

Learning Individual Intrinsic Reward in Multi-Agent Reinforcement Learning via Incorporating Generalized Human Expertise

Xuefei Wu^{*1}, Xiao Yin^{*1}, Yuanyang Zhu^{†2} and Chunlin Chen¹

Abstract—Efficient exploration in multi-agent reinforcement learning (MARL) is a challenging problem when receiving only a team reward, especially in environments with sparse rewards. A powerful method to mitigate this issue involves crafting dense individual rewards to guide the agents toward efficient exploration. However, individual rewards generally rely on manually engineered shaping-reward functions that lack high-order intelligence, thus it behaves ineffectively than humans regarding learning and generalization in complex problems. To tackle these issues, we combine the above two paradigms and propose a novel framework, LIGHT (Learning Individual Intrinsic reward via Incorporating Generalized Human expertise), which can integrate human knowledge into MARL algorithms in an end-to-end manner. LIGHT guides each agent to avoid unnecessary exploration by considering both individual action distribution and human expertise preference distribution. Then, LIGHT designs individual intrinsic rewards for each agent based on actionable representational transformation relevant to Q-learning so that the agents align their action preferences with the human expertise while maximizing the joint action value. Experimental results demonstrate the superiority of our method over representative baselines regarding performance and better knowledge reusability across different sparse-reward tasks on challenging scenarios.

I. INTRODUCTION

Cooperative multi-agent reinforcement learning (MARL) is an important branch in the field of artificial intelligence (AI), playing a crucial role in sequential challenging decision-making problems, such as in autonomous driving [1], sensor networks [2], [3] and robotics control [4]. Centralized training with decentralized execution (CTDE) paradigm has gained substantial attention in cooperative MARL that aims to facilitate agent cooperation by providing global state information during training and executing only based on local observations during execution [5], [6], [7]. Recent research has witnessed extensive investigation into value decomposition methods under the paradigm of CTDE. Since the advances in these MARL approaches [8], [9], [10], [11], [12], well-designed auxiliary rewards are indispensable

to guide agent collaboration or competition. Unfortunately, many cooperative multi-agent tasks currently only offer common team rewards [13].

In real-world multi-agent systems, sparse team rewards present a significant challenge [14]. Existing algorithms rely on dense-reward environments for guiding efficient cooperation strategies [15], but these conditions rarely hold in real-world scenarios. In MARL, handling sparse-reward environments often involves enhancing agent exploration, a common approach demonstrated to be effective in various tasks [16], [17]. However, relying solely on exploration can be insufficient for determining which specific actions trigger rare non-zero rewards. It highlights the need for methods that promote discovery and aid in identifying and reinforcing these critical action-reward connections. Individual intrinsic rewards can be a promising solution to this problem [18], [19], [20], typically by distributing a combination of team and individual rewards among agents [21], [22]. It can change agents' learning objectives, potentially leading to unexpected behaviors divergent from desired team outcomes. An important line of work that leverages human knowledge is a promising method to improve the learning process efficiently [23], [24], [25], [26]. However, a major challenge in leveraging human knowledge is how to obtain the representation of the provided knowledge.

To solve this problem, we propose a novel method called LIGHT, Learning Individual Intrinsic reward via Incorporating Generalized Human expertise, which can plug into the value decomposition algorithms to promote the learning efficiency, especially in sparse-reward setting tasks under the widely-used assumption of CTDE. Our key insight is that instead of using human knowledge to directly guide the agent to interact with the environment, we integrate human knowledge to produce the intrinsic reward to induce the agent to achieve better exploration. Specifically, at each time step, LIGHT learns a parameterized intrinsic reward function by considering the action distribution and human preference that outputs an intrinsic reward for each agent to implicitly induce diversified behaviors.

We evaluate LIGHT on two representative benchmarks: Level-Based Foraging (LBF) and StarCraft Multi-Agent Challenge (SMAC), where empirical results demonstrate our method's superior performance over other baselines. We conduct further component studies to show the effectiveness of individual intrinsic reward incorporated with generalized human expertise for agent learning, which confirms that each component is a key part of LIGHT. We also find that the behavior of LIGHT obtains better alignment with human

^{*} Equal contribution, [†] Corresponding author.

¹Xuefei Wu, Xiao Yin, and Chunlin Chen are with the Department of Control Science and Intelligent Engineering, School of Management and Engineering, Nanjing University, Nanjing 210093, China. (email: xuefei.wu, xiaoyin@smail.nju.edu.cn, clchen@nju.edu.cn).

²Yuanyang Zhu is with the Laboratory of Data Intelligence and Interdisciplinary Innovation, School of Information Management, Nanjing University, Nanjing 210023, China. (email: yuanyangzhu@nju.edu.cn)

This work was supported in part by the China Postdoctoral Science Foundation under Grant Number 2025T180877, the National Key Research and Development Program of China under Grant 2023YFD2001003, Major Science and Technology Project of Jiangsu Province under Grant BG2024041 and the Fundamental Research Funds for the Central Universities under Grant 011814380048.

knowledge, indicating it provides an efficient method to incorporate human preference into the learning process.

