

LINDA: Multi-Agent Local Information Decomposition for Awareness of Teammates

Jiahan CAO^{†1}, Lei YUAN^{†2}, Jianhao WANG³, Shaowei ZHANG¹,
Chongjie ZHANG³, Yang YU^{1,4} & De-Chuan ZHAN^{1,4*}

¹*School of Artificial Intelligence, Nanjing University, Nanjing 210000, China;*

²*Department of Computer Science and Technology, Nanjing University, Nanjing 210000, China;*

³*Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing 100084, China;*

⁴*Polixir Technologies, Nanjing 210000, China*

Abstract In cooperative multi-agent reinforcement learning (MARL), where agents only have access to partial observations, efficiently leveraging local information is critical. During long-time observations, agents can build *awareness* for teammates to alleviate the restriction of partial observability. However, previous MARL methods usually neglect awareness learning from local information for better collaboration. To address this problem, we propose a novel framework, multi-agent *Local Information Decomposition for Awareness of teammates* (LINDA), with which agents learn to decompose local information and build awareness for each teammate. We model the awareness as stochastic random variables and perform representation learning to ensure the informativeness of awareness representations by maximizing the mutual information between awareness and the actual trajectory of the corresponding agent. LINDA is agnostic to specific algorithms and can be flexibly integrated with different MARL methods. Sufficient experiments show that the proposed framework learns informative awareness from local partial observations for better collaboration and significantly improves the learning performance, especially on challenging tasks.

Keywords Multi-agent System, Reinforcement Learning, Teammates Awareness, Centralized Training with Decentralized Execution (CTDE), StarCraft II

Citation Jiahan Cao, Lei Yuan, Jianhao Wang, Shaowei Zhang, Chongjie Zhang, Yang Yu, De-Chuan Zhan. LINDA: Multi-Agent Local Information Decomposition for Awareness of Teammates. *Sci China Inf Sci*, for review

1 Introduction

Learning how to achieve effective collaboration is a significant problem in cooperative multi-agent reinforcement learning (MARL) [7, 9, 33]. However, non-stationarity and partial observability are two major challenges in MARL [22]. Non-stationarity arises from frequent interactions among multiple agents, with which the changes in the policy of an agent will affect the optimal policy of other agents. Partial observability occurs in many practical MARL applications such as autonomous vehicle teams [4, 47] and intelligent warehouse systems [6, 19], where agents' sensor inputs are limited by their field of view.

For the problems of non-stationarity and partial observability, an appealing paradigm is Centralized Training and Decentralized Execution (CTDE) [18]. Over the course of training, the global state of all the agents is shared in a central controller. During execution, each agent makes individual decisions by local observations to achieve collaboration. Many CTDE methods have been proposed recently [27, 32, 36–38, 42] and worked effectively. Besides, many methods [26, 27, 36] adopt a GRU [5] cell to encode historical observations and actions into a local trajectory to alleviate the problem of partial observability. They tend to tacitly assume that neural networks can automatically extract specific information from trajectories for better policy learning, but it is not very easy in practice.

Despite that agents only have a limited view of their surroundings and communicating with other agents is infeasible or unreliable [30], we can still utilize global information, including states of the

* Corresponding author (email: zhandc@nju.edu.cn)

[†]These authors contributed equally to this work.

environment and other agents during training under the CTDE paradigm. Some previous works [10,11,24] use a modeling network to model the states and actions of opponents (teammates). Still, they suffer from unstable opponent modeling when opponents’ policies change. [25] handles the problem with fixed teammates and only controls one agent, but in practice, we need to train multiple agents’ policies simultaneously.

To further alleviate the problem of partial observability and stabilize the modeling of changing teammates, instead of directly modeling the changeable teammates’ policies, we propose to learn awareness for agents to extract knowledge of each teammate from local information. To clearly show our motivation, we start from a 5-agent scenario. The awareness is generated locally but we analyze all the agents’ awareness in a unified global representation space to introduce what we expect a well-learned awareness representation space to be like. As shown in Figure 1, A4 is in the field of view of A1 and A2, and thus A1 and A2 both have sufficient knowledge of A4’s state. Therefore, A1 and A2 should have relatively consistent awareness for A4. On the other hand, even though A1 and A2 cannot observe A5 at the current timestep t , they still have knowledge of A5 because A5 is visible at the previous timestep $t - 1$. They can build awareness based on their past observations. Besides, different observing timesteps lead to a different knowledge of the target agent. Therefore, different agents’ awareness of the same target should also vary and complement each other. In summary, intuitively, we expect that a well-learned awareness encoder should generate consistent and complementary awareness embeddings in the representation space.

