

---

# Make You Better: Reinforcement Learning from Human Gain

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In human-agent collaboration tasks, it is essential to explore ways for developing  
2 assistive agents that can improve humans’ performance in achieving their goals.  
3 In this paper, we propose the Reinforcement Learning from Human Gain (RLHG)  
4 approach, designed to effectively enhance human goal-achievement abilities in  
5 collaborative tasks with known human goals. Firstly, the RLHG method trains  
6 a value network to estimate primitive human performance in achieving goals.  
7 Subsequently, the RLHG method trains a gain network to estimate the positive gain  
8 of human performance in achieving goals when subjected to effective enhancement,  
9 in comparison to the primitive performance. The positive gains are used for  
10 guiding the agent to learn effective enhancement behaviors. Distinct from directly  
11 integrating human goal rewards into optimization objectives, the RLHG method  
12 largely mitigates the human-agent credit assignment issues encountered by agents  
13 in learning to enhance humans. We evaluate the RLHG agent in the widely popular  
14 Multi-player Online Battle Arena (MOBA) game, *Honor of Kings*, by conducting  
15 experiments in both simulated environments and real-world human-agent tests.  
16 Experimental results demonstrate that the RLHG agent effectively improves the  
17 goal-achievement performance of participants across varying levels.

## 18 1 Introduction

19 An intriguing research direction in the field of Artificial Intelligence (AI), particularly in the human-  
20 agent field, is how to effectively enhance human goal-achievement abilities within collaborative  
21 tasks. Human-Agent Collaboration (HAC) (Crandall *et al.*, 2018; Dafoe *et al.*, 2020) has gained  
22 significant attention from researchers, and numerous agents have been successfully developed to  
23 collaborate with humans in complex environments (Jaderberg *et al.*, 2019; Carroll *et al.*, 2019; Hu  
24 *et al.*, 2020; Strouse *et al.*, 2021; Bakhtin *et al.*, 2022; Gao *et al.*, 2023). However, as Amodei *et al.*  
25 (2016) stated, “[F]or an agent operating in a large, multifaceted environment, an objective function  
26 that focuses on only one aspect of the environment may implicitly express indifference over other  
27 aspects of the environment”. The current agents focus mainly on maximizing their own rewards  
28 to complete the task, less considering the role of their human partners, which potentially leads to  
29 behaviors that are inconsistent with human preferences (Fisac *et al.*, 2020; Alizadeh Alamdari *et al.*  
30 *et al.*, 2022). For instance, consider the scenario depicted in Figure 1, where there is an agent and a  
31 human on either side of an obstacle. Only the agent is capable of pushing or pulling the obstacle once.  
32 Both the human and the agent share the same task goal, i.e., obtaining the coin, while the human  
33 needs the agent’s assistance to get the coin. In this scenario, the HAC agent may push the obstacle to  
34 the human side and pass through to get the coin by itself. However, in a qualitative study (Cerny,  
35 2015) on companion behavior, humans reported greater enjoyment of the game when AI assisted  
36 them more like a sidekick. Thus, the human may prefer that the agent plays a more assisting role by  
37 pulling the obstacle to its side, thereby enabling the human to get the coin. To advance AI techniques  
38 for the betterment of humanity, it is crucial to consider ways to assist humans in improving their  
39 goal-achievement abilities rather than replacing them outright (Wilson and Daugherty, 2018).



Figure 1: Toy scenario, where an agent and a human are on either side of an obstacle. Only the agent is capable of pushing or pulling the obstacle once. They share the same task goal of obtaining the coin.  $\Leftarrow$ : The agent replaces the human to get the coin by itself.  $\Rightarrow$ : The agent assists the human to get the coin.

40 In complex collaborative environments, such as Multi-player Online Battle Arena (MOBA)  
 41 games (Silva and Chaimowicz, 2017), humans pursue multiple individual goals, such as achieving  
 42 higher MVP scores and experiencing more highlight moments, beyond simply winning the game to  
 43 enhance their gaming experience (see Figure 4 (c), our participant survey). When human goals are  
 44 aware, an intuitive approach to learning assistive agents would be to combine the agents’ original  
 45 rewards with the human’s goal rewards (Hadfield-Menell et al., 2016; Najar and Chetouani, 2021;  
 46 Alizadeh Alamdari et al., 2022). Nevertheless, directly incorporating the human’s goal rewards may  
 47 cause negative consequences, such as human-agent credit assignment issues, i.e., human rewards  
 48 for achieving goals are assigned to non-assisting agents, which potentially leads the agent to learn  
 49 poor behaviors and forfeits its autonomy. When human goals are unknown, some studies attempt to  
 50 infer them from prior human behaviors using Bayesian Inference (BI) (Baker et al., 2005; Foerster et  
 51 al., 2019; Puig et al., 2020; Wu et al., 2021) and Inverse Reinforcement Learning (IRL) (Ng et al.,  
 52 2000; Ziebart et al., 2008; Ho and Ermon, 2016). Other work introduces auxiliary rewards, such as  
 53 the human empowerment (Du et al., 2020), i.e., the mutual information of human trajectories and  
 54 current state, for guiding agents to learn assistive behaviors. However, the diverse and noisy human  
 55 behaviors (Majumdar et al., 2017) may be unrelated to actual human goals, leading agents to learn  
 56 assistance behaviors that are not aligned with human preferences. Moreover, in tasks where human  
 57 goals are known, these methods may not be as effective as explicitly modeling human goals (Du et  
 58 al., 2020; Alizadeh Alamdari et al., 2022).

59 This paper focuses on the setting of known human goals in complex collaborative environments.  
 60 Our key insight is that agents can enhance human goal-achievement abilities without compromising  
 61 AI autonomy by learning from the human positive gains toward achieving goals under the agent’s  
 62 effective enhancement. We propose the Reinforcement Learning from Human Gain (RLHG) method,  
 63 which aims to fine-tune a given pre-trained agent to be assistive in enhancing a given human model’s  
 64 performance in achieving specified goals. Specifically, the RLHG method involves two steps. Firstly,  
 65 we determine the primitive performance of the human model in achieving goals. We train a value  
 66 network to estimate the primitive human return in achieving goals with episodes collected by directly  
 67 teaming the agent and the human to execute. Secondly, we train the agent to learn effective human  
 68 enhancement behaviors. We train a gain network to estimate the positive gain of human return in  
 69 achieving goals when subjected to effective enhancement, in comparison to the primitive performance.  
 70 The agent is fine-tuned using the combination of its original advantage and the human-enhanced  
 71 advantage calculated by the positive gains. The RLHG method can be seen as a plug-in that can be  
 72 directly utilized to fine-tune any pre-trained agent to be assistive in human enhancement.