II. PRELIMINARIES

Multi-agent Markov Decision Process. A fully cooperative multi-agent task can be modeled as an extension of a Decentralized Partially Observable Markov Decision Process (Dec-POMDP), which is defined by a tuple $G = \langle \mathcal{N}, \mathcal{S}, \mathcal{A}, \mathcal{P}, \Omega, O, r, \gamma \rangle$, where \mathcal{N} represents a finite set of agents with $\mathcal{N} \equiv \{1, 2, \dots, n\}$, $s \in \mathcal{S}$ denotes the global state of the environment. At each time step t , each agent $i \in \mathcal{N}$ selects an action $a_i \in \mathcal{A}$ to formulate a joint action $\mathbf{a} \equiv [a_i]_{i=1}^n \in \mathcal{A}^n$. The joint action leads to a shared reward according to the reward function $r(s, \mathbf{a})$ and a transition to a new state based on the transition probability function $s' \sim \mathcal{P}(\cdot | s, \mathbf{a})$. Given the partial observability, each agent i receives an individual observation $o_i \in \Omega$, associated with the observation probability function $O(o_i | s, a_i)$. The action-observation history for each agent i is denoted as $\tau_i \in \mathcal{T} \equiv (\Omega \times \mathcal{A})^*$, and the joint action-observation history is $\boldsymbol{\tau} \in \mathcal{T}^n$. The objective for all agents is to find an optimal joint policy $\boldsymbol{\pi} = \langle \pi_1, \dots, \pi_n \rangle$ to maximize the expected cumulative reward $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r^t]$, where $\gamma \in [0, 1)$ is a discount factor.

Centralized Training with Decentralized Execution. CTDE is a prevalent paradigm in the MARL, where each agent learns a policy only on its own action observations, and the centralized critic provides a global perspective, offering gradient updates that are informed by the joint state and action space. A promising enhancement of the CTDE framework is the application of value decomposition. This technique enables agents to individually learn utility functions that collectively optimize the joint action-value function, thereby providing a clear framework for credit assignment among agents. For the integrity of multi-agent value decomposition methods, adherence to the Individual-Global-Max (IGM) principle is essential. To ensure consistency [27] for multi-agent value decomposition methods, it should satisfy the IGM principle:

$$\operatorname{argmax}_{\mathbf{a}} Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) = \begin{pmatrix} \operatorname{argmax}_{a_1} Q_1(\tau_1, a_1) \\ \vdots \\ \operatorname{argmax}_{a_n} Q_n(\tau_n, a_n) \end{pmatrix}, \quad (1)$$

Individual Intrinsic Reward. Individual intrinsic rewards have become pivotal in MARL settings, particularly where extrinsic rewards are sparse. The LIIR [18] incorporates individual intrinsic rewards into the Actor-Critic (AC) algorithm by designing an attentive reward mechanism. It can be formally defined as

$$R_{i,t}^{\text{proxy}} = \sum_{l=0}^{\infty} \gamma^l (r_{t+l}^{\text{ex}} + \lambda r_{i,t+l}^{\text{in}}), \quad (2)$$

where λ is a tunable hyperparameter that balances the influence of global extrinsic rewards and the intrinsic reward. The introduction of λ allows for a dynamic adjustment between

learning from the environment and fostering behaviors motivated by the agent's own experiences. The resultant proxy value function for each agent is expressed as

$$V_i^{\text{proxy}}(s_{i,t}) = \mathbb{E}_{a_{i,t}, s_{i,t+1}, \dots} [R_{i,t}^{\text{proxy}}], \quad (3)$$

where a_i is the action space of each agent i at time step t . Then the proxy value function is applied to optimize the agents' policy. ICQL [28] introduces a local uncertainty measure to enhance learning in decentralized agents through intrinsic motivation. This technique fosters a nuanced understanding of the environment by encouraging exploration through uncertainty-driven intrinsic rewards, thereby complementing the LIIR framework's strategy for balancing intrinsic and extrinsic rewards.

III. METHOD

In this section, we propose a novel end-to-end cooperative MARL framework called LIGHT, solving the MARL with sparse reward effectively via leveraging human knowledge to generate individual intrinsic rewards. This section introduces the methodology of LIGHT. We begin with our motivation and then provide a detailed explanation of the implementation of LIGHT.

A. Motivation

Sparse-reward scenarios are typical in RL applications, where agents may not have enough information to develop an optimal behavior and may learn to exploit suboptimal but easily accessible solutions. In CTDE, each agent acts independently with local observability and receives only the factorization global reward. This shared reward structure complicates the learning of cooperative policies, as it can be tough to discern which actions contribute to the success of the group. It makes it difficult for an algorithm to successfully learn a cooperative team policy in such a setting. One could also consider a manual specification of dense rewards. However, designing a useful reward function is notoriously difficult and time-consuming. This naturally leads to the fundamental question: Can we design informative rewards that will guide the agent to efficiently explore and accelerate the agent's learning process?

Recalling the learning process of humans, they rarely approach the acquisition of new skills in a vacuum. They adeptly draw upon a wealth of prior knowledge derived from analogous tasks to formulate an initial strategy. Human cognition is characterized by the ability to extract explainable, task-solving heuristics that exhibit a degree of generalizability across related domains. If this essence of human knowledge could be distilled into the fabric of RL agents by encoding logical inferences into the neural architectures that underpin their learning processes. These agents could bypass the initial trial-and-error phase and embark immediately on the refinement of effective strategies.

