Cramped Room	91.33 ± 4.08	39.24 ± 0.35	5.13 ± 1.31	69.26

6 Discussion

One of the most important problems in AI today is equipping AI agents with the capacity to assist humans achieve their goals. While much of the prior work in this area requires inferring the human’s intention, our work builds on prior work in studying how an assistive agent can *empower* a human user without inferring their intention. Relative to prior methods, we demonstrate how empowerment can be readily estimated using contrastive learning, paving the way for deploying these techniques on high-dimensional problems.

Limitations. One of the main limitations of our approach is the assumption that the assistive agent has access to the human’s actions, which could be challenging to observe in practice. Automatically inferring the human’s actions remains an important problem for future work. A second limitation is that the method is currently an on-policy method, in the sense that the assistive agent has to learn by trial and error. A third limitation is that the ESR formulation assumes that both agents share the same state space. In many cases the empowerment objective will still lead to desirable behavior, however, care must be taken in cases where the agent can restrict the information in its own observations, which could lead to reward hacking. Finally, our experiments do not test our method against real humans, whose policies may differ from the simulated policies. In the future, we plan to investigate techniques from off-policy evaluation and cooperative game theory to enable faster learning of assistive agents with fewer trials. We also plan to test the ESR objective in environments with partial observability over the human’s state.

Safety risks. Perhaps the main risk involved with maximizing empowerment is that it may be at odds with a human’s agents goal, especially in contexts where the pursuit of that goal limits the human’s capacity to pursue other goals. For example, a family choosing to have a kid has many fewer options over where they can travel for vacation, yet we do not want assistive agents to stymie families from having children.

One key consideration is *whom* should be empowered. The present paper assumes there is a single human agent. Equivalently, this can be seen as maximizing the empowerment of all exogenous agents. However, it is easy to adapt the proposed method to maximize the empowerment of a single target individual. Given historical inequities in the distribution of power, practitioners must take care when considering whose empowerment to maximize. Similarly, while we focused on *maximizing* empowerment, it is trivial to change the sign so that an “assistive” agent minimizes empowerment. One could imagine using such a tool in policies to handicap one’s political opponents.

Acknowledgments. We would like to thank Micah Carroll and Cameron Allen for their helpful feedback, as well as Niklas Lauffer for suggesting JaxMARL. We especially thank the fantastic NeurIPS reviewers for their constructive comments and suggestions. This research was partly supported by ARL DCIST CRA W911NF-17-2-0181 and ONR N00014-22-1-2773, as well as NSF HCC 2310757, the Jump Cocosys Center, Princeton Research Computing, and the DoD through the NDSEG Fellowship Program.

References

- [1] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open Problems and Fundamental Limitations of Reinforcement Learning From Human Feedback. 2023. arXiv:2307.15217.
- [2] Dylan Hadfield-Menell, Stuart J. Russell, Pieter Abbeel, and Anca Dragan. Cooperative Inverse Reinforcement Learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- [3] Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J. Russell, and Anca Dragan. Inverse Reward Design. *Advances in Neural Information Processing Systems*, 30, 2017.
- [4] Micah Carroll, Dylan Hadfield-Menell, Stuart Russell, and Anca Dragan. Estimating and Penalizing Preference Shift in Recommender Systems. In *ACM Conference on Recommender Systems*, RecSys '21, pp. 661–667. 2021.
- [5] Alexander Matt Turner, Logan Smith, Rohin Shah, Andrew Critch, and Prasad Tadepalli. Optimal Policies Tend to Seek Power. In *Neural Information Processing Systems*. 2023.
- [6] Yuqing Du, Stas Tiomkin, Emre Kiciman, Daniel Polani, Pieter Abbeel, and Anca Dragan. Ave: Assistance via Empowerment. In *Advances in Neural Information Processing Systems*, volume 33, pp. 4560–4571. 2020.
- [7] Stuart Russell. Learning Agents for Uncertain Environments (Extended Abstract). In *Annual Conference on Computational Learning Theory*, pp. 101–103. 1998.
- [8] Saurabh Arora and Prashant Doshi. A Survey of Inverse Reinforcement Learning: Challenges, Methods and Progress. *Artificial Intelligence*, 297:103500, 2021.
- [9] Christian Wirth, Riad Akrou, Gerhard Neumann, and Johannes Fürnkranz. A Survey of Preference-based Reinforcement Learning Methods. *Journal of Machine Learning Research*, 18(136):1–46, 2017.
- [10] Micah Carroll, Rohin Shah, Mark K. Ho, Thomas L. Griffiths, Sanjit A. Seshia, Pieter Abbeel, and Anca Dragan. on the Utility of Learning About Humans for Human-ai Coordination. In *Conference on Neural Information Processing Systems*. 2019.
- [11] Alexander Pan, Kush Bhatia, and Jacob Steinhardt. the Effects of Reward Misspecification: Mapping and Mitigating Misaligned Models. In *International Conference on Learning Representations*, arXiv:2201.03544. 2022.
- [12] Jeremy Tien, Jerry Zhi-Yang He, Zackory Erickson, Anca D. Dragan, and Daniel S. Brown. Causal Confusion and Reward Misidentification in Preference-based Reward Learning. In *International Conference on Learning Representations*. 2023.
- [13] Cassidy Laidlaw, Shivam Singhal, and Anca Dragan. Preventing Reward Hacking With Occupancy Measure Regularization. 2024. arXiv:2403.03185.
- [14] Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment – an Introduction. 2013. arXiv:1310.1863.
- [15] Ildefons Magrans de Abril and Ryota Kanai. A Unified Strategy for Implementing Curiosity and Empowerment Driven Reinforcement Learning. 2018. arXiv:1806.06505.
- [16] Sean Chen, Jensen Gao, Siddharth Reddy, Glen Berseth, Anca D. Dragan, and Sergey Levine. ASHA: Assistive Teleoperation via Human-in-the-loop Reinforcement Learning. In *International Conference on Robotics and Automation*. 2022.
- [17] Siddharth Reddy, Sergey Levine, and Anca Dragan. First Contact: Unsupervised Human-machine Co-adaptation via Mutual Information Maximization. *Advances in Neural Information Processing Systems*, 35:31542–31556, 2022.
- [18] Andrew G. Barto. Intrinsic Motivation and Reinforcement Learning. In Gianluca Baldassarre and Marco Mirolli, editors, *Intrinsically Motivated Learning in Natural and Artificial Systems*,