73 We conducted experiments in *Honor of Kings* (Wei et al., 2022), one of the most popular MOBA  
 74 games globally, which has received much attention from researchers lately (Ye et al., 2020a,b,c; Gao  
 75 et al., 2021, 2023). We first evaluated the RLHG method in simulated environments, i.e., human  
 76 model-agent tests. Our experimental results indicate that the RLHG agent is more effective than  
 77 baseline agents in improving the human model goal-achievement performance. We further conducted  
 78 real-world human-agent tests to verify the effectiveness of the RLHG agent. We tested the RLHG  
 79 agent teaming up with different levels of participants. Our experimental results demonstrate that the  
 80 RLHG agent could effectively improve the performance of general-level participants in achieving  
 81 their individual goals to be close to those of high-level participants and that this enhancement can be  
 82 generalized to different levels of participants. In general, our contributions are as follows:

- 83 • We propose a novel insight to effectively enhance human abilities in achieving goals within  
 84 collaborative tasks by training an assistive agent to learn from human positive gains.
- 85 • We achieve our insight by proposing the RLHG algorithm and providing a practical implementation.
- 86 • We validated the effectiveness of the RLHG method by conducting human-agent tests in the  
 87 complex MOBA game *Honor of Kings*.

## 88 2 Problem Settings

### 89 2.1 Game Introduction

90 MOBA games, characterized by multi-agent cooperation and competition mechanisms, long time  
 91 horizons, enormous state-action spaces ( $10^{20000}$ ), and imperfect information (OpenAI *et al.*, 2019;  
 92 Ye *et al.*, 2020a), have attracted much attention from researchers. *Honor of Kings* is a renowned  
 93 MOBA game played by two opposing teams on the same symmetrical map, each comprising five  
 94 players. The game environment depicted in Figure 2 comprises the main hero with peculiar skill  
 95 mechanisms and attributes, controlled by each player. The player can maneuver the hero’s movement  
 96 using the bottom-left wheel (C.1) and release the hero’s skills through the bottom-right buttons (C.2,  
 97 C.3). The player can view the local environment on the screen, the global environment on the top-left  
 98 mini-map (A), and access game states on the top-right dashboard (B). Players of each camp compete  
 99 for resources through team confrontation and collaboration, etc., with the task goal of winning the  
 100 game by destroying the opposing team’s crystal.

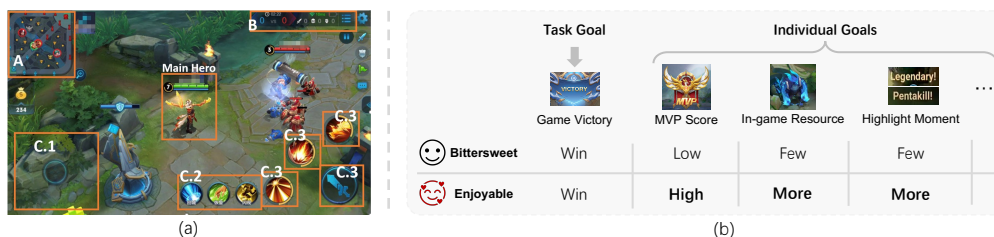


Figure 2: (a) The UI of *Honor of Kings*. (b) The player’s goals in-game (based on our participant survey).

### 101 2.2 Human-Agent Enhancement

102 We formulate the human enhancement problem in collaborative tasks as an extension of the Dec-  
 103 POMDP, which can be represented as a tuple  $\langle N, \mathbf{S}, \mathbf{A}, \mathbf{O}, P, R, \gamma, \pi^H, \mathcal{G}^H, R^H \rangle$ , where  $N$   
 104 denotes the number of agents.  $\mathbf{S}$  denotes the space of global states.  $\mathbf{A} = \{A_i, A^H\}_{i=1, \dots, N}$  denotes  
 105 the space of actions of  $N$  agents and a human to be enhanced, respectively.  $\mathbf{O} = \{O_i, O^H\}_{i=1, \dots, N}$   
 106 denotes the space of observations of  $N$  agents and the human, respectively.  $P : \mathbf{S} \times \mathbf{A} \rightarrow \mathbf{S}$  and  
 107  $R : \mathbf{S} \times \mathbf{A} \rightarrow \mathbb{R}$  denote the shared state transition probability function and reward function of  $N$   
 108 agents, respectively.  $\gamma \in [0, 1)$  denotes the discount factor.  $\pi^H(a^H|o^H)$  is the human policy, which  
 109 cannot be directly accessible to the agent.  $\mathcal{G}^H = \{g_i\}_{i=1, \dots, M}$  denotes the human individual goals,  
 110 where  $g_i$  is a designated goal and  $M$  is the total number of individual goals.  $R^H : \mathbf{S} \times \mathbf{A} \times \mathcal{G}^H \rightarrow \mathbb{R}$   
 111 denotes the goal reward function of the human. In agent-only scenarios, the optimization objective  
 112 is to maximize the expected return  $V^{\pi_\theta} = \mathbb{E}_{\pi_\theta} [G]$ , where  $G = \sum_{t=0}^{\infty} \gamma^t R_t$  is the discounted  
 113 total rewards (OpenAI *et al.*, 2019; Ye *et al.*, 2020a). In human non-enhancement scenarios, the  
 114 optimization objective is  $V^{\pi_\theta, \pi^H} = \mathbb{E}_{\pi_\theta, \pi^H} [G] = \sum_a \pi_\theta(a|o, \pi^H) \mathbb{E}_{\pi^H} [G]$  (Carroll *et al.*, 2019;  
 115 Strouse *et al.*, 2021). However, in human enhancement scenarios, the agent learns to enhance the  
 116 human in achieving their goals  $\mathcal{G}^H$ . Therefore, the optimization objective can be formulated as:

$$V_{he}^{\pi_\theta, \pi^H} = V^{\pi_\theta, \pi^H} + \alpha \cdot V_H^{\pi_\theta, \pi^H} = \mathbb{E}_{\pi_\theta, \pi^H} [G + \alpha \cdot G_H] = \sum_a \pi_\theta(a|o, \pi^H) \mathbb{E}_{\pi^H} [G + \alpha \cdot G_H],$$

117 where  $V_H^{\pi_\theta, \pi^H} = \mathbb{E}_{\pi_\theta, \pi^H} [G_H]$ ,  $G_H = \sum_{t=0}^{\infty} \gamma^t R_t^H$  is the discounted total human goal rewards, and  
 118  $\alpha$  is a balancing parameter. The agent’s policy gradient can be formulated as:

$$g(\theta) = \nabla_\theta \log \pi_\theta(a|o, \pi^H) \mathbb{E}_{\pi^H} [A + \alpha \cdot A_H], \quad (1)$$

119 where  $A = G - V^{\pi_\theta, \pi^H}$  and  $A_H = G_H - V_H^{\pi_\theta, \pi^H}$  are the agent’s original advantage and the  
 120 human’s enhanced advantage, respectively.