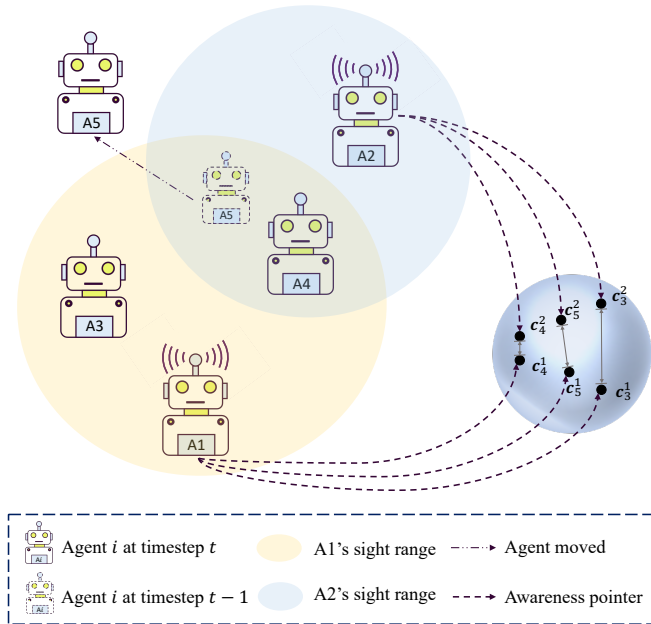


Figure 1 Illustration of building awareness in the representation space. Consider a 5-agent scenario, where A1 (A_i denotes Agent i) observes A3 and A4, A2 observes A4 and A5 at timestep t . A1’s and A2’s awareness for A4 should be consistent in the representation space due to their common observation for A4. A3 is visible to A1 while invisible to A2, making A2 hard to build accurate awareness. A5 is visible to A1 and A2 at timestep $t - 1$, but moves out of both of their sights at timestep t . Though, A1 and A2 are still able to build awareness for A5 based on their past observations.

To learn such informative awareness, in this paper, we propose a novel framework, multi-agent Local INformation Decomposition for Awareness of teammates (LNIDA), with which agents can build awareness for each teammate by learning to decompose local information in their local networks. The awareness incorporates the knowledge about other agents, such as states and strategies. Awareness is modeled as stochastic random variables and generated from the learned decomposition mapping from local trajectories. To facilitate informativeness of the learned awareness, we apply an information-theoretic loss which maximizes the mutual information between awareness and the actual trajectory of the corresponding agent. Based on the popular value-based learning methods [27, 32, 36] under the CTDE paradigm, the auxiliary mutual information loss acts as a regularizer with the global temporal difference (TD) error. During execution, the global information, including other agents’ trajectories and the global state, is removed. Agents infer awareness for teammates from their local trajectories. By learning awareness, agents can effectively exploit the information embedded in the local trajectories to further alleviate the problem of partial observability. In Section 4.2, We reveal that the proposed LINDA framework works to optimize the consistency and complementarity of awareness embeddings in the representation space.

LINDA is not built on a specific algorithm and can be easily integrated with MARL methods that follow the paradigm of CTDE. We apply LINDA to three existing MARL methods, VDN [32], QMIX [27], and QPLEX [36]. We evaluate the effectiveness of LINDA in two benchmark environments frequently used in

multi-agent system research, Level-based foraging (LBF) [1] and StarCraft II^{*}) unit micromanagement benchmark [28, 35]. Experimental results show that LINDA significantly improves learning performance by virtue of awareness learning. We further demonstrate the interpretable characteristics of learned awareness and the relationships among the awareness of different agents.

Our main contributions are:

- We propose a novel framework for MARL, and move a step towards leveraging local information by learning decomposition for awareness of teammates to alleviate the problem of partial observability.
- The LINDA framework is agnostic to specific algorithms, and is applicable to existing MARL methods that follow the paradigm of CTDE.
- Sufficient experimental results demonstrate that awareness learning is robust to diverse tasks of different difficulties, and significantly improves the learning performance, especially on the tough tasks. In the challenging SMAC [28, 35] benchmark, we propose LINDA-QMIX and LINDA-QPLEX, and achieve state-of-the-art results.

