000 INNATE-VALUES-DRIVEN 001 Reinforcement Learning 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Innate values describe agents' intrinsic motivations, which reflect their inherent interests and preferences for pursuing goals and drive them to develop diverse skills that satisfy their various needs. Traditional reinforcement learning (RL) is learning from interaction based on the environment's feedback rewards. However, in real scenarios, the rewards are generated by agents' innate value systems, which differ vastly from individuals based on their needs and requirements. In other words, considering the AI agent as a self-organizing system, developing its awareness through balancing internal and external utilities based on its needs in different tasks is a crucial problem for individuals learning to support others and integrate community with safety and harmony in the long term. To address this gap, we propose a new RL model termed innate-values-driven RL (IVRL) based on combined motivations' models and expected utility theory to mimic its complex behaviors in the evolution through decision-making and learning. Then, we introduce two IVRL-based models: IV-DQN and IV-A2C. By comparing them with benchmark algorithms such as DQN, DDQN, A2C, and PPO in the Role-Playing Game (RPG) reinforcement learning test platform VIZDoom, we demonstrated that the IVRL-based models can help the agent rationally organize various needs, achieve better performance effectively.

1 INTRODUCTION

031

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

In natural systems, motivation is concerned explicitly with the activities of creatures that reflect the 032 pursuit of a particular goal and form a meaningful unit of behavior in this function Heckhausen & 033 Heckhausen (2018). From the neuroscience perspective, intrinsic motivation refers to an agent's 034 spontaneous tendencies to be curious and interested, to seek out challenges, and to exercise and develop their skills and knowledge, even without operationally separable rewards Di Domenico & Ryan (2017). Furthermore, they describe incentives relating to an activity itself, and these incentives 037 residing in pursuing an activity are intrinsic Barto (2013). Moreover, intrinsic motivations deriving 038 from an activity may be driven primarily by interest or activity-specific incentives, depending on whether the object of an activity or its performance provides the main incentive Schiefele (1996). They also fall in the category of cognitive motivation theories, which include theories of the mind 040 that tend to be abstracted from the biological system of the behaving organism Merrick (2013). 041

042 However, natural agents, like humans, often make decisions based on a blend of biological, social, 043 and cognitive motivations, as elucidated by combined motivations' model like Maslow's Hierarchy of 044 Needs Maslow (1958) and Alderfer's Existence-Relatedness-Growth (ERG) theory Alderfer (1972). Fig. 1 illustrates the five human agents with various personalities presenting different amounts of innate values and preferences based on five levels of Maslow's Hierarchy of Needs and Alderfer's 046 ERG theory. On the other hand, the AI agent can be regarded as a self-organizing system that also 047 presents various needs and motivations in its evolution through decision-making and learning to adapt 048 to different scenarios and satisfy their needs Merrick & Maher (2009).

Many researchers regard motivated behavior as behavior that involves the assessment of the con-051 sequences of behavior through learned expectations, which makes motivation theories tend to be intimately linked to theories of learning and decision-making Baldassarre & Mirolli (2013). In 052 particular, intrinsic motivation leads organisms to engage in exploration, play, strategies, and skills driven by expected rewards. The computational theory of reinforcement learning (RL) addresses how

062

063

064

065

066

067

074

054





Figure 1: The illustration five human agents with various personalities presenting different amounts of innate values and preferences based on five levels of Maslow's Hierarchy of Needs and Alderfer's *Existence-Relatedness-Growth* (ERG) theory.

Figure 2: The illustration of the proposed innate-values-driven model.

predictive values can be learned and used to direct behavior, making RL naturally relevant to studying
motivation. For example, development RL is concerned with using deep RL algorithms to tackle a
developmental problem – the intrinsically motivated acquisition of open-ended repertoires of skills
Colas et al. (2022).