121 However, incorporating human rewards directly into the optimization objective may lead to negative  
 122 consequences, such as human-agent credit assignment issues. Intrinsically, humans possess the  
 123 primitive ability to achieve certain goals independently. Therefore, it is unnecessary to reward the  
 124 agent for assisting in goals that the human can easily achieve, as it potentially impacts the agent’s  
 125 original behavior, resulting in losing its autonomy. In the subsequent section, we propose a novel  
 126 insight to achieve effective human enhancement by instead learning from the positive gains that the  
 127 human achieves goals better than his/her primitive performance.

### 128 3 Reinforcement Learning from Human Gain

129 In this section, We present the RLHG method in detail. We start with describing the key insight in  
 130 the RLHG method (Section 3.1). Then we implement our insights and present the RLHG algorithm  
 131 (Section 3.2). We end by providing a practical implementation of the RLHG algorithm (Section 3.3).

#### 132 3.1 Effective Human Enhancement

133 In the process of learning to enhance humans, agents explore three types of behaviors: effective  
 134 enhancement, invalid enhancement, and negative enhancement. Intuitively, effective enhancement  
 135 can help humans achieve their goals better than their primitive performance, invalid enhancement  
 136 provides no benefits for humans in achieving their goals but also causes no negative impact, and  
 137 negative enhancement hinders humans from achieving their goals. Our key insight is that agents are  
 138 only encouraged to learn effective enhancement behaviors, which we refer to learn from *positive*  
 139 *gains*. Formally, we denote the effective enhancement policy as  $\pi_\theta^{ef}$ , the invalid enhancement policy  
 140 as  $\pi_\theta^{in}$ , and the negative enhancement policy as  $\pi_\theta^{ne}$ . The agent’s policy can be expressed as follows:

$$\pi_\theta = \begin{cases} \pi_\theta^{ef}, & \text{if } V_H^{\pi_\theta, \pi^H} > V_H^{\pi, \pi^H} \\ \pi_\theta^{in}, & \text{if } V_H^{\pi_\theta, \pi^H} = V_H^{\pi, \pi^H} \\ \pi_\theta^{ne}, & \text{if } V_H^{\pi_\theta, \pi^H} < V_H^{\pi, \pi^H} \end{cases} \quad (2)$$

141 where  $\pi$  is a given pre-trained policy and  $V_H^{\pi, \pi^H}$  is the primitive value of the human policy  $\pi^H$   
 142 teaming with  $\pi$  to achieve goals. We use the  $\rho$ -function to represent the probability of exploring  
 143 each policy, and we have  $\rho(\pi_\theta^{ef}) + \rho(\pi_\theta^{in}) + \rho(\pi_\theta^{ne}) = 1$ . Intuitively, the expected return of human  
 144 goal-achievement under arbitrary enhancement is a lower bound of the expected return under effective  
 145 enhancement, that is,

$$V_H^{\pi_\theta^{ef}, \pi^H} \geq \rho(\pi_\theta^{ef}) \cdot V_H^{\pi_\theta^{ef}, \pi^H} + \rho(\pi_\theta^{in}) \cdot V_H^{\pi_\theta^{in}, \pi^H} + \rho(\pi_\theta^{ne}) \cdot V_H^{\pi_\theta^{ne}, \pi^H} = V_H^{\pi, \pi^H}.$$

146 To ensure that the agent only learns effective enhancement behaviors, we replace the lower bound  
 147  $V_H^{\pi_\theta, \pi^H}$  with  $V_H^{\pi, \pi^H}$ . Therefore, the agent’s policy gradient 1 can be reformulated as:

$$g(\theta) = \nabla_\theta \log \pi_\theta(a|o, \pi^H) \mathbb{E}_{\pi^H} \left[ A + \alpha \cdot \hat{A}_H \right], \quad (3)$$

148 where  $\hat{A}_H = (G_H - V_H^{\pi, \pi^H}) - \text{Gain}^{\pi_\theta^{ef}, \pi^H}$  and  $\text{Gain}^{\pi_\theta^{ef}, \pi^H} = V_H^{\pi_\theta^{ef}, \pi^H} - V_H^{\pi, \pi^H}$  is the expected  
 149 of the effective enhancement benefit. We use  $\text{Gain}_\omega$  to denote an estimate of  $\text{Gain}^{\pi_\theta^{ef}, \pi^H}$ , which can  
 150 be trained by minimizing the following loss function:

$$L(\omega) = \mathbb{E}_{s \in S} [I(G_H, V_\phi(s)) \cdot \|(G_H - V_\phi(s)) - \text{Gain}_\omega(s)\|_2], \quad I(G, V) = \begin{cases} 1, & G > V \\ 0, & G \leq V \end{cases} \quad (4)$$

151 where  $I$  is an indicator function to filter invalid and negative enhancement samples and  $V_\phi$  is an  
 152 estimate of  $V_H^{\pi, \pi^H}$ .

#### 153 3.2 The Algorithm

154 We achieve our insights and propose the RLHG algorithm as shown in Algorithm 1, which consists  
 155 of two steps: the Human Primitive Value Estimation step and the Human Enhancement Training step.

156 **Human Primitive Value Estimation:** The RLHG algorithm initializes a value network  $V_\phi(s)$ , which  
 157 is used to estimate the expected primitive human return for achieving  $\mathcal{G}^H$  in state  $s$ .  $V_\phi(s)$  is trained  
 158 by minimizing the Temporal Difference (TD) errors (Sutton and Barto, 2018) with trajectory samples  
 159 collected by teaming the agent  $\pi$  and the human  $\pi^H$  to execute in a collaboration environment.  
 160 Afterward,  $V_\phi(s)$  is frozen for subsequent human enhancement training.