2 Related Work

Multi-agent Reinforcement Learning (MARL): Deep MARL has witnessed prominent progress in recent years. Many methods have emerged under the CTDE paradigm. Most of them are roughly divided into two categories: policy-based methods and value-based methods. MADDPG [17], COMA [8] and MAAC [13] are typical policy-based methods that explore the optimization of multi-agent policy gradient methods. Another category of approaches, value-based methods, mainly focus on the factorization of the value function. VDN [32] proposes to decompose the team value function into agent-wise value functions by an additive factorization. Following the Individual-Global-Max (IGM) principle [29], QMIX [27] improves the way of value function decomposition by learning a mixing network, which approximates a monotonic function value decomposition. QPLEX [36] takes a duplex dueling network architecture to factorize the joint value function, which achieves a full expressiveness power of IGM. Weighted QMIX [26] uses a weighted projection to place more importance on the better joint actions, and proposes two algorithms, Centrally-Weighted (CW) QMIX and Optimistically-Weighted (OW) QMIX.

Representation Learning in MARL: Learning an effective representation in MARL is receiving significant attention. ROMA [39] constructs a stochastic role embedding space to lead agents to different policies based on different roles. NDQ [41] learns a message representation to achieve expressive and succinct communication. RODE [40] uses an action encoder to learn action representations and applies clustering methods to decompose joint action spaces into restricted role action spaces to reduce the policy search space. LILI [45] learns latent representations to capture the relationship between its behavior and the other agent's future strategy and use it to influence the other agent. Unlike previous works, our approach focuses on awareness representation learning in the agents' local networks by learning to decompose local information.

Local Information Decomposition in MARL: Recently, researchers have proposed some methods that deal with the decomposition of local information. ASN [43] proposes a new framework named Action Semantics Network to explicitly represent the action semantics between agents. CollaQ [46] learns to decompose the Q-function of an agent into two parts, depending on its own state and nearby observable agents, respectively. UPDeT [12] decomposes the local observations into different parts and then uses Universal Policy Decoupling Transformer to get the policy. However, these approaches need to manually divide the local observations into each agent's part first. Manual observation division may be hard to be directly applied to complex scenarios where each agent's part of the observation is tightly coupled with each other. For example, in MOBA Game AI [44] and multi-agent connected autonomous driving [23], where observational inputs are visual images, parts of each agent in the images are hard to be manually decoupled. Therefore, an automatically learned decomposition is necessary for local information decomposition.

Opponent Modelling in MARL: Modelling opponents (teammates) in MARL is a well-studied field, with which agents learn to predict other agents' mental states (e.g., intentions, beliefs, and desires) for better coordination [2]. DRON [10] learns a modelling network to reconstruct the actions of opponents from full observations. DRIQN [11] appends an extra part to capture the actions of other agents and learns latent representations to improve the policy. OMDDPG [24] takes a further step to use variational

^{*})StarCraft II are trademarks of Blizzard Entertainment™

autoencoders to model other agents with local information during execution, which eliminates the need of opponents’ observations and actions during decentralized execution. LIAM [25] aims to learn latent representations to capture the relationship between the learning agent and modelled agents by encoder-decoder architectures only using agents’ local information. LIAM models fixed teammates’ policies and only trains one protagonist agent, because directly modelling the teammates’ inconstant behaviours leads to unstable training. However, fixing teammates’ policies is not scalable for training multiple agents’ policies simultaneously. Therefore, LINDA takes a further step for multiple agents to learn simultaneously through modelling awareness for teammates. The idea of opponent modelling from partial observations is similar to decomposing local information into awareness for others. However, LINDA starts from awareness alignment and uses information-theoretic tools to derive the learning objective. LINDA focuses more on constructing a meaningful awareness latent space instead of directly modelling the teammates’ changeable policies to stabilize the training process.

3 Preliminaries

Multi-agent Reinforcement Learning. In our work, we consider a fully cooperative multi-agent task that can be modelled by a Dec-POMDP [20] $G = \langle I, S, A, P, R, \Omega, O, n, \gamma \rangle$, where I is the finite set of n agents, $s \in S$ is the true state of the environment, A is the finite action set, and $\gamma \in [0, 1)$ is the discount factor. We consider partially observable settings, where agent i is only accessible to a local observation $o^i \in \Omega$ according to the observation function $O(s, i)$. Each agent has a observation history $\tau^i \in T \equiv (\Omega \times A)^*$. At each timestep, each agent i selects an action $a^i \in \pi^i(a | \tau^i)$, forming a joint action $\mathbf{a} = \langle a^1, \dots, a^n \rangle \in \mathcal{A}$, results in the next state s' according to the transition function $P(s'|s, \mathbf{a})$ and a shared reward $r = R(s, \mathbf{a})$ for each agent. The joint policy $\boldsymbol{\pi}$ induces a joint action-value function: $Q_{tot}^{\boldsymbol{\pi}}(\boldsymbol{\tau}, \mathbf{a}) = \mathbb{E}_{s_0: \infty, \mathbf{a}_0: \infty} [\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, \mathbf{a}_0 = \mathbf{a}, \boldsymbol{\pi}]$, where $\boldsymbol{\tau}$ is the joint action-observation history.