079 In artificial intelligence, researchers propose various abstract computational structures to form the fundamental units of cognition and motivations, such as states, goals, actions, and strategies. For 081 intrinsic motivation modeling, the approaches can be generally classified into three categories: prediction-based Schmidhuber (1991; 2010), novelty-based Marsland et al. (2000); Merrick & 083 Maher (2009), and competence-based Barto et al. (2004); Schembri et al. (2007). Furthermore, the 084 concept of intrinsic motivation was introduced in machine learning and robotics to develop artificial systems learning diverse skills autonomously Yang & Parasuraman (2020a; 2023; 2024). The idea 085 is that intelligent machines and robots could autonomously acquire skills and knowledge under the guidance of intrinsic motivations and later exploit such knowledge and skills to accomplish tasks 087 more efficiently and faster than if they had to acquire them from scratch Baldassarre & Mirolli (2013). 880

In other words, by investigating intrinsically motivated learning systems, we would clearly improve the utility and autonomy of intelligent artificial systems in dynamic, complex, and dangerous environ-090 ments Yang & Parasuraman (2020b; 2021). Specifically, compared with the traditional RL model, 091 intrinsically motivated RL refines it by dividing the environment into an external environment and 092 an internal environment Aubret et al. (2019), which clearly generates all reward signals within the organism¹ Baldassarre & Mirolli (2013). Although the extrinsic reward signals are triggered by 094 the objects and events of the external environment, and activities of the internal environment cause 095 the intrinsic reward signals, it is hard to determine the complexity and variability of the intrinsic 096 rewards (innate values) generating mechanism. Specifically, traditional RL model is learning from interaction based on the environment's feedback rewards. However, in real world, the rewards are 098 generated by agents' innate value systems, which differ vastly from individuals based on their needs 099 and requirements. Moreover, the AI agent can be regarded as a self-organizing system that also 100 presents various needs and motivations in its evolution through decision-making and learning to adapt to different scenarios and satisfy those needs. The traditional RL can not reasonably explain its innate 101 values and motivations nor provide a long-term model to support the AI agent's lifelong development. 102

To address those gaps, we introduce the innate-values-driven reinforcement learning (IVRL) model,
 which integrates combined motivations' models and expected utility theory to describe the complex
 behaviors in AI agents' adaptation and evolution. We formalize the innate values and derive the IVRL
 model, then propose two corresponding algorithms: IV-DQN and IV-A2C. Furthermore, we compare

107

¹Here, the organism represents all the components of the internal environment in the AI agent.

108 them with benchmark RL algorithms such as DQN Mnih et al. (2015), DDQN Wang et al. (2016), A2C Mnih et al. (2016), and PPO Schulman et al. (2017) in the Role-Playing Game (RPG) RL test 110 platform VIZDoom Kempka et al. (2016); Wydmuch et al. (2019). The results demonstrate that the 111 proposed IVRL model can achieve convergence and adapt efficiently to complex and challenging 112 tasks.

114 2 APPROACH OVERVIEW

113

115

124 125

127

128

129

130

131 132

133 134 135

136

116 We assume that all the AI agents (like robots) interact in the same working scenario, and their external 117 environment includes all the other group members and mission setting. In contrast, the internal 118 environment consists of individual perception components including various sensors (such as Lidar 119 and camera), the critic module involving intrinsic motivation analysis and innate values generation, 120 the RL brain making the decision based on the feedback of rewards and description of the current state (including internal and external) from the critic module, and actuators relating to all the manipulators 121 and operators executing the RL brain's decisions as action sequence and strategies. Fig. 2 illustrates 122 the proposed innate-values-driven model. 123



Figure 3: The illustration of the IVRL model based on Expected Utility Theory.

137 Compared with the traditional RL model, our model generates the input state and rewards from the 138 critic module instead of directly from the environment, which means that the AI agent receives various 139 utilities from the environment through executing an action or strategy in the IVRL model. Moreover, 140 the individual needs to calculate innate values (expected utility) through its needs weights and current 141 utilities and then select suitable actions or strategies to optimize or maximize its accumulated expected 142 utility (Fig. 3). Specifically, we formalize the IVRL of an AI agent with an external environment using a Markov decision process (MDP) Puterman (2014). The MDP is defined by the tuple $\langle S, A, \mathcal{R}, \mathcal{T}, \gamma \rangle$ 143 where S represents the finite sets of internal state S_i^2 and external states S_e . A represents a finite 144 set of actions. The transition function $\mathcal{T}: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0, 1]$ determines the probability of a 145 transition from any state $s \in S$ to any state $s' \in S$ given any possible action $a \in A$. Assuming 146 the critic function is \mathcal{C} , which describes the individual innate value model. The reward function 147 $\mathcal{R} = \mathcal{C}(S_e) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ defines the immediate and possibly stochastic innate reward $\mathcal{C}(S_e)$ that 148 an agent would receive given that the agent executes action a which in state s and it is transitioned 149 to state $s', \gamma \in [0, 1)$ the discount factor that balances the trade-off between innate immediate and 150 future rewards.