One natural solution is to integrate human knowledge to produce the reward. Following this idea, we introduce our LIGHT framework as illustrated in Fig. 1, which synthesizes human-derived insights into intrinsic rewards,

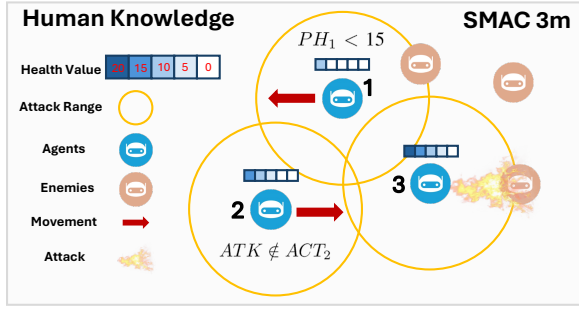


Fig. 2. Example of a task with the rule about human knowledge.

distributions in response to the iterative refinements of the team policy. The purpose of using the intrinsic reward is to optimize and maximize the global reward, as shown in Fig. 1-(b). The sparse team rewards make it hard to guide the joint policy to the optimal policy, like the blue line. Through the factorization of team utility, each agent rarely receives beneficial signals to update its policy towards the optimal policy (orange line). To this end, we introduce the intrinsic reward, which incorporates human knowledge to implicitly guide individual agents to achieve better exploration.

At each time step t , the agent i receives its intrinsic reward by computing the negative Euclidean distance between its action distribution and soft logic rules. When receiving its local observation o_t^i at time step t , we compute the intrinsic reward r_t^i as

$$r_t^i = -\|\phi^i(\mathcal{H}(o_t^i)) - \phi^i(A_t^i)\|_2, t = 1, \dots, T, \quad (4)$$

where $\phi^i(\cdot)$ denotes an operable representational transformation, $\mathcal{H}(o_t^i)$ is the action distribution of the soft logic rules, and A_t^i is a value distribution of each agent. Each agent tries to reach its returns while maximizing the team's returns. It can also be regarded as the difference in actions between human preferences and individual agents.

We redefine the reward function to quantify the contribution of individual agents to the success of the team, combining the environment's extrinsic reward with added intrinsic incentives:

$$R_t = r_t^{ex} + \lambda \frac{1}{N} \sum_{i=1}^N r_t^{i,j}, \quad (5)$$

where R_t is used to update the mixed network parameters θ . To relieve the sparsity of extrinsic rewards, intrinsic rewards are added to make the mixing parameters non-trivially updated at each time step.

Overall Learning. Under the training execution framework of CTDE, LIGHT learns by sampling a multitude of transitions from a replay buffer, and the loss function for the mixing network parameter is represented as

$$L(\theta) = \left(R_t + \gamma \max_{u'} Q_{\theta^-}^{\text{tot}}(s_{t+1}, u') - Q_{\theta}^{\text{tot}}(s_t, u_t) \right)^2, \quad (6)$$

where θ_i^- is the parameter of the target network for the mixing network. The individual Q-values update the intrinsic reward value of each agent. The overall learning objective is

to minimize the following loss:

$$L_i(\theta_i) = \left[r_t^i + \gamma \max_{u^i} Q_{\theta_i^-}^i(o_{t+1}^i, u^i) - Q_{\theta_i}^i(o_t^i, u_t^i) \right]^2, \quad (7)$$

where L_i represents the individual loss value of each agent, and r_t^i stands for the intrinsic reward. The total loss function used in this work is expressed as follows:

$$L = L_{TD}(\theta) + \lambda_K L_i(\theta_i), \quad (8)$$

where λ_K is denoted as a coefficient set for individual loss.

IV. EXPERIMENTS

In this section, we evaluate LIGHT on the widely-used and challenging tasks over Level-Based Foraging (LBF) [30] and StarCraft Multi-Agent Challenge (SMAC) [31] benchmarks, where SMAC includes dense-reward and sparse-reward settings. We compared with five representative MARL algorithms: MASER [19], LIIR [18], VDN [8], QMIX [9] and QTRAN [27]. Then, we conduct ablation studies on LIGHT to better understand each component's effect. To ensure fair evaluation, we conducted all experiments with five random seeds, and the results are plotted using mean \pm std.

V. EXPERIMENTAL ENVIRONMENT SETTINGS

A. Level-Based Foraging (LBF)

LBF is a sparse reward and mixed cooperative-competitive environment, where agents collect randomly-scattered food items by navigating a grid world. Each agent navigates a 10×10 grid world while observing a 5×5 sub-grid centered at its current position. The food collection is successful when the sum of the near-by agents' levels exceeds the food's level. Then, agents can receive rewards that are equal to the level of the food they collect, divided by their level. We construct two scenarios with different quantities of agents and food to evaluate the performance of all methods, including 4 agents with 2 food (4-agent & 2-food) and 3 agents with 3 food (3-agent & 3-food).

B. StarCraft Multi-Agent Challenge (SMAC)

SMAC simulates intricate scenarios from the acclaimed real-time strategy game StarCraft, providing a robust environment for validating our proposed methodologies. Each agent in our setup is equipped with a local observation vector drawing from a wealth of tactical data points, including the proximity, placement, vitality, shield status, and classification of both allied and adversarial units. A noteworthy feature of the simulation is the dynamic shield regeneration that activates after a designated duration of non-combat status, alongside an armor mechanic that necessitates depletion before any reduction in an agent's health pool can occur. Our experiments are conducted against a formidable opposition in-game built-in AI, which operates at a difficulty level of 7.