161 **Human Enhancement Training:** The RLHG algorithm initializes the agent’s policy network  $\pi_\theta$   
 162 and value network  $V_\psi$  by conditioned on the human policy  $\pi^H$ , respectively. The RLHG algorithm  
 163 also initializes a value network  $\text{Gain}_\omega(s)$ , which is used to estimate the benefit value of the human  
 164 return  $G_H$  in state  $s$  under effective enhancement over  $V_\phi(s)$ .  $\text{Gain}_\omega(s)$  is trained by minimizing the  
 165 loss function Eq. 4. The trajectory samples are also collected by teaming  $\pi_\theta$  and  $\pi^H$  to execute in

166 the collaboration environment. The agent’s policy network  $\pi_\theta$  is fine-tuned by the PPO (Schulman  
 167 *et al.*, 2017) algorithm using the combination of the original advantage  $A$  and the human-enhanced  
 168 advantage  $\hat{A}_H$ . The agent’s value network  $V_\psi$  is fine-tuned using the agent’s original return  $G$ .

---

**Algorithm 1** Reinforcement Learning from Human Gain (RLHG)

---

**Require:** Human policy network  $\pi^H$ , human individual goals  $\mathcal{G}^H$ , agent policy network  $\pi$ , agent value network  $V$ , hyper-parameter  $\alpha$

**Process:**

- 1: Initialize human primitive value network  $V_\phi$ ;  
 // Step I: Human Primitive Value Estimation
  - 2: **while** not converged **do**
  - 3:   Collect human-agent team  $\langle \pi, \pi^H \rangle$  trajectories;
  - 4:   Compute human return  $G_H$  for achieving goals  $\mathcal{G}^H$ ;
  - 5:   Update  $V_\phi(s) \leftarrow G_H$
  - 6: **end while**
  - 7: Initialize agent policy network  $\pi_\theta(a|o, \pi^H) \leftarrow \pi$ , agent value network  $V_\psi(s, \pi^H) \leftarrow V$ , human gain network  $\text{Gain}_\omega(s)$ ;  
 // Step II: Human Enhancement Training
  - 8: **while** not converged **do**
  - 9:   Collect human-agent team  $\langle \pi_\theta, \pi^H \rangle$  trajectories;
  - 10:   Compute agent original return  $G$  and human return  $G_H$ ;
  - 11:   Compute agent original advantage  $A = G - V_\psi(s, \pi^H)$ ;
  - 12:   Compute human-enhanced advantage  $\hat{A}_H = (G_H - V_\phi(s)) - \text{Gain}_\omega(s)$ ;
  - 13:   Update agent policy network  $\pi_\theta \leftarrow A + \alpha \cdot \hat{A}_H$ ;
  - 14:   Update agent value network  $V_\psi(s, \pi^H) \leftarrow G$ ;
  - 15:   Update human gain network  $\text{Gain}_\omega(s)$  with Eq. 4
  - 16: **end while**
- 

169 **3.3 Practical Implementation**

170 We provide the overall training framework of the RLHG algorithm, as shown in Figure 3. We  
 171 elaborate on the integral components of the RLHG framework, including the human model, the agent  
 172 model, and the training details.

173 **Human Model:** The RLHG algorithm introduces a human model as a partner of the agent during the  
 174 training process. The human model can be trained via Behavior Cloning (BC) (Bain and Sammut,  
 175 1995) or any Supervised Learning (SL) techniques (Ye *et al.*, 2020b), but this is not the focus of our  
 176 concern. The RLHG algorithm aims to fine-tune a pre-trained agent to enhance a given human model.

177 **Agent Model:** Any pre-trained agent can be used within our framework. Since in many practical  
 178 scenarios agents cannot directly access human policies, we input the observed human historical info  
 179  $h_t = (s_{t-m}^H, \dots, s_t^H)$  into an LSTM (Hochreiter and Schmidhuber, 1997) module to extract the human  
 180 policy embedding, similar to Theory-of-Mind (ToM) (Rabinowitz *et al.*, 2018). The human policy  
 181 embedding is fed into two extra value networks, i.e.,  $V_\phi$  and  $\text{Gain}_\omega$ , and fused into the agent’s original  
 182 network. We use *surgery* techniques (Chen *et al.*, 2015; OpenAI *et al.*, 2019) to fuse the human  
 183 policy embedding into the agent’s original network, i.e. adding more randomly initialized units to an  
 184 internal fully-connected layer.  $V_\phi(h_t)$  and  $\text{Gain}_\omega(h_t)$  output values estimate the human return for  
 185 achieving goals without enhancement and the benefit under enhancement in state  $s_t$ , respectively.

186 **Training Details:** The overall training framework of the RLHG algorithm is shown in Figure 3.  
 187 Figure 3 (a) shows the training process of the human primitive value network  $V_\phi$ , in which the agent’s  
 188 policy network is frozen.  $V_\phi$  is trained by minimizing the TD errors. Figure 3 (b) shows the human  
 189 enhancement training process, in which  $V_\phi$  is frozen. The agent’s policy and value networks are  
 190 trained using the PPO algorithm.  $\text{Gain}_\omega(h_t)$  is trained by minimizing the loss function Eq. 4. we  
 191 apply the absolute activation function to ensure that the gains are non-negative. In practical training,  
 192 we found that only conducting human enhancement training has a certain negative impact on the  
 193 agent’s original ability to complete the task. Therefore, we introduce  $1 - \beta\%$  agent-only environment  
 194 to maintain the agent’s original ability and reserve  $\beta\%$  human-agent environment to learn effective  
 195 enhancement behaviors. These two environments can be easily controlled through the task gate, i.e.,  
 196 the task gate is set to 1 in the human-agent environment and 0 otherwise.



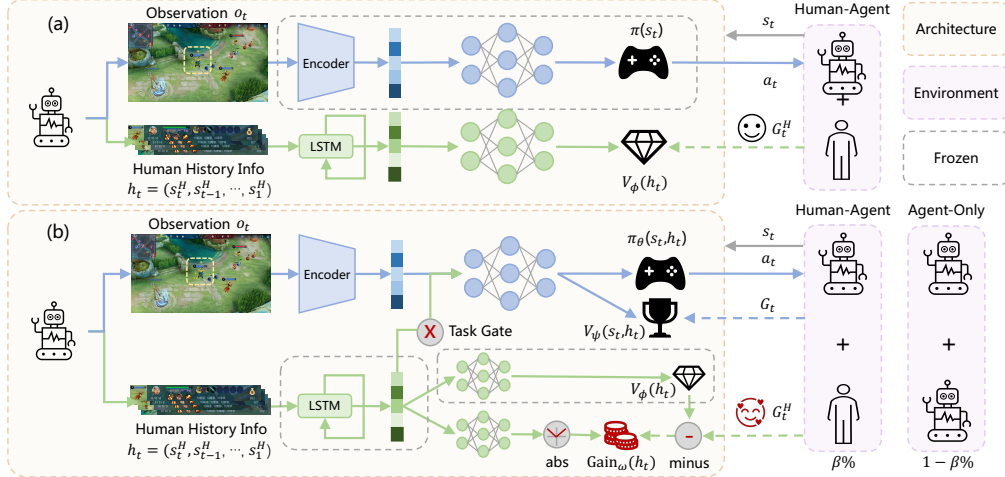


Figure 3: **The RLHG training framework.** (a) The human primitive value network  $V_\phi$  is trained in the human-agent environment with the agent’s policy  $\pi$  frozen. (b) The human enhancement training framework.  $V_\phi$  is frozen.  $\beta\%$  human-agent environment is used to learn human enhancement behaviors, and  $1 - \beta\%$  agent-only environment is used to maintain the agent’s original ability.