Value Function factorization MARL. This paper considers value function factorization in collaborative multi-agent systems (e.g., VDN [32], QMIX [27], QPLEX [36]). These three methods all follow the IGM (Individual-Global-Max) principle proposed by QTRAN [29], which asserts the consistency between joint and local greedy action selections by the joint value function $Q_{tot}(\boldsymbol{\tau}, \mathbf{a})$ and individual value functions $[Q_i(\tau^i, a^i)]_{i=1}^n$:

$$\forall \boldsymbol{\tau} \in \mathcal{T}, \arg \max_{\mathbf{a} \in \mathcal{A}} Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) = \left(\arg \max_{a^1 \in \mathcal{A}} Q_1(\tau^1, a^1), \dots, \arg \max_{a^n \in \mathcal{A}} Q_n(\tau^n, a^n) \right). \quad (1)$$

VDN utilizes the additivity to factorize the global value function $Q_{tot}^{\text{VDN}}(\boldsymbol{\tau}, \mathbf{a})$:

$$Q_{tot}^{\text{VDN}}(\boldsymbol{\tau}, \mathbf{a}) = \sum_{i=1}^n Q_i(\tau^i, a^i). \quad (2)$$

While QMIX constrains the global value function $Q_{tot}^{\text{QMIX}}(\boldsymbol{\tau}, \mathbf{a})$ with monotonicity property:

$$\forall i \in \mathcal{N}, \frac{\partial Q_{tot}^{\text{QMIX}}(\boldsymbol{\tau}, \mathbf{a})}{\partial Q_i(\tau^i, a^i)} > 0. \quad (3)$$

These two structures are sufficient conditions for the IGM principle but not necessary [36]. To achieve a complete IGM function class, QPLEX [36] uses a duplex dueling network architecture by decomposing the global value function $Q_{tot}^{\text{QPLEX}}(\boldsymbol{\tau}, \mathbf{a})$ as:

$$Q_{tot}^{\text{QPLEX}}(\boldsymbol{\tau}, \mathbf{a}) = V_{tot}(\boldsymbol{\tau}) + A_{tot}(\boldsymbol{\tau}, \mathbf{a}) = \sum_{i=1}^n Q_i(\boldsymbol{\tau}, a^i) + \sum_{i=1}^n (\lambda^i(\boldsymbol{\tau}, \mathbf{a}) - 1) A_i(\boldsymbol{\tau}, a^i). \quad (4)$$

The difference among the three methods is in the mixing networks, with increasing representational complexity. Our proposed framework LINDA follows the value factorization learning paradigm but focuses on enhancing the learning ability of agents’ individual local networks. Different global mixing networks in VDN, QMIX, and QPLEX can be freely applied to LINDA.

4 Method

In this section, we will propose multi-agent Local Information Decomposition for Awareness of teammates (LINDA), a novel framework that introduces the concept of awareness to alleviate partial observability and promote collaboration in MARL.

LINDA is a value-based MARL framework under the paradigm of centralized training with decentralized execution (CTDE) [15, 21]. Over the course of training, each agent builds awareness for teammates from local information, which includes historical local observations and actions. Based on the awareness