The Reward Setting for Two Scenarios. In our investigation, we benchmark our LIGHT framework against leading state-of-the-art MARL algorithms across both dense-reward and sparse-reward environments. The specific configurations

TABLE I
THE CONFIGURATIONS OF REWARD SETTINGS.

	Dense reward	Sparse reward
Win	+200	+200
Enemy's death	+10	+10
Ally's death	-5	-5
Enemy's health	-Enemy's remaining health	-
Ally's health	+Ally's remaining health	-
Other elements	+/- with other elements	-

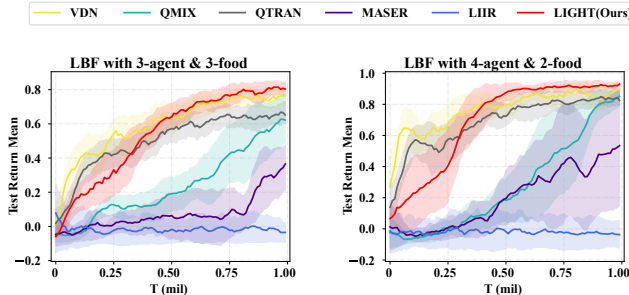


Fig. 3. Performance comparison with different baselines on two constructed scenarios of LBF.

for these reward structures are detailed in Table I. The dense-reward scenario adheres to the conventional reward paradigm, aligning with the standard reward mechanism implemented in the PyMARL framework. In this setup, agents receive frequent feedback, with rewards dispensed during the course of an episode. Conversely, the sparse-reward scenario in SMAC experiments aligns with the parameters set forth in the MASER framework [19]. Within this challenging environment, the reward signal is markedly reduced, with agents receiving feedback exclusively upon the culmination of the task or in response to critical in-game events such as the elimination of enemy units or the loss of allied ones. This comparative study aims to scrutinize the efficacy of our LIGHT framework in learning and decision-making processes under contrasting reward densities, thereby offering comprehensive insights into its performance relative to established MARL benchmarks in strategic game settings.

VI. HYPERPARAMETER SETTINGS

The hyperparameters of configurations follow the source code provided by the authors while keeping it consistent across all baselines for fairness. The detailed hyperparameters for LIGHT and other baselines on LBF and SMAC can be found in Table II. In addition, LIGHT employs greedy action selection, and λ_K is set to 0.02.

A. Overall Performance Comparison

Performance on LBF. It shows in Fig. 3 that LIGHT outperforms all baselines on constructed scenarios of LBF. The poor performance of LIIR can be attributed to its failure to explore cooperative strategies. Although MASER generates subgoals from the experience replay buffer, it fails to perform on all scenarios where rewards are highly sparse. QMIX

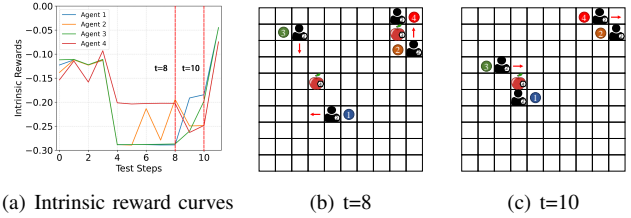


Fig. 4. An example of the intrinsic reward curves and some keyframes on LBF with 4-agent & 2-food.

requires more steps to explore superior policies, which may be due to the lack of intrinsic rewards for each agent, resulting in an inability to guide agents toward efficient exploration with sparse rewards. Despite alleviating QMIX's monotonic constraint, QTRAN still performs poorly due to its insufficient constraint relaxation in challenge scenarios. VDN yields competitive performance that spends fewer time steps to achieve a higher return than LIGHT before 0.5 steps. However, it may be due to neglecting human expertise, which leads to an ultimately suboptimal performance compared with LIGHT. In summary, LIGHT demonstrates impressive performance on effectiveness over all baselines, indicating that it is advantageous to consider human preferences when designing individual intrinsic rewards for MARL.

Performance on SMAC. We first ran all the considered algorithms on the conventional dense-reward setting. As shown in Fig. 5, most algorithms can achieve near 100% win rates on 3m, 4m, and 5m scenarios, and only LIIR can not work well on the 2m_vs_1z map. It indicates that LIGHT and most algorithms are able to accomplish the tasks on scenarios with dense rewards. Note that here we apply our LIGHT architecture to the individual Q-network of the QMIX and denote it as LIGHT. We then evaluate all the algorithms on the difficult sparse-reward setting. As shown in Fig. 6, LIGHT performs better than all baselines on all scenarios. QMIX achieves satisfactory performance, which may contribute to the efficient credit assignment for facilitating collaboration among agents. QTRAN does not yield satisfactory performance, which may be due to the relaxation in practice that is insufficient for challenging domains. Additionally, MASER performs better than VDN and LIIR, which should benefit from its individual intrinsic reward for each agent based on the actionable representation. This can help agents reach their subgoals while maximizing the joint action value. In summary, our method achieves impressive performance on all scenarios, demonstrating the advantage of LIGHT with attentive design through incorporating human expertise.