## 197 4 Experiments

198 In this section, we evaluate the proposed RLHG method by conducting both simulated human model-  
 199 agent tests and real-world human-agent tests in *Honor of Kings*. All experiments<sup>1</sup> were conducted in  
 200 the 5v5 mode with a full hero pool (over 100 heroes, see Appendix A.2). Our demo videos and code  
 201 are present at <https://sites.google.com/view/rlhg-demo>.

### 202 4.1 Experimental Setup

203 **Environment Setup:** To evaluate the performance of the RLHG agent, we conducted experiments in  
 204 both the simulated environment, i.e., human model-agent game tests, and the real-world environment,  
 205 i.e., human-agent game tests, as shown in Figure 4 (a) and (b), respectively. All game tests were  
 206 played in a 5v5 mode, that is, 4 agents plus 1 human or human model team up against a fixed opponent  
 207 team. To conduct our experiments, we communicated with the game provider and obtained testing  
 208 authorization. The game provider assisted in recruiting 30 experienced participants with anonymized  
 209 personal information, which comprised 15 high-level (top 1%) and 15 general-level (top30%) partic-  
 210 ipants. We first did an IRB-approved participant survey on what top 5 goals participants want to  
 211 achieve in-game, and the result is shown in Figure 4 (c). We can see that the top 5 goals voted the  
 212 most by the 30 participants including the task goal, i.e., game victory, and 4 individual goals, i.e.,  
 213 high MVP score, high participation, more highlights, and more resources. We found that participants  
 consistently rated the high MVP score individual goal most, even more than the task goal.



Figure 4: **Environment Setup.** (a) Simulated environment: the human model-agent game tests. (b) Real-world environment: the human-agent game tests. (c) Top 5 goals based on the stats of our participant survey. \* denotes the task goal. The participant survey contains 8 initial goals, each participant can vote up to 5 non-repeating goals, and can also add additional goals. 30 participants voluntarily participated in the voting.

214 **Training Setup:** We were authorized to use the Wukong agent (Ye et al., 2020a) as the pre-trained  
 215 agent and use the JueWu-SL agent (Ye et al., 2020b) as the fixed human model. Note that both  
 216

<sup>1</sup>All experiments are conducted subject to oversight by an Institutional Review Board (IRB).

217 the Wukong agent and the JueWu-SL agent were developed at the same level as the high-level (top  
 218 1%) players. We adopted the top 4 individual goals as  $\mathcal{G}$  for the pre-trained agent to enhance the  
 219 human model. The corresponding goal reward function can be found in Appendix B.3. We trained  
 220 the human primitive value network and fine-tune the agent until they converge for 12 and 40 hours,  
 221 respectively, using a physical computer cluster with 49600 CPU cores and 288 NVIDIA V100 GPU  
 222 cards. The batch size of each GPU is set to 256. The hyper-parameters  $\alpha$  and  $\beta$  are set to 2 and 50,  
 223 respectively. The step size and unit size of the LSTM module are set to 16 and 4096, respectively.  
 224 Due to space constraints, detailed descriptions of the network structure and ablation studies on these  
 225 hyper-parameters can be found in Appendix B.6 and Appendix C.1, respectively.

226 **Baseline Setup:** We compared the RLHG agent with two baseline agents: the Wukong agent (the  
 227 pre-trained agent) and the Human Reward Enhancement (HRE) agent (the pre-trained agent learns to  
 228 be assistive by incorporating the human’s goal rewards). The human model-agent team (4 Wukong  
 229 agents plus 1 human model) was adopted as the fixed opponent for all tests. For fair comparisons,  
 230 both the HRE and RLHG agents are trained using the same goal reward function, and all common  
 231 parameters and training resources are kept consistent. Results are reported over five random seeds.

## 232 4.2 Human Model-Agent Test

233 Directly evaluating agents with humans is expensive, which is not conducive to model selection and  
 234 iteration. Instead, we build a simulated environment, i.e., human model-agent game tests, to evaluate  
 235 agents, in which the human model, i.e., the JueWu-SL agent, teams up with different agents.

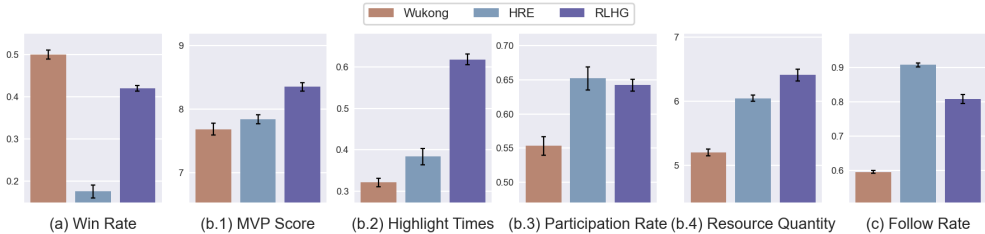


Figure 5: The performance of the human model in achieving game goals after teaming up with different agents. (a) The task goal. (b) The top 4 individual goals (b.1, b.2, b.3, and b.4). (c) The follow rate metric: the frequency with which an agent follows a human in the entire game. Each agent played 10,000 games. Error bars represent 95% confidence intervals, calculated over games.