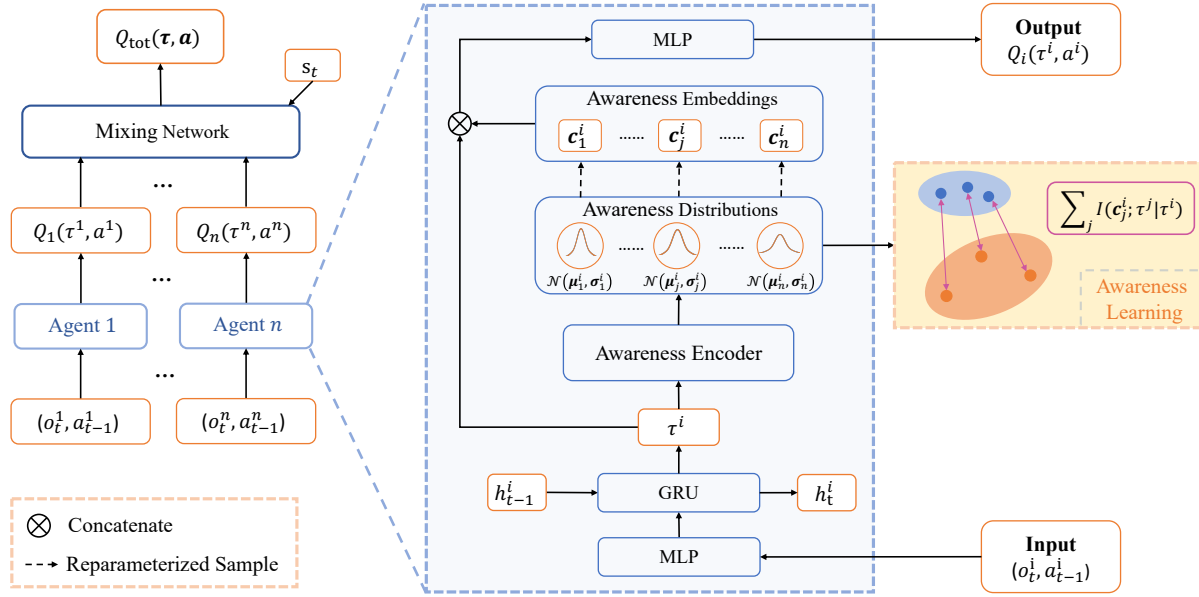


Figure 2 Structure of LINDA. For agent i , the GRU cell encodes the current observation and hidden historical state into an embedding vector of the local trajectory, denoted as τ^i . τ^i is then fed into the awareness encoder to generate n awareness distributions, the awareness representations c_1^i, \dots, c_n^i for teammates are sampled from the respective distributions using reparameterization trick for gradient flow. The awareness representations and the trajectory are concatenated and fed into an MLP network to get the local action-value function $Q_i(\tau^i, a^i)$. During centralized training, a mixing network is used to estimate the global action value $Q_{tot}(\tau, a)$ and compute the TD error. We propose an additional information-theoretic regularizer to facilitate awareness representation learning.

for each teammate, the agent produces its local action value function. In the global mixing network, the action value functions of all the agents are gathered to estimate the global action value and compute the Temporal Difference (TD) error for optimization. Another information-theoretic loss function is used to optimize the awareness distributions. During decentralized execution, the mixing network and trajectories of other agents are removed. Each agent builds awareness for teammates from local historical observations and makes individual decisions dependent on the awareness for the teammates.

4.1 The LINDA Architecture

As Figure 2 show, The LINDA framework focuses on learning awareness of teammates based on local information in each agent's individual network. For agent i , LINDA uses a GRU [5] cell to encode historical observations and actions into trajectory τ^i . τ^i is fed into an awareness encoder f_i with parameter θ_c^i to build awareness of each agent. The awareness encoder learns a decomposition mapping for each agent, and outputs n multivariate Gaussian distributions $\mathcal{N}(\mu_1^i, \sigma_1^i), \dots, \mathcal{N}(\mu_n^i, \sigma_n^i)$, where n is the number of agents, and μ_j^i, σ_j^i are the awareness mean and awareness variance for agent j respectively. Awareness representation embeddings c_1^i, \dots, c_n^i are sample from the corresponding multivariate Gaussian distributions, where c_j^i denotes agent i 's awareness for agent j . To ensure the gradient is tractable for the sampling operation, we apply the reparameterization trick [14]. To sample from a Gaussian distribution $z \sim \mathcal{N}(\mu, \sigma)$, it converts the random variable z into $z = \mu + \sigma\epsilon$, where $\epsilon \sim \mathcal{N}(0, 1)$, such that the gradient descent can be backpropagated through the sampling operation. Formally, for agent i , its n awareness

embeddings are generated by:

$$\begin{aligned} (\boldsymbol{\mu}^i, \boldsymbol{\sigma}^i) &= f_i(\tau^i; \boldsymbol{\theta}_c^i) \\ \mathbf{c}_j^i &= \boldsymbol{\mu}_j^i + \boldsymbol{\sigma}_j^i \odot \boldsymbol{\epsilon}_j^i, \quad \boldsymbol{\epsilon}_j^i \sim \mathcal{N}(0, 1), \quad \text{for } j = 1, 2, \dots, n \end{aligned} \quad (5)$$

where \odot is element-wise production, \mathbf{c}_j^i denotes agent i 's awareness for agent j , and $\boldsymbol{\theta}_c^i$ is the parameters of the awareness encoder f_i . The awareness encoder f_i inputs local trajectory τ^i and outputs the mean and variance of n awareness distributions, $(\boldsymbol{\mu}^i, \boldsymbol{\sigma}^i) = (\boldsymbol{\mu}_1^i, \dots, \boldsymbol{\mu}_n^i; \boldsymbol{\sigma}_1^i, \dots, \boldsymbol{\sigma}_n^i)$.