- pp. 17–47. Springer, 2013.
- [19] Arthur Aubret, Laetitia Matignon, and Salima Hassas. A Survey on Intrinsic Motivation in Reinforcement Learning. In *Entropy*. 2023.
 - [20] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity Is All You Need: Learning Skills Without a Reward Function. In *International Conference on Learning Representations (ICLR)*. 2019.
 - [21] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by Random Network Distillation. In *International Conference on Machine Learning*. 2023.
 - [22] Karl Friston. the Free-energy Principle: A Unified Brain Theory? *Nature Reviews Neuroscience*, 11(2):127–138, 2010.
 - [23] Glen Berseth, Daniel Geng, Coline Devin, Nicholas Rhinehart, Chelsea Finn, Dinesh Jayaraman, and Sergey Levine. Smirl: Surprise Minimizing Reinforcement Learning in Unstable Environments. In *International Conference on Learning Representations*. 2019.
 - [24] Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. on Variational Bounds of Mutual Information. In *International Conference on Machine Learning*. 2019.
 - [25] Jongwook Choi, Archit Sharma, Honglak Lee, Sergey Levine, and Shixiang Shane Gu. Variational Empowerment as Representation Learning for Goal-conditioned Reinforcement Learning. In *International Conference on Machine Learning*, pp. 1953–1963. 2021.
 - [26] Shakir Mohamed and Danilo Jimenez Rezende. Variational Information Maximisation for Intrinsically Motivated Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 28. 2015.
 - [27] Tobias Jung, Daniel Polani, and Peter Stone. Empowerment for Continuous Agent—environment Systems. *Adaptive Behavior*, 19(1):16–39, 2011.
 - [28] Michael Laskin, Hao Liu, Xue Bin Peng, Denis Yarats, Aravind Rajeswaran, and Pieter Abbeel. CIC: Contrastive Intrinsic Control for Unsupervised Skill Discovery. 2022. arXiv:2202.00161.
 - [29] Seohong Park, Jongwook Choi, Jaekyeom Kim, Honglak Lee, and Gunhee Kim. Lipschitz-constrained Unsupervised Skill Discovery. In *International Conference on Learning Representations*. 2021.
 - [30] Susanne Still and Doina Precup. an Information-theoretic Approach to Curiosity-driven Reinforcement Learning. *Theory in Biosciences*, 131:139–148, 2012.
 - [31] Nikolay Nikolov, Johannes Kirschner, Felix Berkenkamp, and Andreas Krause. Information-Directed Exploration for Deep Reinforcement Learning. In *International Conference on Learning Representations*. 2019.
 - [32] Erdem Biyik, Dylan P. Losey, Malayandi Palan, Nicholas C. Landolfi, Gleb Shevchuk, and Dorsa Sadigh. Learning Reward Functions From Diverse Sources of Human Feedback: Optimally Integrating Demonstrations and Preferences. In *Int. J. Robotics Res.* 2022.
 - [33] Vivek Myers, Erdem Biyik, Nima Anari, and Dorsa Sadigh. Learning Multimodal Rewards From Rankings. In *Conference on Robot Learning*, pp. 342–352. 2021.
 - [34] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian Active Learning for Classification and Preference Learning. 2011. arXiv:1112.5745.
 - [35] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation Learning With Contrastive Predictive Coding. 2019. arXiv:1807.03748.
 - [36] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. In *International Conference on Machine Learning*. 2020.
 - [37] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised Feature Learning via Non-parametric Instance Discrimination. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3733–3742. 2018.
 - [38] Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Curl: Contrastive Unsupervised Representations for Reinforcement Learning. In *International Conference on Machine Learning*, pp. 5639–5650. 2020.
 - [39] Benjamin Eysenbach, Tianjun Zhang, Sergey Levine, and Russ R. Salakhutdinov. Contrastive Learning as Goal-conditioned Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 35, pp. 35603–35620. 2022.