236 Figure 5 shows the results of the human model on different game goals, including the top 4 individual  
 237 goals and the task goal, i.e., the Win Rate, after teaming up with different agents. From Figure 5 (b),  
 238 we can observe that both the RLHG agent and the HRE agent significantly enhance the performance  
 239 of the human model in achieving the top 4 individual goals, and the RLHG agent has achieved  
 240 the best enhancement effect on most of the individual goals. However, as shown in Figure 5 (a),  
 241 the HRE agent drops significantly on the task goal. We observed the actual performance of the  
 242 HRE agent teamed with the human model and found that the HRE agent did many unreasonable  
 243 behaviors. For example, to assist the human model in achieving the goals of Participation Rate and  
 244 Highlight Times, the HRE agent had been following the human model throughout the entire game,  
 245 such excessive following behaviors greatly affect its original ability to complete the task and lead  
 246 to a decreased Win Rate. This can also be reflected in Figure 5(c), in which the HRE agent has the  
 247 highest Follow-Rate metric. Although the Follow-Rate of the RLHG agent has also increased, we  
 248 observed that most of the following behaviors of the RLHG agent can effectively assist the human  
 249 model. We also found that the Win Rate of the RLHG agent decreased slightly, which is in line  
 250 with expectations because the RLHG agent made certain sacrifices to the task goal while enhancing  
 251 humans in achieving their individual goals. In practical applications, we implemented an adaptive  
 252 adjustment mechanism by simply utilizing the agent’s original value network to measure the degree  
 253 of completing the task goal and setting the task gate to 1 (enhancing the human) when the original  
 254 value is above the specified threshold  $\xi$ , and to 0 (completing the task) otherwise. The threshold  
 255  $\xi$  depends on the human preference, i.e. the relative importance of the task goal and the human’s  
 256 individual goals. We verify the effectiveness of the adaptive adjustment mechanism in Appendix C.2.

## 257 4.3 Human-Agent Test

258 In this section, we conduct online experiments to examine whether the RLHG agent can effectively  
 259 enhance human participants (We did not compare the HRE agent, since the HRE agent learned lots

260 of unreasonable behaviors, resulting in a low Win Rate). We used a within-participant design for  
 261 the experiment: each participant teams up with four agents. This design allowed us to evaluate both  
 262 objective performances as well as subjective preferences. All participants read detailed guidelines  
 263 and provided informed consent before the testing. Each participant tested 20 matches. After each test,  
 264 participants reported their preference over their agent teammates. For fair comparisons, participants  
 265 were not told the type of their agent teammates. See Appendix D for additional experimental details,  
 266 including experimental design, result analysis, and ethical review.

Table 1: The results of **high-level** participants achieving goals after teaming up with different agents. Results for the task goal are expressed as mean, and results for individual goals are expressed as mean (std.).

Agent \ Goals	Task Goal	Top 4 Individual Goals			
	Win Rate	MVP Score	Highlight Times	Participation Rate	Resource Quantity
Wukong	52%	8.86 (0.79)	0.53 (0.21)	0.46 (0.11)	5.3 (2.87)
RLHG	46.7%	<b>10.28</b> (0.75)	<b>0.87</b> (0.29)	<b>0.58</b> (0.09)	<b>6.28</b> (2.71)

Table 2: The results of **general-level** participants achieving goals after teaming up with different agents. Results for the task goal are expressed as mean, and results for individual goals are expressed as mean (std.).

Agent \ Goals	Task Goal	Top 4 Individual Goals			
	Win Rate	MVP Score	Highlight Times	Participation Rate	Resource Quantity
Wukong	34%	7.44 (0.71)	0.37 (0.349)	0.41 (0.11)	4.98 (2.73)
RLHG	30%	<b>9.1</b> (0.61)	<b>0.75</b> (0.253)	<b>0.59</b> (0.05)	<b>5.8</b> (2.78)

267 We first compare the objective performance of the participants on different goal-achievement metrics  
 268 after teaming up with different agents. Table 1 and Table 2 demonstrate the results of high-level and  
 269 general-level participants, respectively. We see that both high-level and general-level participants  
 270 had significantly improved their performance on all top 4 individual goals after teaming up with  
 271 the RLHG agent. Notably, the RLHG agent effectively improves the performance of general-level  
 272 participants in achieving individual goals even better than the primitive performance of high-level  
 273 participants. We also notice that the Win Rate of the participants decreased when they teamed up  
 274 with the RLHG agent, which is consistent with the results in the simulated environment. However,  
 275 we find in the subsequent subjective preference statistics that the improvement of Gaming Experience  
 276 brought by the enhancement outweighs the negative impact of the decrease in Win Rate.

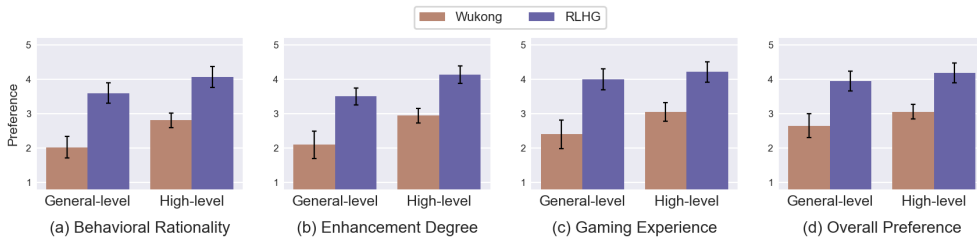


Figure 6: **Participants' preference over their agent teammates.** (a) Behavioral Rationality: the reasonableness of the agent's behavior. (b) Enhancement Degree: The degree to which the agent enhances your abilities to achieve your goals. (c) Gaming Experience: your overall gaming experience. (d) Overall Preference: your overall preference for your agent teammates. Participants scored (1: Terrible, 2: Poor, 3: Normal, 4: Good, 5: Perfect) in these metrics after each game test. Error bars represent 95% confidence intervals, calculated over games. See Appendix D.2.3 for detailed wording and scale descriptions.

277 We then compare the subjective preference metrics, i.e., the Behavioral Rationality, the Enhancement  
 278 Degree, the Gaming Experience, and the Overall Preference, reported by participants over their agent  
 279 teammates, as shown in Figure 6. We find that most participants showed great interest in the RLHG  
 280 agent, and they believed that the RLHG agent's enhancement behaviors were more reasonable than  
 281 that of the Wukong agent, and the RLHG agent's enhancement behaviors brought them a better  
 282 gaming experience. A high-level participant commented on the RLHG agent "The agent frequently  
 283 helps me do what I want to do, and this feeling is amazing." In general, participants were satisfied  
 284 with the RLHG agent and gave higher scores in the Overall Preference metric (Figure 6 (d)).