Since the awareness is designed for agents, awareness alone will cause the loss of environmental information. Therefore we concatenate agent i 's awareness embeddings $\mathbf{c}_1^i, \dots, \mathbf{c}_n^i$ together with its trajectory τ^i to compute the local action value Q_i by the local utility network. During the centralized training, action values of all the agents together with the global state s_t are fed into a mixing network to produce the global action value Q_{tot} and compute TD error for gradient descent. In our implementation, we try three different kinds of mixing networks, VDN [32], QMIX [27], and QPLEX [36] for their monotonic approximation. To facilitate awareness learning, our approach additionally applies an information-theoretic regularization loss $\mathcal{L}_c(\boldsymbol{\theta}_c^i)$ for agent i 's local network. The overall objective to minimize is

$$\mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}_{TD}(\boldsymbol{\theta}) + \lambda \sum_{i=1}^n \mathcal{L}_c(\boldsymbol{\theta}_c^i), \quad (6)$$

where

$$\mathcal{L}_{TD}(\boldsymbol{\theta}) = \left[r + \gamma \max_{\mathbf{a}'} Q_{tot}(\boldsymbol{\tau}', \mathbf{a}'; \boldsymbol{\theta}^-) - Q_{tot}(\boldsymbol{\tau}, \mathbf{a}; \boldsymbol{\theta}) \right]^2$$

($\boldsymbol{\theta}^-$, are the parameters of a periodically updated target network) is the TD loss, $\boldsymbol{\theta}$ is the all parameters of the framework, and λ is a scaling factor. We will then discuss the definition and optimization of the regularization loss $\mathcal{L}_c(\boldsymbol{\theta}_c^i)$.

4.2 Optimized Awareness Objective and Variational Bound

The latent awareness representations for each agent are hard to be learned automatically. Therefore an auxiliary loss function is necessary. Intuitively, we expect the learned awareness to be informative, which means that the learned awareness needs to incorporate actual information about others. Since our framework works under the CTDE paradigm, other agents' trajectories are accessible during training but inaccessible during execution. We can utilize other agents' trajectories for awareness centralized training, but awareness representations are entirely generated from the individual local trajectory for decentralized execution. To this end, we establish the relationship between agent i 's awareness for agent j , denoted as \mathbf{c}_j^i , and agent j 's actual trajectory τ^j , by maximizing their mutual information conditioned on agent i 's local trajectory τ^i . For agent i , the objective for optimizing its awareness for all the agents is to maximize

$$J_c(\boldsymbol{\theta}_c^i) = \sum_{j=1}^n I(\mathbf{c}_j^i; \tau^j | \tau^i). \quad (7)$$

During centralized training, the overall objective for awareness learning is:

$$\max_{\boldsymbol{\theta}_c} \sum_{i=1}^n J_c(\boldsymbol{\theta}_c^i) = \max_{\boldsymbol{\theta}_c} \sum_{i=1}^n \sum_{j=1}^n I(\mathbf{c}_j^i; \tau^j | \tau^i). \quad (8)$$

By unfolding the mutual information term and swapping the summation order, we can rewrite the objective Eq. 8 into:

$$\max_{\boldsymbol{\theta}_c} \sum_{j=1}^n \left(\sum_{i=1}^n H(\mathbf{c}_j^i | \tau^i) - \sum_{i=1}^n H(\mathbf{c}_j^i | \tau^i, \tau^j) \right). \quad (9)$$

For the same target agent j , the objective maximizes all the agents' entropy of their awareness distributions for agent j conditioned on their local trajectories, while minimizing the entropy conditioned on local trajectories and the target trajectory τ^j . Due to the partial observability, awareness for agent j is

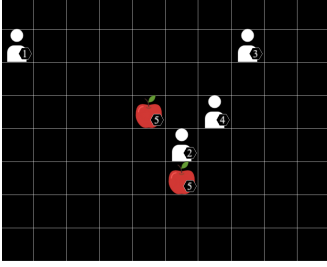


Figure 3 The Level-Based Foraging environment.



Figure 4 The StarCraft II Micromanagement environment.