151 152

153

158

159

161

2.1 THE EXPECTED UTILITY AND SOURCE OF RANDOMNESS

In the IVRL model, the reward is regarded as the *expected utility* Fishburn et al. (1979); Fishburn 154 (1988) generated by the agent's utility values $u(x_k)$ and the corresponding probability p_k (equation 1). 155 It is equal to the sum of each category's needs weight n_k times its current utility u_k in the IVRL 156 model (equation 3). 157

$$R_t = \sum_{i=1}^k u_k \times n_k = \mathbb{E}\left[U(p)\right] = \sum_{i=1}^k u(x_k)p_k \tag{1}$$

²The internal state S_i describes an agent's innate value distribution and presents the dominant intrinsic motivation based on the external state S_e .

Furthermore, in the IVRL model, the randomness comes from three sources. The randomness in action is from the policy function: $A \sim \pi(\cdot|s)$; the needs weight function: $W \sim \omega(\cdot|s)$ makes the randomness of innate values; the state-transition function: $S' \sim p(\cdot|s, a)$ causes the randomness in state.

Figure 4: Illustration of the trajectory of state S, needs weight W, action A, and reward R in the IVRL model.

Supposing at current state s_t an agent has a needs weight matrix N_t (equation 2) in a mission, which presents its innate value weights for different levels of needs. Correspondingly, it has a utility matrix U_t (equation 2) for specific needs resulting from action a_t . Then, we can calculate its reward R_t for a_t through equation 3 at the state s_t .

$$N_{t} = \begin{bmatrix} n_{11} & n_{12} & \cdots & n_{1m} \\ n_{21} & n_{22} & \cdots & n_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ n_{n1} & n_{n2} & \cdots & n_{nm} \end{bmatrix}; \quad U_{t} = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1m} \\ u_{21} & u_{22} & \cdots & u_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ u_{n1} & u_{n2} & \cdots & u_{nm} \end{bmatrix}$$
(2)

$$R_t = \sum_{i=1}^m \sum_{j=1}^n N_t \times U_t^T \tag{3}$$

In the process, the agent will first generate the needs weight and action based on the current state, then, according to the feedback utilities and the needs weights, calculate the current reward (expected utility) and iterate the process until the end of an episode. Fig. 4 illustrates the trajectory of state S, needs weight W, action A, and reward R, and Fig. 3 presents the corresponding IVRL model.

2.2 RANDOMNESS IN DISCOUNTED RETURNS

According to the above discussion, we define the discounted return G at t time as cumulative discounted rewards in the IVRL model (equation 4) and γ is the discount factor.

$$G_{t} = R_{t} + \gamma R_{t+1} + \gamma^{2} R_{t+2} + \dots + \gamma^{n-t} R_{n}$$
(4)

At time t, the randomness of the return G_t comes from the the rewards R_t, \dots, R_n . Since the reward R_t depends on the state S_t , action A_t , and needs weight W_t , the return G_t also relies on them. Furthermore, we can describe their randomness as follows:

State transition: $\mathbb{P}[A = a | S = s, A = a] = p(s' | s, a);$ (5)

Needs weights function:
$$\mathbb{P}[W = w | S = s] = \omega(w | s);$$

Policy function:
$$\mathbb{P}[A = a | S = s] = \pi(a | s).$$
 (7)

(6)

207 2.3 ACTION-INNATE-VALUE FUNCTION

Based on the discounted return equation 4 and its random factors – equation 5, equation 6, and equation 7, we can define the *Action-Innate-Value* function as the expectation of the discounted return G at t time (equation 8).

$$Q_{\pi,\omega}(s_t, w_t, a_t) = \mathbb{E}[G_t | S_t = s_t, W_t = w_t, A_t = a_t]$$
(8)

 $Q_{\pi,\omega}(s_t, w_t, a_t)$ describes the quality of the action a_t taken by the agent in the state s_t , using the 215 needs weight w_t generating from the needs weight function ω as the innate value judgment to execute the policy π .



Figure 5: Illustration of the IV-DQN network generating Needs-Behavior distribution.