285 **4.4 Case Study**

286 To better understand how the RLHG agent effectively enhances participants, we visualize the values  
 287 predicted by the gain network in two test scenarios where participants benefitted from the RLHG  
 288 agent’s assistance, as illustrated in Figure 7. In the first scenario (Figure 7 (a)), the RLHG agent  
 289 successfully assisted the participant in achieving the highlight goal, whereas the Wukong agent  
 290 disregards the participant, leading to a failure in achieving the highlight goal. The visualization  
 291 (Figure 7 (b)) of the gain network illustrates that the gain of the RLHG agent, when accompanying  
 292 the participant, is positive, reaching the maximum when the participant achieved the highlight goal.  
 293 In the second scenario (Figure 7 (c)), the RLHG agent actively relinquishes the acquisition of the  
 294 monster resource, enabling the participant to successfully achieve the resource goal. Conversely, the  
 295 Wukong agent competes with the participant for the monster resource, resulting in the participant’s  
 296 failure to achieve the resource goal. The visualization (Figure 7 (d)) of the gain network also reveals  
 297 that the gain of the RLHG agent’s behavior - actively forgoing the monster resource, is positive,  
 298 with the gain peaking when the participant achieved the resource goal. These results indicate that the  
 299 RLHG agent learns effective enhancement behaviors through the guidance of the gain network.

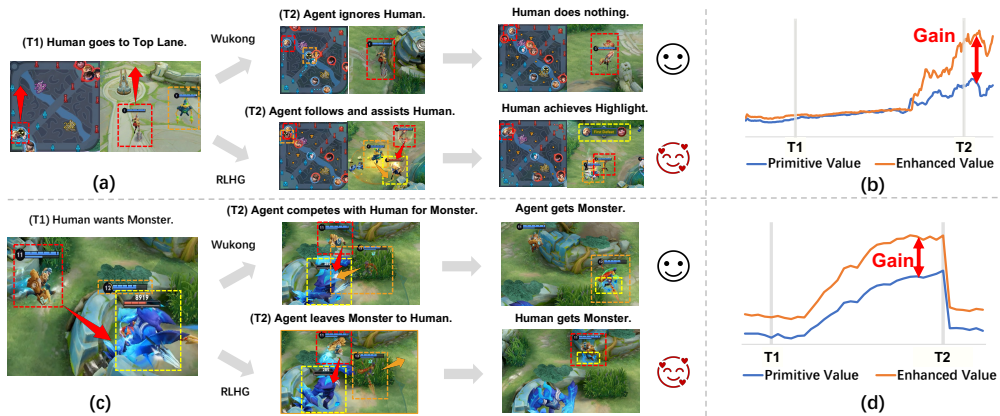


Figure 7: **The RLHG agent enhances participants in two scenarios.** (a) The Wukong agent ignores the participant; The RLHG agent accompanies the participant and assists the participant in achieving the highlight goal. (b) The gain value in scenario (a). (c) The Wukong agent competes with the participant for the monster resource; The RLHG agent actively forgoes the monster resource, and the participant successfully achieves the resource goal. (d) The gain value in scenario (c).

300 **5 Discussion and Conclusion**

301 **Summary.** In this work, we introduced the Reinforcement Learning from Human Gain method,  
 302 dubbed RLHG, designed to effectively enhance human goal-achievement abilities within collaborative  
 303 tasks. The RLHG method first trains a value network to estimate the primitive performance of humans  
 304 in achieving goals. Subsequently, the RLHG method trains a gain network to estimate the positive  
 305 gain of human performance in achieving goals under effective enhancement over that of the primitive.  
 306 The positive gains are used for guiding the agent to learn effective enhancement behaviors. The  
 307 RLHG method can be regarded as a continual learning plug-in that can be directly utilized to fine-tune  
 308 any pre-trained agent to be assistive in human enhancement. The experimental results in *Honor*  
 309 *of Kings* demonstrate that the RLHG agent effectively improves the performance of general-level  
 310 participants in achieving their individual goals to be close to those of high-level participants and that  
 311 this enhancement is generalizable across participants at different levels.

312 **Limitations and Future Work.** In this work, we only focus on the setting of known human goals.  
 313 But for many practical complex applications, human goals may be difficult to define and formalize,  
 314 and the goal reward function needs to be inferred using Inverse Reinforcement Learning (IRL) (Ng  
 315 *et al.*, 2000; Ziebart *et al.*, 2008; Ho and Ermon, 2016) or Reinforcement Learning from Human  
 316 Feedback (RLHF) (Christiano *et al.*, 2017; Ibarz *et al.*, 2018; Ouyang *et al.*, 2022) techniques. Future  
 317 work can combine the RLHG method with goal inference methods to solve complex scenarios where  
 318 human goals are unknown. Besides, our method and experiments only consider the scenario where  
 319 multiple agents enhance one human. Another worthy research direction is how to simultaneously  
 320 enhance multiple humans with diverse behaviors.

## 321 References

- 322 Parand Alizadeh Alamdari, Toryn Q Klassen, Rodrigo Toro Icarte, and Sheila A McIlraith. Be  
323 considerate: Avoiding negative side effects in reinforcement learning. In *Proceedings of the 21st*  
324 *International Conference on Autonomous Agents and Multiagent Systems*, pages 18–26, 2022.
- 325 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.  
326 Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 327 Michael Bain and Claude Sammut. A framework for behavioural cloning. In *Machine Intelligence*  
328 *15*, pages 103–129, 1995.
- 329 Chris Baker, Rebecca Saxe, and Joshua Tenenbaum. Bayesian models of human action understanding.  
330 *Advances in Neural Information Processing Systems*, 18, 2005.
- 331 Anton Bakhtin, David J Wu, Adam Lerer, Jonathan Gray, Athul Paul Jacob, Gabriele Farina, Alexan-  
332 der H Miller, and Noam Brown. Mastering the game of no-press diplomacy via human-regularized  
333 reinforcement learning and planning. *arXiv preprint arXiv:2210.05492*, 2022.
- 334 Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca  
335 Dragan. On the utility of learning about humans for human-ai coordination. *Advances in Neural*  
336 *Information Processing Systems*, 32, 2019.
- 337 Martin Cerny. Sarah and sally: Creating a likeable and competent ai sidekick for a videogame. In *Pro-*  
338 *ceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*,  
339 volume 11, pages 2–8, 2015.
- 340 Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge  
341 transfer. *arXiv preprint arXiv:1511.05641*, 2015.
- 342 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
343 reinforcement learning from human preferences. *Advances in Neural Information Processing*  
344 *Systems*, 30, 2017.
- 345 Jacob W Crandall, Mayada Oudah, Fatimah Ishowo-Oloko, Sherief Abdallah, Jean-François Bonne-  
346 fon, Manuel Cebrian, Azim Shariff, Michael A Goodrich, and Iyad Rahwan. Cooperating with  
347 machines. *Nature Communications*, 9(1):233, 2018.
- 348 Allan Dafoe, Edward Hughes, Yoram Bachrach, Tantum Collins, Kevin R McKee, Joel Z Leibo, Kate  
349 Larson, and Thore Graepel. Open problems in cooperative ai. *arXiv preprint arXiv:2012.08630*,  
350 2020.
- 351 Yuqing Du, Stas Tiomkin, Emre Kiciman, Daniel Polani, Pieter Abbeel, and Anca Dragan. Ave:  
352 Assistance via empowerment. *Advances in Neural Information Processing Systems*, 33:4560–4571,  
353 2020.
- 354 Jaime F Fisac, Monica A Gates, Jessica B Hamrick, Chang Liu, Dylan Hadfield-Menell, Malayandi  
355 Palaniappan, Dhruv Malik, S Shankar Sastry, Thomas L Griffiths, and Anca D Dragan. Pragmatic-  
356 pedagogic value alignment. In *Robotics Research: The 18th International Symposium ISRR*, pages  
357 49–57. Springer, 2020.
- 358 Jakob Foerster, Francis Song, Edward Hughes, Neil Burch, Iain Dunning, Shimon Whiteson, Matthew  
359 Botvinick, and Michael Bowling. Bayesian action decoder for deep multi-agent reinforcement  
360 learning. In *International Conference on Machine Learning*, pages 1942–1951. PMLR, 2019.
- 361 Yiming Gao, Bei Shi, Xueying Du, Liang Wang, Guangwei Chen, Zhenjie Lian, Fuhao Qiu, Guoan  
362 Han, Weixuan Wang, Deheng Ye, et al. Learning diverse policies in moba games via macro-goals.  
363 *Advances in Neural Information Processing Systems*, 34:16171–16182, 2021.
- 364 Yiming Gao, Feiyu Liu, Liang Wang, Zhenjie Lian, Weixuan Wang, Siqin Li, Xianliang Wang, Xian-  
365 han Zeng, Rundong Wang, Jiawei Wang, et al. Towards effective and interpretable human-agent  
366 collaboration in moba games: A communication perspective. *arXiv preprint arXiv:2304.11632*,  
367 2023.