We apply our LIGHT architecture to the individual Q-network of the fine-tuned QMIX and VDN and denote them as LIGHT-QMIX and LIGHT-VDN, respectively. As shown in Fig. 8, our LIGHT-QMIX surpasses the fine-tuned QMIX by a large margin in almost all scenarios, especially in 3m, 4m, 5m, and 2m_vs_1z scenarios. Our LIGHT-VDN also significantly improves the performance of the fine-tuned VDN, and it even surpasses the fine-tuned QMIX in most scenarios and achieves close performance to the LIGHT-QMIX, which minimizes the gaps between the VDN and

TABLE II
THE CONFIGURATIONS OF HYPERPARAMETER SETTINGS FOR LIGHT AND OTHER BASELINES.

	LIGHT	MASER	LIIR	QMIX	VDN	QTRAN
Buffer Size	5000	5000	32	5000	5000	5000
Batch Size	32	32	32	32	32	32
Test Interval (SMAC)	2000	2000	2000	2000	2000	2000
Test Interval (LBF)	1000	1000	1000	1000	1000	1000
Test Episodes	32	32	32	32	32	32
Optimizer	RMSProp	RMSProp	RMSProp	RMSProp	RMSProp	RMSProp
Agent Runner	episode	episode	parallel	episode	episode	episode
Learning Rate	0.0005	0.0005	0.0005	0.0005	0.0005	0.0005
TD Discounted Factor	0.99	0.99	0.99	0.99	0.99	0.99
Start Exploration Rate	1	1	0.5	1	1	1
End Exploration Rate	0.05	0.05	0.01	0.05	0.05	0.05
Epsilon Anneal Step	50000	50000	50000	50000	50000	50000
Target Update Interval	200	200	200	200	200	200
Mixing Embed Dimension	32	32	-	32	-	64

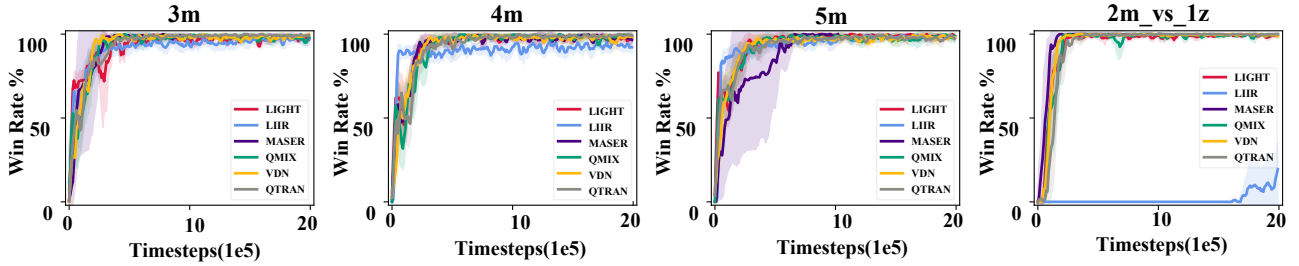


Fig. 5. Performance comparison with baselines on dense-reward setting scenarios.

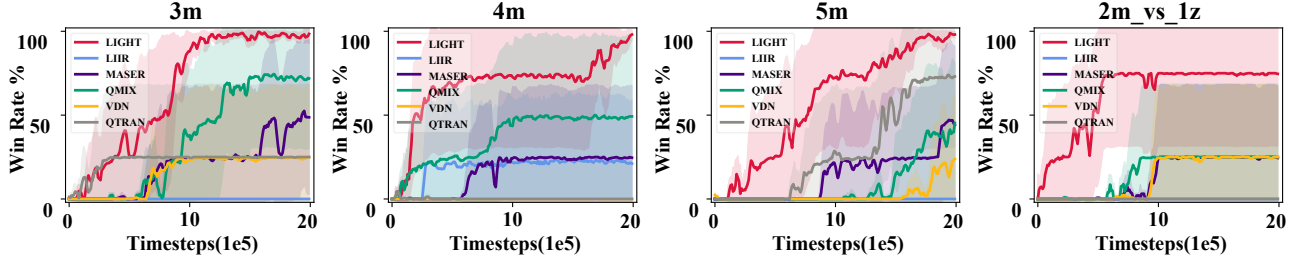


Fig. 6. Performance comparison with baselines on sparse-reward setting scenarios.

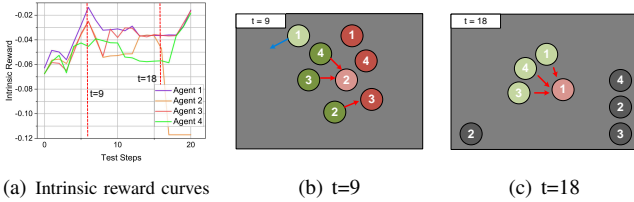


Fig. 7. An example of the intrinsic reward curves and auxiliary illustration on 4m map. Green circles and red circles represent allies and enemies, respectively, where the darker color indicates the higher health value of the agent. Gray circles indicate that the agent was killed. Blue arrows represent move, red arrows represent attack.

QMIX algorithms.