2.4 STATE-INNATE-VALUE FUNCTION

Furthermore, we can define the *State-Innate-Value* function as equation 9, which calculates the expectation of $Q_{\pi,\omega}(s_t, w_t, a_t)$ for action A and reflects the situation in the state s_t with the innate value judgment w_t .

$$V_{\pi}(s_t, w_t) = \mathbb{E}_A[Q_{\pi,\omega}(s_t, w_t, A)] \tag{9}$$

2.5 APPROXIMATE THE ACTION-INNATE-VALUE FUNCTION

The agent's goal is to interact with the environment by selecting actions to maximize future rewards based on its innate value judgment. We make the standard assumption that a factor of γ per time-step discounts future rewards and define the future discounted return at time t as equation 4. Moreover, we can define the optimal action-value function $Q^*(s, a, w)$ as the maximum expected return achievable by following any strategy after seeing some sequence s, making corresponding innate value judgment w, and then taking action a, where ω is a needs weight function describing sequences about innate value weights and π is a policy mapping sequences to actions.

$$Q^{*}(s, w, a) = \max_{\omega, \pi} \mathbb{E}[G_{t} | S_{t} = s_{t}, W_{t} = w_{t}, A_{t} = a_{t}, \omega, \pi]$$
(10)

Since the optimal action-innate-value function obeys the Bellman equation, we can estimate the function by using the Bellman equation as an iterative update. This is based on the following intuition: if the optimal innate-value $Q^*(s', w', a')$ of sequence s' at the next time-step was known for all possible actions a' and needs weights w', then the optimal strategy is to select the reasonable action a' and rational innate value weight w', maximising the expected innate value of $r + \gamma Q^*(s', w', a')$,

$$Q^*(s, w, a) = \mathbb{E}_{s' \sim \epsilon} \left[r + \gamma \max_{w', a'} Q^*(s', w', a') \middle| s, w, a \right]$$
(11)

Furthermore, the same as the DQN Mnih et al. (2015), we use a function approximator (equation 12) to estimate the action-innate-value function.

$$Q(s, w, a; \theta) \approx Q^*(s, w, a) \tag{12}$$

We refer to a neural network function approximator with weights θ as a Q-network. It can be trained by minimising a sequence of loss function $L_i(\theta_i)$ that changes at each iteration *i*,

$$L_i(\theta_i) = \mathbb{E}_{s,w,a\sim\sigma(\cdot)}\left[(y_i - Q(s,w,a;\theta_i))^2\right]$$
(13)

$$y_i = \mathbb{E}_{s' \sim \epsilon} \left[r + \gamma \max_{w',a'} Q(s', w', a'; \theta_{i-1}) \middle| s, w, a \right]$$
(14)

267 268

230 231 232

233 234

235

236 237 238

239

247 248 249

255 256 257

258

259 260

261

262

269 Where equation 14 is the target for iteration *i* and $\sigma(s, w, a)$ is a probability distribution over sequences *s*, needs weights *w*, and action *a* that we refer to as the *needs-behavior distribution*. We

270 Algorithm 1: Innate-Values-driven DQN (IV-DQN) 271 1 Initialize replay memory \mathcal{D} to capacity N; 272 2 Initialize the action-innate-value function Q with random neural network weights; 273 3 for each episode do 274 for each environment step t do 4 275 With probability ϵ select a random action a_t and w_t , otherwise select 5 $a_t = \max_{w,a} Q(\phi(s), w, a; \theta);$ 276 Execute action a_t in emulator, calculate reward $r_t = w_t \times u_t$ based on agent current needs weights 6 277 w_t and utilities u_t , and image x_{t+1} ; 278 Set $s_{t+1} = s_t$, a_t , w_t , x_{t+1} and preprocess $\phi_{t+1} = \phi(s_{t+1})$; 7 279 8 Store transition $(\phi_t, a_t, w_t, r_t, \phi_{t+1})$ in \mathcal{D} ; Sample random minibatch of transitions $(\phi_i, a_i, w_i, r_i, \phi_{i+1})$ from \mathcal{D} ; 9 281 10 Set $y_j = \begin{cases} r_j, & \text{for terminal } \phi_{j+1}; \\ r_j + \gamma \max_{w,a} Q(\phi_{j+1}, w', a'; \theta), & \text{else.} \end{cases}$ 282 283 Perform a gradient descent step on $(y_i - Q(\phi_i, w_i, a_i; \theta))^2$ according to equation 15. 284 11 285

can approximate the gradient as follows:

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,w,a \sim \sigma(\cdot);s' \sim \epsilon} \left[\left(r + \gamma \max_{w',a'} Q(s',w',a';\theta_{i-1}) - Q(s,w,a;\theta_i) \right) \nabla_{\theta_i} Q(s,w,a;\theta_i) \right]$$
(15)

Instead of computing the full expectations in the equation 15, stochastic gradient descent is usually computationally expedient to optimize the loss function. Here, the weights θ are updated after every time step, and single samples from the *needs-behavior distribution* σ and the emulator ϵ replace the expectations, respectively.