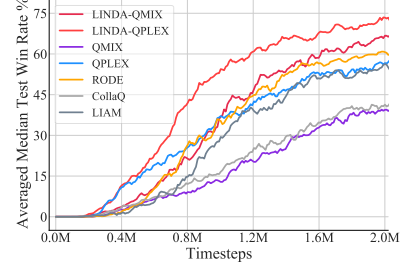


Figure 5 The median test win rate %, averaged across 3 hard scenarios and 6 super hard scenarios.

not entirely precise. Conditioned on the local trajectories alone, the uncertainty of the awareness distributions for agent j should be large. By pushing the entropy higher, it prevents the awareness embeddings from collapsing to the same point in the representation space. It means that conditioned on multiple views from different agents, the awareness for the same object should be diverse and complementary.

The second term suggests that when the target trajectory τ_j is given, the awareness distributions for agent j should get more deterministic. It means that conditioned on sufficient information, the agents' awareness for the same object should get aligned and consistent in the representation space.

However, directly optimizing the objective Eq. 7 is difficult because computation involving mutual information is intractable. Inspired by [3], we introduce a variational estimator to derive a lower bound for the mutual information term:

$$\begin{aligned} I(\mathbf{c}_j^i; \tau^j | \tau^i) &= \mathbb{E}_{\tau, \mathbf{c}_j^i} \left[\log \frac{p(\mathbf{c}_j^i, \tau^j | \tau^i)}{p(\mathbf{c}_j^i | \tau^i) p(\tau^j | \tau^i)} \right] \\ &\geq -\mathbb{E}_{\tau} [\mathcal{CE} [p(\mathbf{c}_j^i | \tau^i) \| q_{\xi}(\mathbf{c}_j^i | \tau^i, \tau^j)]] + \mathbb{E}_{\tau^i} [H(\mathbf{c}_j^i | \tau^i)], \end{aligned} \quad (10)$$

where $q_{\xi}(\mathbf{c}_j^i | \tau^i, \tau^j)$ is the variational posterior estimator with parameter ξ , and \mathcal{CE} is the Cross-Entropy operator. A detailed derivation is shown in Appendix A. By using a replay buffer \mathcal{D} , we can rewrite the lower bound in Eq. 10 and derive its loss function:

$$\mathcal{L}_c(\theta_c) = \sum_{j=1}^n \mathbb{E}_{\tau \sim \mathcal{D}} [D_{\text{KL}} [p(\mathbf{c}_j^i | \tau^i) \| q_{\xi}(\mathbf{c}_j^i | \tau^i, \tau^j)]], \quad (11)$$

where KL is the Kullback-Leibler divergence operator. Finally, the overall loss function is

$$\begin{aligned} \mathcal{L}(\theta) &= \left[r + \gamma \max_{\mathbf{a}'} Q_{\text{tot}}(\tau', \mathbf{a}'; \theta^-) - Q_{\text{tot}}(\tau, \mathbf{a}; \theta) \right]^2 && \text{(Temporal Difference Loss)} \\ &+ \lambda \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_{\tau \sim \mathcal{D}} [D_{\text{KL}} [p(\mathbf{c}_j^i | \tau^i) \| q_{\xi}(\mathbf{c}_j^i | \tau^i, \tau^j)]], && \text{(Awareness Learning Loss)} \end{aligned} \quad (12)$$

where λ is an adjustable hyper-parameter to achieve a trade-off between the temporal difference loss and the summation of all agents' awareness learning loss. Over the course of training, global information including trajectories of all the agents is used to compute the mutual information loss. During execution, the awareness learning module is removed, and each agent infers awareness for others conditioned on its local trajectory.

5 Experiments

In this section, we design experiments to answer the following questions: (1) Can LINDA be applied to multiple existing MARL methods and improve their performances? (Section 5.2) (2) Does the superiority of LINDA come from awareness learning? (Section 5.3) (3) How do the learned awareness embeddings distribute in the representation space and how do they influence the team cooperation? (Section 5.4) (4) Do the awareness distributions have interpretable characteristics? (Section 5.5)

5.1 Environments

We choose two multi-agent benchmark environments: Level-based Foraging (LBF) [1] and StarCraft II micromangement (SMAC)[†] [28] as the testbed.

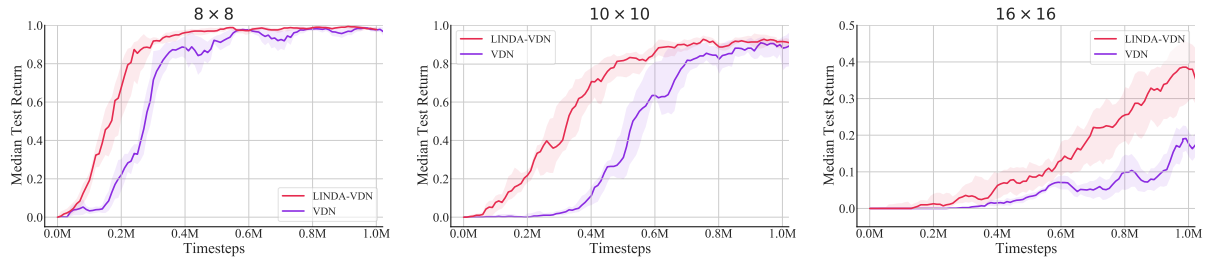


Figure 6 Test episodic return for LINDA-VDN and VDN on three different LBF configurations.