B. Visualizing the Learned Intrinsic Reward

Interpretability of LIGHT on LBF. We visualize these rewards to satisfy curiosity about how powerful the learned intrinsic reward function is for the policy learning. As shown in Fig. 4, we display some keyframes on 3-agent & 3-food map, where the red arrows represent the direction of movement. As shown in Fig. 4(a), we find that the intrinsic reward for agent-2 increases a lot, implying that this action is

good at this time. Fig. 4(b) provides visualization evidence that agent-4 cannot eat food of a higher level than itself, which requires waiting to cooperate with other agents. At this stage, agent-2 moves north and cooperates with agent-4 to eat the food item in the next step. As seen from Fig. 4(a), at time step 10, we see that agent-1 and agent-3 achieve the supreme intrinsic rewards compared to others, which indicates that they all select great behaviors in this step. In Fig. 4(c), agent-3 selects the *East* action, and then agent-1 and agent-3 form a cooperative alliance to jointly acquire the food item exceeding their individual capability levels. The above analysis shows the superiority of the designed intrinsic reward function that plays a pivotal role in providing important feedback for agents, enabling real-time behavior evaluation and policy optimization in complex cooperative or competitive environments.

Interpretability of LIGHT on SMAC. To better understand the impact of the learned intrinsic reward function on policy training, we propose a direct visualization of these rewards. Specifically, we plot the intrinsic rewards assigned to each agent at every step of a complete trajectory during

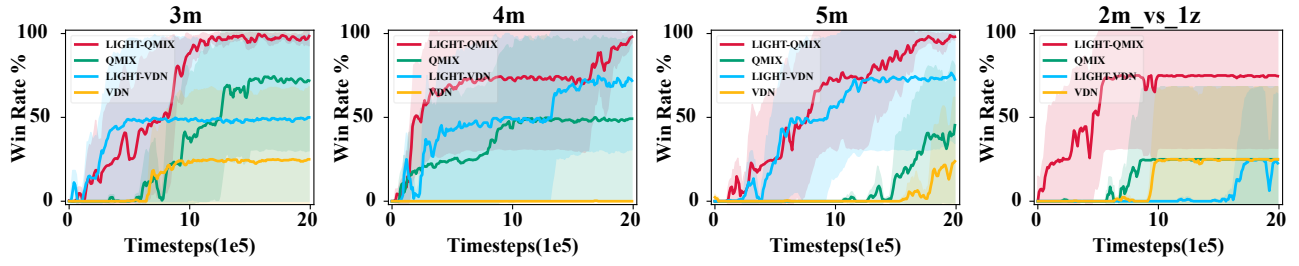


Fig. 8. LIGHT is plugged into two different baselines on sparse-reward scenarios.

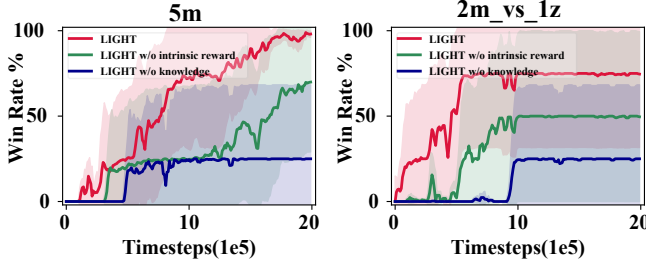


Fig. 9. The ablation study results of LIGHT, LIGHT w/o intrinsic reward, and LIGHT w/o knowledge on 5m and 2m_vs_1z scenarios.

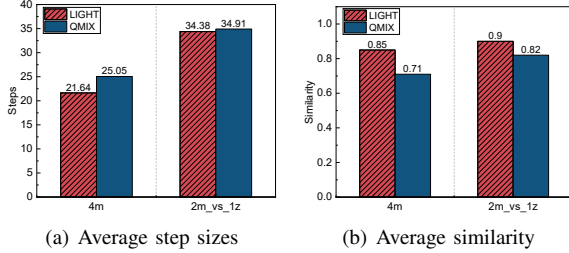


Fig. 10. A comparison of alignment with human knowledge preference behaviors of LIGHT and QMIX on 4m and 2m_vs_1z maps.

testing. It is important to note that the intrinsic rewards do not influence the learned policy and are not utilized in trajectory generation. For better visualization and clarity, we select two test replays from scenarios 4m and 2m_vs_1z and chart the intrinsic rewards for all agents involved. Figures 4 and 11 illustrate these intrinsic rewards in the 4m and 2m_vs_1z scenarios, respectively.

In Fig. 7-(a), at time step 9, the intrinsic reward of agent-1 rises to near 0 because it has the lowest health and chooses *move* rather than *attack* as revealed in Fig. 7-(b). It indicates that the *attack* action is not a good behavior for agent-1 at this time. After time step 18, agent-2 has been attacked until it dies, and it receives a large negative intrinsic reward.

On the 2m_vs_1z map, we find that agent-1 receives a larger reward at time step 21 in Fig. 11-(a). This may be because it takes on more of the task of drawing fire and attacking at moments of higher health than its companions, while agent-2 stops firing and falls behind its team with a lower intrinsic reward. After several steps, agent-2 obtains cooperative skills, and the intrinsic rewards gradually increase. Meanwhile, agent-1 is killed and receives a lower intrinsic reward. By visualizing the intrinsic reward curves for the two maps, the results illustrate that the generated intrinsic reward can effectively provide diverse feedback signals, which are highly informative in assessing the agents' immediate behaviors.