Our approach is a model-free and off-policy algorithm, which learns about the greedy strategy a = $\max_{w,a} Q(s, w, a; \theta)$ following a *needs-behavior distribution* to ensure adequate state space exploration. Moreover, the *needs-behavior distribution* selects action based on an ϵ -greedy strategy that follows the greedy strategy with probability $1 - \epsilon$ and selects a random action with probability ϵ . Fig. 5 illustrates the action-innate-value network generating Needs-Behavior distribution.

Moreover, we utilize the *experience replay* technique Lin (1992), which stores the agent's experiences at each time-step, $e_t = (s_t, w_t, a_t, r_t, s_{t+1})$ in a data-set $\mathcal{D} = e_1, \ldots, e_N$, pooled over many episodes into a *replay memory*. During the algorithm's inner loop, we apply Q-learning updates, or minibatch updates, to samples of experience, $e \sim \mathcal{D}$, drawn at random from the pool of stored samples. After performing experience replay, the agent selects and executes an action according to an ϵ -greedy policy, as we discussed. Since implementing arbitrary length histories as inputs to a neural network is difficult, we use a function ϕ to produce our action-innate-value Q-function. Alg. 1 presents the algorithm of the IV-DQN.

313

287

289

295

314 315

316

2.6 IVRL ADVANTAGE ACTOR-CRITIC (A2C) MODEL

Furthermore, we extend our IVRL method to the Advantage Actor-Critic (A2C) version. Specifically, our IV-A2C maintains a policy network $\pi(a_t|s_t;\theta)$, a needs network $\omega(w_t|s_t;\delta)$, and a utility value network $u(s_t, a_t; \varphi)$. Since the reward in each step is equal to the current utilities $u(s_t, a_t)$ multiplying the corresponding weight of needs, the state innated-values function can be approximated by presenting it as equation 16. Then, we can get the policy gradient equation 17 and needs gradient equation 18 of the equation 16 deriving $V(s; \theta, \delta)$ according to the *Multi-variable Chain Rule*, respectively. We can update the policy network θ and needs network δ by implementing policy gradient and needs gradient, and using the temporal difference (TD) to update the value network φ .



377 Since the IVRL model differs from the traditional RL model, it uses need weights and corresponding utilities from the internal and external environment to calculate innate value rewards. The traditional

378 379 380 381 382 384 (b) Defend the Line (c) Deadly Corridor (d) Arena

(a) Defend the Center

Figure 7: The four scenarios of experiments in the VIZDoom

environment feedback-based reward RL platform can not be used in the IVRL experiments. Considering that the VIZDoom testbed Kempka et al. (2016); Wydmuch et al. (2019) can customize the experiment environment and define various utilities based on different tasks and cross-platform, we selected it to evaluate our IVRL model. We choose four scenarios: Defend the Center, Defend the Line, Deadly Corridor, and Arens (Fig. 7), and compare our models with several benchmark algorithms, such as DQN Mnih et al. (2015), DDQN Wang et al. (2016), A2C Mnih et al. (2016), and PPO Schulman et al. (2017). These models were trained on an NVIDIA GeForce RTX 3080Ti GPU with 16 GiB of RAM.

396 397 398

385 386

387 388 389

390

391

392

393

394

395

3.1 Environment Setting

399 In our experiments, we define four categories of utilities (health points, amount of ammo, environment 400 rewards, and number of killed enemies), presenting three different levels of needs: low-level safety 401 and basic needs, medium-level recognition needs, and high-level achievement needs. When the agent 402 executes an action, it receives all the corresponding innate utilities, such as health points and ammo 403 costs, and external utilities, such as environment rewards (living time) and the number of killed 404 enemies. At each step, the agent can calculate the rewards for the action by multiplying the current 405 utilities and the needed weight for them. In our experiments, we assume that the agent has no bias or 406 preference at the beginning of the game. Therefore, the initial needs weight for each utility category 407 is 0.25, fixed in the benchmark DRL algorithms' training, such as DQN, DDQN, and PPO. For more 408 details about the experiment code, please check the supplementary materials.