- [40] Peter Dayan. Improving Generalization for Temporal Difference Learning: the Successor Representation. *Neural Computation*, 5(4):613–624, 1993.
- [41] Ida Momennejad, Evan M Russek, Jin H Cheong, Matthew M Botvinick, Nathaniel Douglass Daw, and Samuel J Gershman. the Successor Representation in Human Reinforcement Learning. *Nature Human Behaviour*, 1(9):680–692, 2017.
- [42] Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment—an Introduction. *Guided Self-organization: Inception*, pp. 67–114, 2014.
- [43] Alexander S Klyubin, Daniel Polani, and Chrystopher L Nehaniv. Empowerment: A Universal Agent-centric Measure of Control. In *2005 IEEE Congress on Evolutionary Computation*, volume 1, pp. 128–135. 2005.
- [44] Siddharth Reddy, Anca D Dragan, and Sergey Levine. Shared Autonomy via Deep Reinforcement Learning. 2018. arXiv:1802.01744.
- [45] Nathan D. Ratliff, J. Andrew Bagnell, and Martin A. Zinkevich. Maximum Margin Planning. In *23rd International Conference on Machine Learning - ICML '06*, pp. 729–736. 2006.
- [46] Pieter Abbeel and Andrew Y. Ng. Apprenticeship Learning via Inverse Reinforcement Learning. In *Twenty-first International Conference on Machine Learning - ICML '04*, p. 1. 2004.
- [47] Thomas M Cover. *Elements of Information Theory*. John Wiley & Sons, 1999.
- [48] Nihat Ay, Jürgen Jost, Hông Vân Lê, and Lorenz Schwachhöfer. *Information Geometry*, volume 64 of *Ergebnisse Der Mathematik Und Ihrer Grenzgebiete 34*. Springer International Publishing, 2017.
- [49] Frank Nielsen. an Elementary Introduction to Information Geometry. *Entropy*, 22(10):1100, 2020.
- [50] Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. the Information Geometry of Unsupervised Reinforcement Learning. In *International Conference on Learning Representations*. 2022.
- [51] Robert G. Gallager. Source Coding With Side Information and Universal Coding. 1979.
- [52] Boris Yakovlevich Ryabko. Coding of a Source With Unknown but Ordered Probabilities. *Problems of Information Transmission*, 15(2):134–138, 1979.
- [53] R. Duncan Luce. *Individual Choice Behavior*. Individual Choice Behavior. John Wiley, 1959.
- [54] Brian D Ziebart, Andrew Maas, J Andrew Bagnell, and Anind K Dey. Maximum Entropy Inverse Reinforcement Learning. In *AAAI Conference on Artificial Intelligence*. 2008.
- [55] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft Actor-critic Algorithms and Applications. 2018. arXiv:1812.05905.
- [56] Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, and Stefano Ermon. Iq-learn: Inverse Soft-Q Learning for Imitation. In *Advances in Neural Information Processing Systems*, volume 34, pp. 4028–4039. 2021.
- [57] Cassidy Laidlaw and Anca Dragan. the Boltzmann Policy Distribution: Accounting for Systematic Suboptimality in Human Models. In *International Conference on Learning Representations*. 2022.
- [58] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. 2017. arXiv:1412.6980.
- [59] Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Gardar Ingvarsson, Timon Willi, Akbir Khan, et al. Jaxmarl: Multi-agent RL Environments in JAX. In *International Conference on Autonomous Agents and Multiagent Systems*. 2024.
- [60] Andrew L. Maas, Awni Y. Hannun, and Andrew Y. Ng. Rectifier Nonlinearities Improve Neural Network Acoustic Models. In *International Conference on Machine Learning*, volume 30, p. 3. 2013.
- [61] Tongzhou Wang and Phillip Isola. Understanding Contrastive Representation Learning Through Alignment and Uniformity on the Hypersphere. In *International Conference on Machine Learning*. 2020.
- [62] Benjamin Eysenbach, Vivek Myers, Ruslan Salakhutdinov, and Sergey Levine. Inference via Interpolation: Contrastive Representations Provably Enable Planning and Inference. In *Advances in Neural Information Processing Systems*. 2024.

A Experimental Details

We ran all our experiments on NVIDIA RTX A6000 GPUs with 48GB of memory within an internal cluster. Each evaluation seed took around 5-10 hours to complete. Our losses (Eqs. 10 and 13) were computed and optimized in JAX with Adam [58]. We used a hardware-accelerated version of the Overcooked environment from the JaxMARL package [59]. The experimental results described in Section 5 were obtained by averaging over 5 seeds for the Overcooked coordination ring layout, 15 for the cramped room layout, and 20 for the obstacle gridworld environment. Specific hyperparameter values can be found in our code, which is available at https://github.com/vivekmyers/empowerment_successor_representations.

A.1 Network Architecture

In the obstacle grid environment, we used a network with 2 convolutional and 2 fully connected layers and SiLU activations. In Overcooked, we adapted the policy architecture from past work [4], using 3 convolutional layers followed by 4 MLP layers with Leaky ReLU activations [60]. We concatenate in $a^{\mathbf{R}}$ and $a^{\mathbf{H}}$ to the state as one-hot encoded channels, i.e. if the action is 5, 6 additional channels will be concatenated to the state with all set to 0s except the 5th channel which is set to 1s.

B Theoretical Analysis of Empowerment

To connect our empowerment objective to reward, we will extend the notation in Section 3.1 to include a distribution over possible tasks the human might be trying to solve, \mathcal{R} , such that each $R \sim \mathcal{R}$ defines a distinct reward function $R : \mathcal{S} \rightarrow \mathbb{R}$. We assume $\pi_{\mathbf{R}}$ tries to maximize the γ -discounted empowerment” of the human, defined as

$$\mathcal{E}_{\gamma}(\pi_{\mathbf{H}}, \pi_{\mathbf{R}}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t I(\mathfrak{s}_t^{\gamma}; \mathfrak{a}_t^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t) \right] \quad (\text{Eq. 6})$$

for

$$\mathfrak{s}_t^{\gamma} \triangleq \{\mathfrak{s}_k \text{ for } k \sim \text{Geom}(1 - \gamma)\}. \quad (14)$$

We additionally define $\bar{\mathfrak{s}}_t$ to be the full history of states up to time t and $\bar{\mathfrak{a}}_t^{\mathbf{H}}$ to be the full history of human actions up to time t ,

$$\begin{aligned} \bar{\mathfrak{s}}_t &= \{\mathfrak{s}_i\}_{i=0}^t, \\ \bar{\mathfrak{a}}_t^{\mathbf{H}} &= \{\mathfrak{a}_i^{\mathbf{H}}\}_{i=0}^t. \end{aligned} \quad (15)$$

Then, $\tilde{\mathfrak{s}}_t$ is the full history of states and past human actions up to time t ,

$$\tilde{\mathfrak{s}}_t = \bar{\mathfrak{s}}_t \cup \bar{\mathfrak{a}}_{t-1}^{\mathbf{H}}. \quad (16)$$

Note that the definition of empowerment in Eq. (6) differs slightly from the original construction Eq. (2)—we condition on the full history of human actions, not just the most recent one. This distinction becomes irrelevant in practice if our MDP maintains history in the state, in which case we can equivalently use \mathfrak{s}_t in place of $\tilde{\mathfrak{s}}_t$.