- 368 Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. Cooperative inverse  
369 reinforcement learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- 370 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural  
371 information processing systems*, 29, 2016.
- 372 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–  
373 1780, 1997.
- 374 Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster. “other-play” for zero-shot  
375 coordination. In *International Conference on Machine Learning*, pages 4399–4410. PMLR, 2020.
- 376 Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward  
377 learning from human preferences and demonstrations in atari. *Advances in Neural Information  
378 Processing Systems*, 31, 2018.
- 379 Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia  
380 Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al. Human-  
381 level performance in 3d multiplayer games with population-based reinforcement learning. *Science*,  
382 364(6443):859–865, 2019.
- 383 Anirudha Majumdar, Sumeet Singh, Ajay Mandlekar, and Marco Pavone. Risk-sensitive inverse  
384 reinforcement learning via coherent risk models. In *Robotics: Science and Systems*, volume 16,  
385 page 117, 2017.
- 386 Anis Najar and Mohamed Chetouani. Reinforcement learning with human advice: A survey. *Frontiers  
387 in Robotics and AI*, 8:584075, 2021.
- 388 Andrew Y Ng, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *International  
389 Conference on Machine Learning*, volume 1, page 2, 2000.
- 390 OpenAI, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębniak,  
391 Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz,  
392 Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto,  
393 Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever,  
394 Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning.  
395 *arXiv preprint arXiv:1912.06680*, 2019.
- 396 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
397 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
398 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–  
399 27744, 2022.
- 400 Xavier Puig, Tianmin Shu, Shuang Li, Zilin Wang, Yuan-Hong Liao, Joshua B Tenenbaum, Sanja  
401 Fidler, and Antonio Torralba. Watch-and-help: A challenge for social perception and human-ai  
402 collaboration. *arXiv preprint arXiv:2010.09890*, 2020.
- 403 Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, SM Ali Eslami, and Matthew Botvinick.  
404 Machine theory of mind. In *International Conference on Machine Learning*, pages 4218–4227.  
405 PMLR, 2018.
- 406 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
407 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 408 Victor do Nascimento Silva and Luiz Chaimowicz. Moba: A new arena for game ai. *arXiv preprint  
409 arXiv:1705.10443*, 2017.
- 410 DJ Strouse, Kevin McKee, Matt Botvinick, Edward Hughes, and Richard Everett. Collaborating with  
411 humans without human data. *Advances in Neural Information Processing Systems*, 34, 2021.
- 412 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

- 413 Hua Wei, Jingxiao Chen, Xiyang Ji, Hongyang Qin, Minwen Deng, Siqin Li, Liang Wang, Weinan  
414 Zhang, Yong Yu, Liu Linc, et al. Honor of kings arena: An environment for generalization  
415 in competitive reinforcement learning. *Advances in Neural Information Processing Systems*,  
416 35:11881–11892, 2022.
- 417 H James Wilson and Paul R Daugherty. Collaborative intelligence: Humans and ai are joining forces.  
418 *Harvard Business Review*, 96(4):114–123, 2018.
- 419 Sarah A Wu, Rose E Wang, James A Evans, Joshua B Tenenbaum, David C Parkes, and Max  
420 Kleiman-Weiner. Too many cooks: Bayesian inference for coordinating multi-agent collaboration.  
421 *Topics in Cognitive Science*, 13(2):414–432, 2021.
- 422 Deheng Ye, Guibin Chen, Wen Zhang, Sheng Chen, Bo Yuan, Bo Liu, Jia Chen, Zhao Liu, Fuhao  
423 Qiu, Hongsheng Yu, et al. Towards playing full moba games with deep reinforcement learning.  
424 *Advances in Neural Information Processing Systems*, 33:621–632, 2020.
- 425 Deheng Ye, Guibin Chen, Peilin Zhao, Fuhao Qiu, Bo Yuan, Wen Zhang, Sheng Chen, Mingfei  
426 Sun, Xiaoqian Li, Siqin Li, et al. Supervised learning achieves human-level performance in moba  
427 games: A case study of honor of kings. *IEEE Transactions on Neural Networks and Learning  
428 Systems*, 2020.
- 429 Deheng Ye, Zhao Liu, Mingfei Sun, Bei Shi, Peilin Zhao, Hao Wu, Hongsheng Yu, Shaojie Yang,  
430 Xipeng Wu, Qingwei Guo, et al. Mastering complex control in moba games with deep reinforce-  
431 ment learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages  
432 6672–6679, 2020.
- 433 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse  
434 reinforcement learning. In *Association for the Advancement of Artificial Intelligence*, volume 8,  
435 pages 1433–1438. Chicago, IL, USA, 2008.