409 a. Defend the Center – Fig. 7(a): For this scenario, the map is a large circle where the agent is in the 410 middle, and monsters spawn along the wall. The agent's basic actions are turn-left, turn-right, and 411 attack, and the action space is 8. It needs to survive in the scenario as long as possible.

412 b. Defend the Line. - Fig. 7(b): The agent is located on a rectangular map, and the monsters are on the 413 opposite side. Similar to the defend the center scenario, the agent needs to survive as long as possible. 414 Its basic actions are move-left, move-right, turn-left, turn-right, and attack, and the action space is 32. 415

c. Deadly Corridor. – Fig. 7(c): In this scenario, the map is a corridor. The agent is spawned at one 416 end of the corridor, and a green vest is placed at the other end. Three pairs of monsters are placed on 417 both sides of the corridor. The agent needs to pass the corridor and get the vest. Its basic actions are 418 move-left, move-right, move-forward, turn-left, turn-right, and attack, and the action space is 64. 419

420 d. Arena. – Fig. 7(d): This scenario is the most challenging map compared with the other three. The 421 agent's start point is in the middle of the map, and it needs to eliminate various enemies to survive as long as possible. Its basic actions are move-left, move-right, move-forward, move-backward, turn-left, 422 turn-right, and attack, and the action space is 128. 423

424 3.2 EVALUATION 425

426 The performance of the proposed IV-DQN and IV-A2C models is shown in the Fig. 8(a), 8(d), 427 8(g), and 8(j) demonstrate that IVRL models can achieve higher average scores than traditional 428 RL benchmark methods. Especially for the IV-A2C algorithm, it presents more robust, stable, and 429 efficient properties than other models. 430

Moreover, we analyze their corresponding tendencies in different scenarios to compare the needs 431 weight differences between the IV-DQN and IV-A2C models. In the defend-the-center and defend-the-



Figure 8: The performance comparison of IV-DQN and IV-A2C agents with DQN, DDQN, PPO, and A2C in the VIZDoom.

474 475 476

473

476

477 line experiments, each category of the need weight in the IV-DQN model does not split and converges 478 to a specific range compared with its initial setting in our training (Fig. 8(b) and 8(e)). In contrast, the 479 weights of health depletion, ammo cost, and sub-goal (environment rewards) shrink to approximately 480 zero, and the weight of the number of killed enemies converges to one in the IV-A2C model. This 481 means that the top priority of the IV-A2C agent is to eliminate all the threats or adversaries in those 482 scenarios so that it can survive, which is similar to the Arena task. According to the performance in 483 those three scenarios (Fig. 8(c), 8(f), and 8(l)), the IV-A2C agent represents the characteristics of bravery and fearlessness, much like the human hero in a real battle. However, in the deadly corridor 484 mission, the needs weight of the task goal (getting the vest) becomes the main priority, and the killing 485 enemy weight switches to the second for the IV-A2C agent (Fig. 8(i)). They converge to around 0.6 and 0.4, respectively. In training, by adjusting its different needs weights to maximize rewards, the
 IV-A2C agent develops various strategies and skills to kill the encounter adversaries and get the vast
 efficiently, much like a military spy.

489 In our experiments, we found that selecting the suitable utilities to consist of the agent innate-values 490 system is critically important for building its reward mechanism, which decides the training speed 491 and sample efficiency. Moreover, the difference in the selected utility might cause some irrelevant 492 experiences to disrupt the learning process, and this perturbation leads to high oscillations of both 493 innate-value rewards and needs weight. Furthermore, the IV-DQN performs better in the Arena than 494 any other algorithm (Fig. 8(j)). However, in other experiments, the IV-A2C's performance is better 495 than the IV-DQN. It reflects that, due to different task scenarios, the small perturbation introduced 496 by the innate-values utilities may have made it difficult for the network weights in some topologies to reach convergence. Generally speaking, the performances of IV-DQN and IV-A2C are generally 497 better than traditional A2C, DQN, and PPO. 498

The innate value system serves as a unique reward mechanism driving agents to develop diverse actions or strategies satisfying their various needs in the systems. It also builds different personalities and characteristics of agents in their interaction. From the environmental perspective, due to the various properties of the tasks, agents need to adjust their innate value system (needs weights) to adapt to different tasks' requirements. These experiences also shape their intrinsic values in the long term, similar to humans building value systems in their lives.