Meanwhile, for any fixed $\pi_{\mathbf{R}}$ and $\beta > 0$, the human is Boltzmann-rational with respect to the robot’s policy:

$$\pi_{\mathbf{H}}(\mathfrak{a}_t^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t) \propto \exp(\beta Q_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(s_t, \mathfrak{a}_t^{\mathbf{H}})) \quad (17)$$

$$\text{where } Q_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(s_t, \mathfrak{a}_t^{\mathbf{H}}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R(s_{t+k}) \mid s_t, \mathfrak{a}_t^{\mathbf{H}} \right]. \quad (18)$$

Equivalently, we can define the human’s (soft) Q-function and value as

$$\begin{aligned} Q_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(s_t, \mathfrak{a}_t^{\mathbf{H}}) &= R(s_t) + \gamma \mathbb{E} \left[V_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(s_{t+1}) \mid s_t, \mathfrak{a}_t^{\mathbf{H}} \right] \\ \text{for } V_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(s_t) &= \mathbb{E} \left[R(s_t) + \gamma R(s_{t+1}) + \gamma^2 R(s_{t+2}) + \dots \mid s_t, \mathfrak{a}_t^{\mathbf{H}} \right]. \end{aligned} \quad (19)$$

The overall human objective is to maximize the expected soft value:

$$\mathcal{J}_{\pi_{\mathbf{R}}}^{\gamma}(\pi_{\mathbf{H}}) = \mathbb{E}_{R \sim \mathcal{R}} \mathbb{E}_{s_0 \sim p_0} \left[V_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(s_0) \right]. \quad (\text{Eq. 5})$$

Note that this definition of $\pi_{\mathbf{H}}$ depends on R and $\pi_{\mathbf{R}}$ and is bounded $0 \leq \mathcal{J}_{\pi_{\mathbf{R}}}^{\gamma}(\pi_{\mathbf{H}}) \leq 1$. As in the CIRL setting [2], we assume robot is unable to access the true human reward $R : \mathcal{S} \rightarrow \mathbb{R}$. One way to think of the robot’s task is as finding a Nash equilibrium between the objectives Eq. (6) and the human best response in Eq. (17).

For convenience, we will also define a multistep version of $Q_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}$,

$$Q_{R,\gamma}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(s_t, \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R(\mathfrak{s}_{t+k}) \mid s_t, \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \right]. \quad (20)$$

B.1 Connecting Empowerment to Reward

Our approach will be to first relate the empowerment (influence of $\mathbf{a}_t^{\mathbf{H}}$ on \mathfrak{s}_+^{γ}) to the mutual information between $\mathbf{a}_t^{\mathbf{H}}$ and the reward R .

Then, we will connect this quantity to a notion of “advantage” for the human (Eq. 27), which in turn can be related to the expected reward under the human’s policy. In its simplest form, this argument will require an assumption over the reward distribution:

Assumption 3.1 (Skill Coverage). *The rewards $R \sim \mathcal{R}$ are uniformly distributed over the scaled $|\mathcal{S}|$ -simplex $\Delta^{|\mathcal{S}|}$ such that:*

$$\left(R + \frac{1}{|\mathcal{S}|} \right) \left(\frac{1}{1-\gamma} \right) \sim \text{Unif}(\Delta^{|\mathcal{S}|}) = \text{Dirichlet}(\underbrace{1, 1, \dots, 1}_{|\mathcal{S}| \text{ times}}). \quad (7)$$

In other words, Assumption 3.1 says our prior over the human’s reward function is uniform with zero mean. This is not the only prior for which this argument works, but for general \mathcal{R} we will need a correction term to incentivize states that are more likely across the distribution of \mathcal{R} . Another way to view Assumption 3.1 is that the human is trying to execute diverse “skills” $z \sim \text{Unif}(\Delta^{|\mathcal{S}|})$.

We also assume ergodicity (Assumption 3.2). In the special case of an MDP that resets to some distribution with full support over \mathcal{S} , this assumption is automatically satisfied.

Assumption 3.2 (Ergodicity). *For some $\pi_{\mathbf{H}}, \pi_{\mathbf{R}}$, we have*

$$\mathbb{P}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}(\mathfrak{s}_+^{\gamma} = s \mid s_0) > 0 \quad \text{for all } s \in \mathcal{S}, \gamma \in (0, 1). \quad (8)$$

Our main result connects empowerment directly to a (lower bound on) the human’s expected reward.

Theorem 3.1. *Under Assumption 3.1 and Assumption 3.2, for sufficiently large γ and any $\beta > 0$,*

$$\mathcal{E}_{\gamma}(\pi_{\mathbf{H}}, \pi_{\mathbf{R}})^{1/2} \leq (\beta/e) \mathcal{J}_{\pi_{\mathbf{R}}}^{\gamma}(\pi_{\mathbf{H}}). \quad (9)$$

Theorem 3.1 says that for a long enough horizon (i.e., γ close to 1), the robot’s empowerment objective will lower bound the (squared, MaxEnt) human objective.