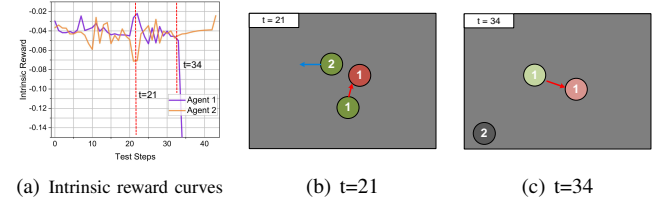


Fig. 11. An example of the intrinsic reward curves and auxiliary illustration on 2m_vs_1z map. Green circles and red circles represent allies and enemies, respectively, where the darker color indicates the higher health value of the agent. Gray circles indicate that the agent was killed. Blue arrows represent *move*, red arrows represent *attack*.

C. Ablation Studies

To understand the impact of each component of LIGHT, we conduct ablation studies to answer the following questions: (1) How does the intrinsic reward influence performance? (2) How does human knowledge influence performance? To study components (1) and (2), LIGHT w/o intrinsic reward represents removing intrinsic rewards, and LIGHT w/o human knowledge represents human knowledge with randomly generated distributions, respectively. We carry out ablation studies on 5m and 2m_vs_1z maps. As shown in Fig. 9, the ablation of each part of LIGHT brings a noticeable decrease in performance. Specifically, the performance of LIGHT w/o human knowledge decreases, which indicates that human knowledge is beneficial in guiding the agent to explore better. Besides, the performance of LIGHT w/o intrinsic reward is lower than LIGHT, which indicates that the intrinsic rewards can ultimately induce better exploration in sparse-reward environments. To summarize, LIGHT, conditioned on all components, gives the best performance, which could improve exploration with the given limited learning time steps.

D. Behavior analysis

In addition to evaluating the performance of LIGHT, we are more curious about whether the behavior aligns with the given human knowledge. To study how human knowledge influences the behavior of LIGHT, we compare the consistent behavior of LIGHT and QMIX on 4m and 2m_vs_1z maps with sparse-reward settings. Specifically, we make a statistical analysis of whether the action of the agent is consistent with the given human knowledge at each time step during testing for 100 episodes. Here, we consider the behavior to be consistent if the agent produces an action that is consistent with human knowledge. On the 4m and 2m_vs_1z maps, as shown in Fig. 10-(a), we can find that

the average steps of LIGHT are 21.64 and 34.38 over 100 episodes, while QMIX requires 25.05 and 34.91 time steps, respectively. It indicates that LIGHT requires fewer time steps to solve the task with a more optimal policy. As shown in Fig. 10-(b), LIGHT has more behaviors aligning with the given human knowledge than QMIX across both scenarios, which indicates that our method can efficiently capture the given human knowledge to facilitate the learning process. The above case studies demonstrate that LIGHT can not only promote learning efficiency but also provide a novel method to incorporate the given human preference into the behavior of the agents.

VII. CONCLUSION

In this work, we propose a novel value decomposition framework called LIGHT to leverage human expertise to accelerate the learning process of MARL agents. LIGHT produces intrinsic rewards to induce the agent to explore efficiently by considering both each agent's action distribution and human preference at an early stage. This end-to-end framework can be combined with existing value decomposition algorithms to deal with the sparse-reward setting tasks. Experiments on the challenging LBF and SMAC benchmarks show that our method obtains the best performance on almost all sparse-reward maps, and the intrinsic reward module of LIGHT can help the behavior of agents better align with human preferences. In the future, this simple yet effective method further motivates us to explore more effective ways to utilize intrinsic rewards by incorporating human knowledge in more challenging tasks.