505 506

507

4 CONCLUSION

508 This paper introduces a new RL model from individual intrinsic motivations perspectives termed 509 innate-values-driven reinforcement learning (IVRL). It is based on the expected utility theory to model 510 mimicking the complex behaviors of agent interactions in its evolution. By adjusting needs weights 511 in its innate-values system, it can adapt to different tasks representing corresponding characteristics to maximize the rewards efficiently. For theoretical derivation, we formulated the IVRL model and 512 proposed two types of IVRL models: IV-DQN and IV-A2C. Furthermore, we compared them with 513 benchmark algorithms such as DQN, DDQN, A2C, and PPO in the RPG reinforcement learning 514 test platform VIZDoom. The results prove that rationally organizing various individual needs can 515 effectively achieve better performance. Moreover, in the multi-agent setting, organizing agents with 516 similar interests and innate values in the mission can optimize the group utilities and reduce costs 517 effectively, just like "Birds of a feather flock together." in human society. Especially combined with 518 AI agents' capacity to aid decision-making, it will open up new horizons in human-multi-agent 519 collaboration. This potential is crucially essential in the context of interactions between human 520 agents and intelligent agents when considering establishing stable and reliable relationships in their 521 cooperation, particularly in adversarial and rescue mission environments.

522 For future work, we want to improve the IVRL further and develop a more comprehensive system to 523 personalize individual characteristics to achieve various tasks testing in several standard MAS testbeds, 524 such as StarCraft II, OpenAI Gym, Unity, etc. Especially in multi-object and multi-agent interaction 525 scenarios, building the awareness of AI agents to balance the group utilities and system costs and 526 satisfy group members' needs in their cooperation is a crucial problem for individuals learning 527 to support their community and integrate human society in the long term. Moreover, integrating 528 efficient deep RL algorithms with the IVRL can help agents evolve diverse skills to adapt to complex environments in MAS cooperation. Furthermore, implementing the IVRL in real-world systems, such 529 as human-robot interaction, multi-robot systems, and self-driving cars, would be challenging and 530 exciting. 531

- 532
- 533
- 534
- 535
- 536
- 537
- 538
- 539

540 REFERENCES 541

542 543	Clayton P Alderfer. Existence, relatedness, and growth: Human needs in organizational settings. 1972.
544 545 546	Arthur Aubret, Laetitia Matignon, and Salima Hassas. A survey on intrinsic motivation in reinforce- ment learning. <i>arXiv preprint arXiv:1908.06976</i> , 2019.
547 548	Gianluca Baldassarre and Marco Mirolli. Intrinsically motivated learning systems: an overview. <i>Intrinsically motivated learning in natural and artificial systems</i> , pp. 1–14, 2013.
550 551	Andrew G Barto. Intrinsic motivation and reinforcement learning. <i>Intrinsically motivated learning in natural and artificial systems</i> , pp. 17–47, 2013.
552 553 554 555	Andrew G Barto, Satinder Singh, Nuttapong Chentanez, et al. Intrinsically motivated learning of hierarchical collections of skills. In <i>Proceedings of the 3rd International Conference on Development and Learning</i> , volume 112, pp. 19. Citeseer, 2004.
556 557 558	Cédric Colas, Tristan Karch, Olivier Sigaud, and Pierre-Yves Oudeyer. Autotelic agents with intrinsically motivated goal-conditioned reinforcement learning: a short survey. <i>Journal of Artificial Intelligence Research</i> , 74:1159–1199, 2022.
559 560 561	Stefano I Di Domenico and Richard M Ryan. The emerging neuroscience of intrinsic motivation: A new frontier in self-determination research. <i>Frontiers in human neuroscience</i> , 11:145, 2017.
562 563	P.C. Fishburn. <i>Nonlinear Preference and Utility Theory</i> . Johns Hopkins series in the mathematical sciences. Wheatsheaf Books, 1988. ISBN 9780745005461.
564 565	Peter C Fishburn, Peter C Fishburn, et al. Utility theory for decision making. Krieger NY, 1979.
566 567	Jutta Heckhausen and Heinz Heckhausen. Motivation and action. Springer, 2018.
568 569 570 571	Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaśkowski. ViZ- Doom: A Doom-based AI research platform for visual reinforcement learning. In <i>IEEE Conference</i> <i>on Computational Intelligence and Games</i> , pp. 341–348, Santorini, Greece, Sep 2016. IEEE. doi: 10.1109/CIG.2016.7860433. The Best Paper Award.
572 573 574	Long-Ji Lin. Reinforcement learning for robots using neural networks. Carnegie Mellon University, 1992.
575 576 577	Stephen Marsland, Ulrich Nehmzow, and Jonathan Shapiro. A real-time novelty detector for a mobile robot. <i>EUREL European Advanced Robotics Systems Masterclass and Conference</i> , 2000.
578	Abraham Harold Maslow. A dynamic theory of human motivation. 1958.
579 580 581 582	Kathryn E Merrick. Novelty and beyond: Towards combined motivation models and integrated learning architectures. <i>Intrinsically motivated learning in natural and artificial systems</i> , pp. 209–233, 2013.
583 584 585	Kathryn E Merrick and Mary Lou Maher. <i>Motivated reinforcement learning: curious characters for multiuser games</i> . Springer Science & Business Media, 2009.
586 587 588	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. <i>nature</i> , 518(7540):529–533, 2015.
589 590 591	Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In <i>International conference on machine learning</i> , pp. 1928–1937. PMLR, 2016.
592 593	Martin L Puterman. Markov decision processes: discrete stochastic dynamic programming. John Wiley & Sons, 2014.