We make use of the following lemmas in the proof.

Lemma B.1. *For $t \sim \text{Geom}(1 - \gamma)$ and any $K \geq 0$,*

$$\liminf_{\gamma \rightarrow 1} I(\mathfrak{s}_+^{\gamma}; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t) \leq I(R; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t). \quad (21)$$

Proof. For sufficiently large γ , \mathfrak{s}_+^{γ} will approach the stationary distribution of $\mathbb{P}^{\pi_{\mathbf{H}},\pi_{\mathbf{R}}}$ for a fixed $\pi_{\mathbf{H}}, \pi_{\mathbf{R}}$, irrespective of $\tilde{\mathfrak{s}}_t$ and $\mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}}$ from Assumption 3.2. So,

$$\liminf_{\gamma \rightarrow 1} I(\mathfrak{s}_+^{\gamma}; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t) \leq I \left(\lim_{\gamma \rightarrow \infty} \mathfrak{s}_+^{\gamma}; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t \right) \quad (22)$$

Since each $R, \pi_{\mathbf{R}}, \gamma$ defines a human policy $\pi_{\mathbf{H}}$ via Eq. (17), we can express the dependencies as the following Markov chain:

$$\hat{\mathbf{a}}_t \longrightarrow R \longrightarrow \lim_{\gamma \rightarrow 1} \mathfrak{s}_+^{\gamma}. \quad (23)$$

Applying the data processing inequality [47], we get

$$I \left(\lim_{\gamma \rightarrow \infty} \mathfrak{s}_+^{\gamma}; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t \right) \leq I(R; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathfrak{s}}_t), \quad (24)$$

from which Eq. (21) follows. \square

Lemma B.2. Suppose we have k logits, denoted by the map $\alpha : \{1 \dots k\} \rightarrow [0, 1]$. For any $\beta > 0$, we can construct the (softmax) distribution

$$p_\beta(i) \propto \exp(\beta\alpha(i)).$$

Then,

$$\mathcal{H}(p_\beta) \geq \log k - \left(\frac{\beta}{e}\right)^2. \quad (25)$$

Proof. We lower bound the “worst-case” of the RHS, $\alpha = (1, 0, \dots, 0)$:

$$\begin{aligned} \mathcal{H}(p_\beta) &= \frac{(1-n)\log\left(\frac{1}{k+e^\beta-1}\right) - e^\beta \log\left(\frac{e^\beta}{k+e^\beta-1}\right)}{k+e^\beta-1} \\ &= \frac{(k+e^\beta-1)\log(k+e^\beta-1) - e^\beta \log(e^\beta)}{k+e^\beta-1} \\ &= \log(k+e^\beta-1) - \frac{e^\beta \log(e^\beta)}{k+e^\beta-1} \\ &\geq \log k - (\beta/e)^2. \end{aligned} \quad (26)$$

□

Lemma B.3. For any t and $K \geq 0$,

$$I(R; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathbf{s}}_t) \leq \lim_{\gamma \rightarrow 1} \left(\frac{\beta}{e} \mathbb{E} [Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_t, \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}})] \right)^2. \quad (27)$$

Proof. Denote by $\hat{\mathbf{a}}_t^{\mathbf{H}} \dots \hat{\mathbf{a}}_{t+K}^{\mathbf{H}} \sim \text{Unif}(\mathcal{A}^{\mathbf{H}})$ a sequence of K random actions. From Lemma B.2:

$$\begin{aligned} I(R; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathbf{s}}_t) &= \mathcal{H}(\mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathbf{s}}_t) - \mathcal{H}(\mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid R, \tilde{\mathbf{s}}_t) \\ &\leq \log(K|\mathcal{A}|) - \mathcal{H}(\mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid R, \tilde{\mathbf{s}}_t) \\ &\leq \lim_{\gamma \rightarrow 1} \left(\frac{\beta}{e} \mathbb{E} [Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_t, \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}}) - Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_t, \hat{\mathbf{a}}_t^{\mathbf{H}}, \dots, \hat{\mathbf{a}}_{t+K}^{\mathbf{H}})] \right)^2, \end{aligned} \quad (28)$$

where the last inequality follows from Lemma B.2 and $Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(\dots) \leq 1$. We also have

$$0 \leq Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_t, \hat{\mathbf{a}}_t^{\mathbf{H}}, \dots, \hat{\mathbf{a}}_{t+K}^{\mathbf{H}}) \leq Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_t, \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}}) \leq 1, \quad (29)$$

which lets us conclude from Eq. (28) that

$$I(R; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathbf{s}}_t) \leq \left(\frac{\beta}{e} \mathbb{E} [Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_t, \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}})] \right)^2. \quad (\text{Eq. 27})$$

□

We can now prove Theorem 3.1 directly by combining Lemmas B.1 and B.3.