REFERENCES

- [1] Y. Cao, W. Yu, W. Ren, and G. Chen, "An overview of recent progress in the study of distributed multi-agent coordination," *IEEE Trans. Industr. Inform.*, pp. 427–438, 2012.
- [2] C. Zhang and V. Lesser, "Coordinated multi-agent reinforcement learning in networked distributed pomdps," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2011.
- [3] A. Ghorbani, A. Abid, and J. Zou, "Interpretation of neural networks is fragile," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, pp. 3681–3688.
- [4] M. Hüttenrauch, A. Šošić, and G. Neumann, "Guided deep reinforcement learning for swarm systems," *arXiv:1709.06011*, 2017.
- [5] Y. Yang, J. Hao, G. Chen, H. Tang, Y. Chen, Y. Hu, C. Fan, and Z. Wei, "Q-value path decomposition for deep multiagent reinforcement learning," in *Proceedings of the International Conference on Machine Learning*, 2020, pp. 10 706–10 715.
- [6] T. Wang, T. Gupta, A. Mahajan, B. Peng, S. Whiteson, and C. Zhang, "RODE: Learning roles to decompose multi-agent tasks," in *Proceedings of the International Conference on Learning Representations*, 2020, pp. 1–20.
- [7] T. Wang, H. Dong, V. Lesser, and C. Zhang, "ROMA: Multi-agent reinforcement learning with emergent roles," in *Proceedings of the International Conference on Machine Learning*, 2020, pp. 9876–9886.
- [8] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, et al., "Value-decomposition networks for cooperative multi-agent learning based on team reward," in *Proceedings of the International Conference on Autonomous Agents and MultiAgent Systems*, 2018, pp. 2085–2087.
- [9] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson, "QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning," in *Proceedings of the International Conference on Machine Learning*, 2018, pp. 4295–4304.
- [10] J. Wang, Z. Ren, T. Liu, Y. Yu, and C. Zhang, "QPLEX: Duplex dueling multi-agent Q-learning," in *Proceedings of the International Conference on Learning Representations*, 2020, pp. 1–27.
- [11] Z. Liu, Y. Zhu, Z. Wang, Y. Gao, and C. Chen, "Mixrts: Toward interpretable multi-agent reinforcement learning via mixing recurrent soft decision trees," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 47, no. 5, pp. 4090–4107, 2025.
- [12] Z. Liu, Y. Zhu, and C. Chen, "NA²Q: Neural attention additive model for interpretable multi-agent q-learning," in *Proceedings of the International Conference on Machine Learning*, vol. 202, 2023, pp. 22 539–22 558.
- [13] L. Busoniu, R. Babuska, and B. De Schutter, "A comprehensive survey of multiagent reinforcement learning," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 38, no. 2, pp. 156–172, 2008.
- [14] M. Sadeghlou, M. R. Akbarzadeh-T, and M. B. Naghibi-S, "Dynamic agent-based reward shaping for multi-agent systems," in *Iranian Conference on Intelligent Systems*, 2014, pp. 1–6.
- [15] A. Wong, T. Bäck, A. V. Kononova, and A. Plaat, "Multiagent deep reinforcement learning: Challenges and directions towards human-like approaches," *arXiv:2106.15691*, 2021.
- [16] I.-J. Liu, U. Jain, R. A. Yeh, and A. Schwing, "Cooperative exploration for multi-agent deep reinforcement learning," in *Proceedings of the International conference on machine learning*, 2021, pp. 6826–6836.
- [17] A. Mahajan, T. Rashid, M. Samvelyan, and S. Whiteson, "Maven: Multi-agent variational exploration," *Advances in neural information processing systems*, vol. 32, 2019.
- [18] Y. Du, L. Han, M. Fang, J. Liu, T. Dai, and D. Tao, "Liir: Learning individual intrinsic reward in multi-agent reinforcement learning," in *Advances in neural information processing systems*, vol. 32, 2019.
- [19] J. Jeon, W. Kim, W. Jung, and Y. Sung, "MASER: Multi-agent reinforcement learning with subgoals generated from experience replay buffer," in *Proceedings of the International Conference on Machine Learning*, 2022, pp. 10 041–10 052.
- [20] L. Wang, Y. Zhang, Y. Hu, W. Wang, C. Zhang, Y. Gao, J. Hao, T. Lv, and C. Fan, "Individual reward assisted multi-agent reinforcement learning," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2022, pp. 23 417–23 432.
- [21] Z. Xu, Y. Bai, B. Zhang, D. Li, and G. Fan, "Haven: hierarchical cooperative multi-agent reinforcement learning with dual coordination mechanism," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 10, 2023, pp. 11 735–11 743.
- [22] X. Xu, T. Huang, P. Wei, A. Narayan, and T.-Y. Leong, "Hierarchical reinforcement learning in starcraft ii with human expertise in subgoals selection," *arXiv preprint arXiv:2008.03444*, 2020.
- [23] M. Fischer, M. Balunovic, D. Drachler-Cohen, T. Gehr, C. Zhang, and M. Vechev, "DI2: training and querying neural networks with logic," in *Proceedings of the International Conference on Machine Learning*, 2019, pp. 1931–1941.
- [24] Y. Zhu, Z. Wang, C. Chen, and D. Dong, "Rule-based reinforcement learning for efficient robot navigation with space reduction," *IEEE/ASME Transactions on Mechatronics*, vol. 27, no. 2, pp. 846–857, 2022.
- [25] Y. Zhu, X. Yin, and C. Chen, "Extracting decision tree from trained deep reinforcement learning in traffic signal control," *IEEE Transactions on Computational Social Systems*, vol. 10, no. 4, pp. 1997–2007, 2023.
- [26] P. Zhang, J. Hao, W. Wang, H. Tang, Y. Ma, Y. Duan, and Y. Zheng, "Kogun: accelerating deep reinforcement learning via integrating human suboptimal knowledge," *arXiv preprint arXiv:2002.07418*, 2020.
- [27] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi, "QTRAN: Learning to factorize with transformation for cooperative multi-agent reinforcement learning," in *Proceedings of the International Conference on Machine Learning*, 2019, pp. 5887–5896.
- [28] S. W. Wendelin Böhmer, Tabish Rashid, "Exploration with unreliable intrinsic reward in multi-agent reinforcement learning," in *Proceedings of the International Conference on Machine Learning*, 2019.
- [29] F. Zhang, C. Jia, Y.-C. Li, L. Yuan, Y. Yu, and Z. Zhang, "Discovering generalizable multi-agent coordination skills from multi-task offline data," in *Proceedings of the International Conference on Learning Representations*, 2023.
- [30] F. Christianos, L. Schäfer, and S. Albrecht, "Shared experience actor-critic for multi-agent reinforcement learning," in *Advances in Neural Information Processing Systems*, 2020, pp. 10 707–10 717.
- [31] M. Samvelyan, T. Rashid, C. Schroeder de Witt, G. Farquhar, N. Nardelli, T. G. Rudner, C.-M. Hung, P. H. Torr, J. Foerster, and S. Whiteson, "The StarCraft Multi-Agent Challenge," in *Proceedings of the International Conference on Autonomous Agents and MultiAgent Systems*, 2019, pp. 2186–2188.