11

594 595 596 597	Massimiliano Schembri, Marco Mirolli, and Gianluca Baldassarre. Evolution and learning in an intrinsically motivated reinforcement learning robot. In <i>Advances in Artificial Life: 9th European Conference, ECAL 2007, Lisbon, Portugal, September 10-14, 2007. Proceedings 9</i> , pp. 294–303. Springer, 2007.
598 599	Ulrich Schiefele. Motivation und Lernen mit Texten. Hogrefe Göttingen, 1996.
600 601	Jürgen Schmidhuber. Curious model-building control systems. In <i>Proc. international joint conference</i> on neural networks, pp. 1458–1463, 1991.
603 604	Jürgen Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). <i>IEEE transactions on autonomous mental development</i> , 2(3):230–247, 2010.
605 606	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
608 609 610	Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In <i>International conference on machine learning</i> , pp. 1995–2003. PMLR, 2016.
611 612 613	Marek Wydmuch, Michał Kempka, and Wojciech Jaśkowski. ViZDoom Competitions: Playing Doom from Pixels. <i>IEEE Transactions on Games</i> , 11(3):248–259, 2019. doi: 10.1109/TG.2018.2877047. The 2022 IEEE Transactions on Games Outstanding Paper Award.
614 615 616 617	Qin Yang and Ramviyas Parasuraman. Hierarchical needs based self-adaptive framework for coopera- tive multi-robot system. In 2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC), pp. 2991–2998. IEEE, 2020a.
618 619 620	Qin Yang and Ramviyas Parasuraman. Needs-driven heterogeneous multi-robot cooperation in rescue missions. In 2020 IEEE International Symposium on Safety, Security, and Rescue Robotics (SSRR), pp. 252–259. IEEE, 2020b.
621 622 623 624	Qin Yang and Ramviyas Parasuraman. How can robots trust each other for better cooperation? a relative needs entropy based robot-robot trust assessment model. In 2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC), pp. 2656–2663. IEEE, 2021.
625 626 627 628 629	Qin Yang and Ramviyas Parasuraman. A hierarchical game-theoretic decision-making for co- operative multiagent systems under the presence of adversarial agents. In <i>Proceedings of</i> <i>the 38th ACM/SIGAPP Symposium on Applied Computing</i> , SAC '23, pp. 773–782, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9781450395175. doi: 10.1145/3555776.3577642. URL https://doi.org/10.1145/3555776.3577642.
630 631 632	Qin Yang and Ramviyas Parasuraman. Bayesian strategy networks based soft actor-critic learning. ACM Transactions on Intelligent Systems and Technology, 2024.
633	
634	
635	
636	
637	
638	
639	
640	
641	
642	
643	
644	
645	
646	
647	

A APPENDIX

We provide the code of IV-QDN and IV-A2C models for the corresponding experiments. Please check the supplemental material.