Proof of Theorem 3.1. Simplifying the limit in Eq. (9), we get

$$\begin{aligned} \liminf_{\gamma \rightarrow 1} \mathcal{E}_\gamma(\pi_{\mathbf{H}}, \pi_{\mathbf{R}}) &\leq \liminf_{\gamma \rightarrow 1} \left(\sum_{t=0}^{\infty} \gamma^t I(\mathbf{s}_+^\gamma; \mathbf{a}_t^{\mathbf{H}} \mid \tilde{\mathbf{s}}_t) \right) \\ &\leq \liminf_{\gamma \rightarrow 1} I(\mathbf{s}_+^\gamma; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathbf{s}}_t) && \text{(chain rule)} \\ &\leq I(R; \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}} \mid \tilde{\mathbf{s}}_t) && \text{(Lemma B.1)} \\ &\leq \lim_{\gamma \rightarrow 1} \left(\frac{\beta}{e} \mathbb{E} [Q_{R,\gamma}^{\pi_{\mathbf{H}}, \pi_{\mathbf{R}}}(s_t, \mathbf{a}_t^{\mathbf{H}}, \dots, \mathbf{a}_{t+K}^{\mathbf{H}})] \right)^2 && \text{(Lemma B.3)} \\ &\leq \lim_{\gamma \rightarrow 1} \left(\frac{\beta \mathcal{J}_{\pi_{\mathbf{R}}}^\gamma(\pi_{\mathbf{H}})}{e} \right)^2. \end{aligned} \quad (30)$$

It follows that for sufficiently large γ ,

$$\mathcal{E}_\gamma(\pi_{\mathbf{H}}, \pi_{\mathbf{R}})^{1/2} \leq (\beta/e) \mathcal{J}_{\pi_{\mathbf{R}}}^\gamma(\pi_{\mathbf{H}}). \quad (\text{Eq. 9})$$

□

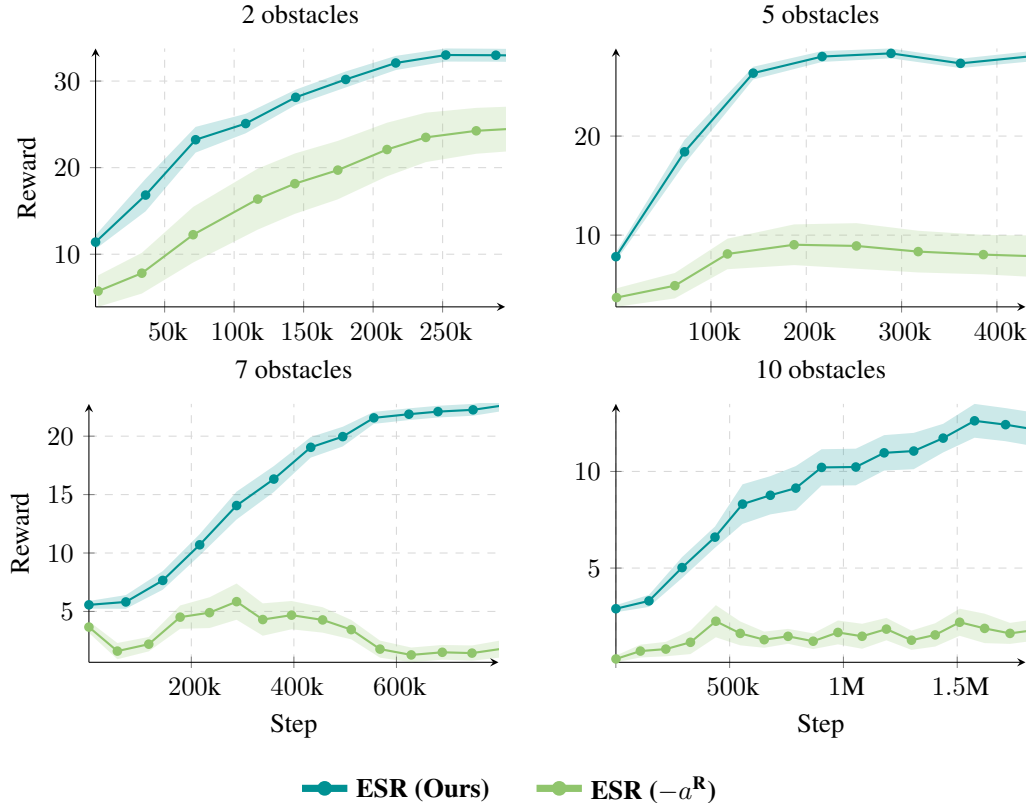


Figure 6: We evaluate our method with and without conditioning on the robot action a^R . Conditioning aids learning significantly, which we theorize is because it removes uncertainty in the classification.

C Additional Ablations and Qualitative Results

In this section we evaluate additional ablations and qualitative results for the ESR method.

C.1 Learning Representations without the Robot Action

In our estimation of empowerment (Eq. 12) we supply the robot action a^R when learning both ϕ and ϕ , however, the theoretical empowerment formulation in Section 3.3 does not require it.

To evaluate the impact of including a^R , we run an additional ablation without it on the gridworld environment, shown in Fig. 6. This ablation shows that the inclusion of a^R is most impactful in more challenging (higher number of boxes) environments. We hypothesize that conditioning the representations on the robot action reduces the noise in the mutual information estimation, and also reduces the difficulty of classifying true future states.

D Greedy Empowerment Policy

All of our experiments have used Soft-Q learning to learn a policy from the empowerment estimation. Here, we additionally study a greedy empowerment policy which takes the most empowering action at each step. We model this by setting the q-learning gamma to 0 to fully discount future rewards.

Results for this ablation are shown in Fig. 7. Unsurprisingly, the greedy optimization vastly underperforms the policy with $\gamma = 0.9$.

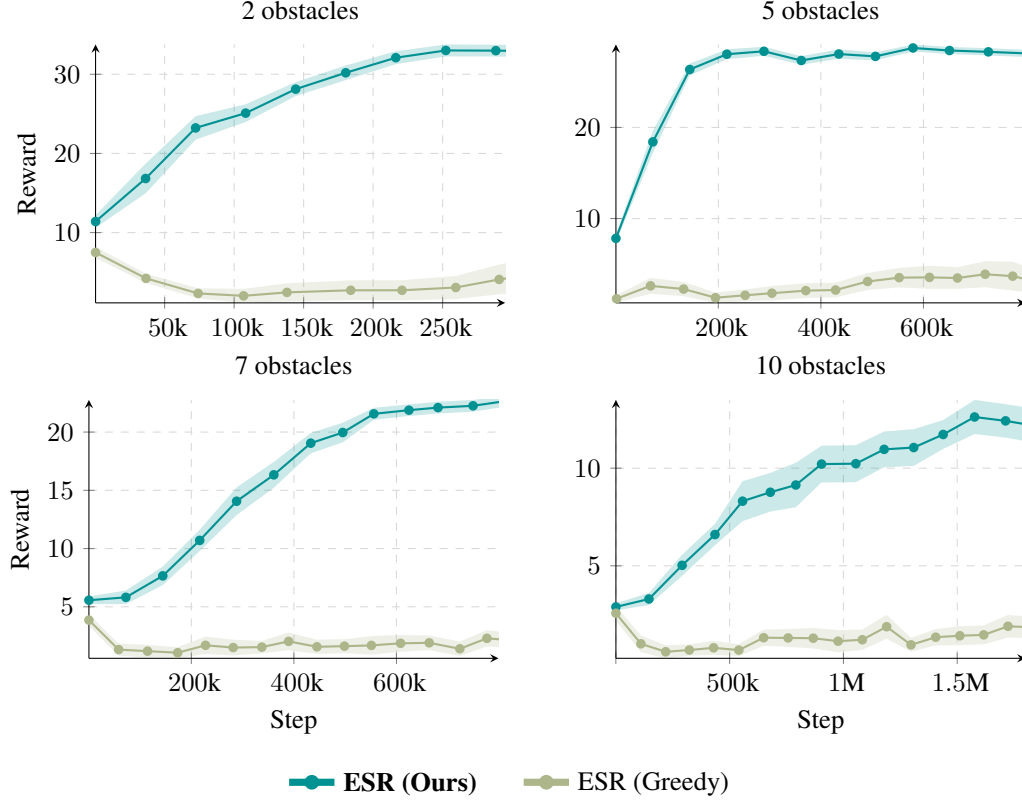


Figure 7: We compare a greedy policy ($\gamma = 0$) against our standard policy ($\gamma = 0.9$).

D.1 ESR Training Example

In Fig. 8, we show the mutual information during training of the ESR agent in the gridworld environment with 5 obstacles. The mutual information quickly becomes positive and remains so throughout training. As long as the mutual information is positive, the classifier is able to reward the agent for taking actions that empower the human.

E Simplifying the Objective

The reward function in Eq. (13) is itself a random variable because it depends on future states g . This subsection describes how this randomness can be removed. To do this, we follow prior work [61, 62] in arguing that the learned representations $\psi(g)$ follow a Gaussian distribution:

Assumption E.1 (Based on Wang and Isola [61]). *The representations of future states $\psi(g)$ learned by contrastive learning have a marginal distribution that is Gaussian:*

$$P(\psi) = \int P(g)\delta(\psi = \psi(g)) dg \stackrel{d}{=} \mathcal{N}(0, I). \quad (31)$$

With this assumption, we can remove the random sampling of g from the reward function. We start by noting that the learned representations tell us the *relative* likelihood of seeing a future state Eq. (12)). Assumption E.1 will allow us to convert these relative likelihoods into likelihoods.

$$\begin{aligned} \mathbb{E}_{P(s^+|s, a^R, a^H)}[r(s, a^R)] &= \mathbb{E}_{P(s^+)} \left[\frac{P(s^+|s, a^R, a^H)}{P(s^+)} r(s, a^R) \right] \\ &= \mathbb{E}_{P(s^+)} \left[C_1 e^{\phi(s, a^R, a^H)^T \phi(s^+)} r(s, a^R) \right] \\ &= C_1 \mathbb{E}_{\psi \sim P(\phi(s^+))} \left[e^{\phi(s, a^R, a^H)^T \psi} (\phi(s, a^R, a^H) - \phi(s, a^R))^T \psi \right] \end{aligned}$$

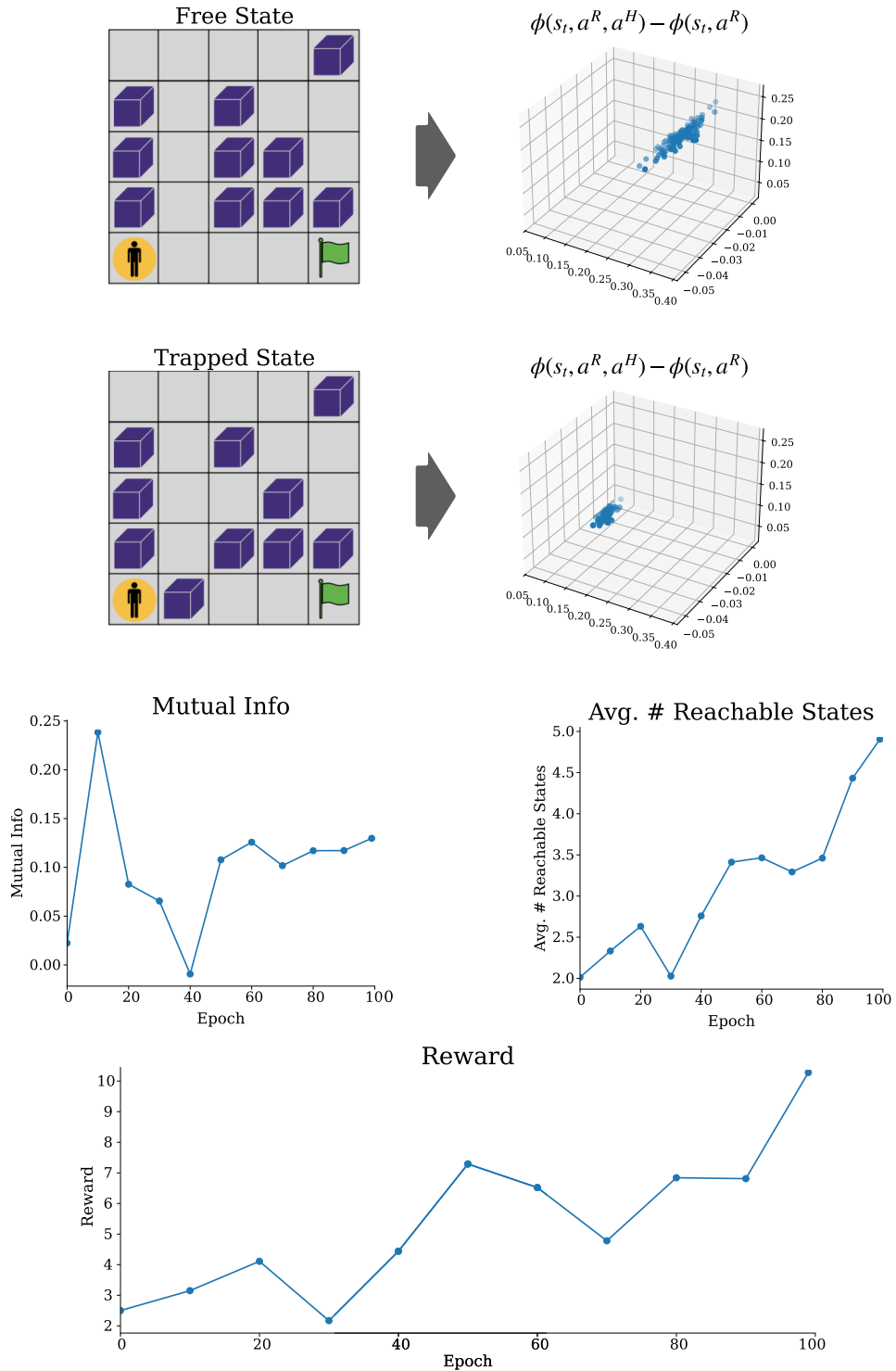


Figure 8: Visualizing training empowerment in a 5x5 Gridworld with 10 obstacles. Our empowerment objective maximizes the influence of the human’s actions on the future state, preferring the state where the human can reach the goal to the trapped state. This corresponds to maximizing the volume of the state marginal polytope, which is proportional to the number of states that the human can reach from their current position. To visualize the representations, we set the latent dimension to 3 instead of 100.

$$\begin{aligned}
&= C_1 (\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}}) - \phi(s, a^{\mathbf{R}}))^T \\
&\quad \int \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2} \|\psi\|_2^2 + \phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})^T \psi} \psi \, d\psi \\
&= C_1 (\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}}) - \phi(s, a^{\mathbf{R}}))^T e^{\frac{1}{2} \|\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})\|_2^2} \\
&\quad \int \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2} \|\psi\|_2^2 + \phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})^T \psi - \frac{1}{2} \|\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})\|_2^2} \psi \, d\psi \\
&= C_1 (\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}}) - \phi(s, a^{\mathbf{R}}))^T \\
&\quad e^{\frac{1}{2} \|\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})\|_2^2} \mathbb{E}_{\psi \sim \mathcal{N}(\mu = \phi(s, a^{\mathbf{R}}, a^{\mathbf{H}}), \Sigma = I)} [\psi] \\
&= C_1 e^{\frac{1}{2} \|\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}})\|_2^2} (\phi(s, a^{\mathbf{R}}, a^{\mathbf{H}}) - \phi(s, a^{\mathbf{R}}))^T \phi(s, a^{\mathbf{R}}, a^{\mathbf{H}}). \quad (32)
\end{aligned}$$

This simplification may be attractive in cases where the computed empowerment bonuses have high variance, or when the empowerment horizon is large (i.e., $\gamma \rightarrow 1$, as in Section 3.3). Empirically, we found this version of the objective to be less effective in practice due to the additional representation structure required by Assumption E.1.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Yes, our goals in the abstract and introduction are well-specified and accurately reflect the rest of the paper and the experimental results.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Yes, we have a limitation subsection that addresses the core limitations of the approach.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Yes, we clearly cite our assumptions and prior work that we are building on for theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Yes, we fully describe how to reproduce our method by detailing the different contrastive networks, the exact loss functions we use, the environments we train on, and the human policies we are training against.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Yes, our code is released at <https://anonymous.4open.science/r/esr-7e94>

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Yes, all of the training details are included in depth to help understand the results, and also the code containing all of this is released.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Yes standard error bars are included and explained, and the number of seeds are given.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Yes, we provide the experimental details including the compute resources needed to reproduce the experiments in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Yes, we comply to the Code of Ethics in every way.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We extensively discuss various safety risks, and how empowerment plays into that, in the Discussion section of the paper.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The data and models are not at risk of misuse, because they are in a gridworld and Overcooked, and policies trained on these environments do not pose a safety issue.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes, we cite all the assets used, including code for environments, experiments, and algorithms, and graphics describing the environments. We additionally credit the authors of the behavioral cloning policies we use as human models, link to their repository/license, and mention that they use an MIT License.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.

- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The code that we release is well documented and the license is provided alongside the code repository.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper involves neither crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.