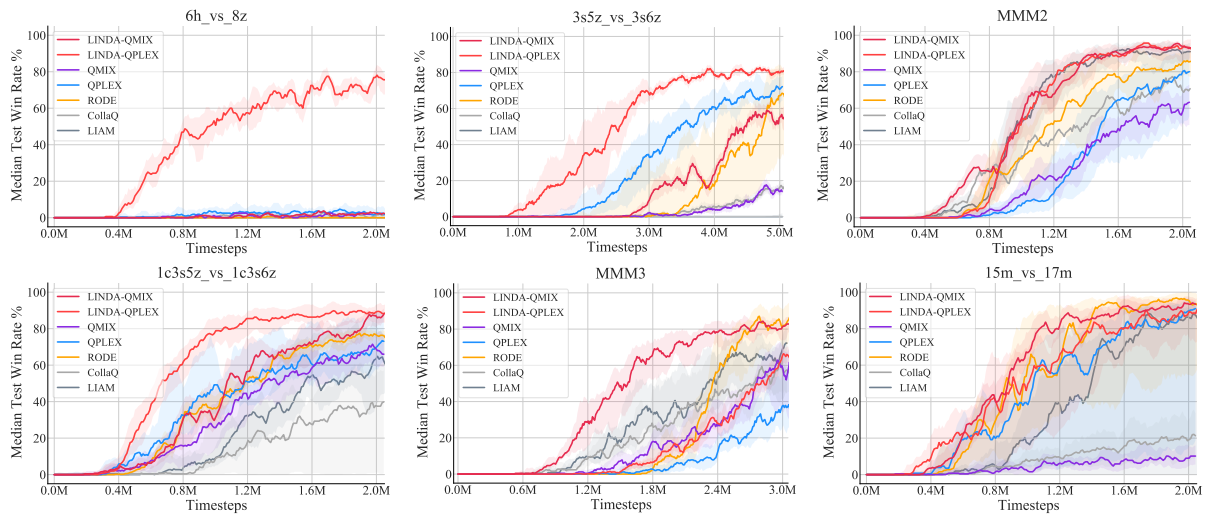


Figure 7 Test win rate for LINDA-QPLEX, LINDA-QMIX, and other baselines on six super hard SMAC maps.

Level-based foraging (LBF): LBF [1] is a grid-world game that focuses on the coordination of the agents (Figure 3). It consists of agents and foods initialized at different positions at the beginning of an episode. Both agents and foods are assigned random levels. The goal is for agents to collect foods to achieve the maximum team score. The food collection is constrained by the level. Food is successfully collected only when the sum of the levels of the agents who pick up the food simultaneously exceeds the level of the food. The action set of agents is composed of picking up action and movement in four directions: up, down, left, and right. Agents receive a team reward only when they successfully collect food, which means the environment has a sparse reward. Besides, the environment is partially observable, where the agents observe up to two grid cells in every direction. The agents’ observation includes the positions of other agents and foods in the visible range. For LBF, we test the algorithms in different sizes of exploration space, with 8×8 , 10×10 , 16×16 grid sizes, and all with 2 agents and 1 food.

StarCraft II micromangement (SMAC): SMAC [28] is an environment for collaborative multi-agent reinforcement learning based on Blizzard’s StarCraft II RTS game (Figure 4). It consists of a set of StarCraft II micro scenarios, each of which is a confrontation between two armies of units. The allied agents learn coordination to beat the enemy units controlled by the built-in game AI. It is a partially observable environment, where agents are only accessible to the status including positions and health of teammates and enemies in their limited field of view. For SMAC, we test the algorithms on maps of different difficulties and classify them as *easy*, *hard* and *super hard* (see Table C1 for detail).

[†]Our experiments are all based on the PYMARL framework which uses SC2.4.6, note that performance is not always comparable between versions.

5.2 Performance on Level-Based Foraging and StarCraft II

To study the effectiveness of LINDA framework, we try three different mixing networks in: (i) VDN [32], (ii) QMIX [27], (iii) QPLEX [36]. Their mixing network structures are from simple to complex, with increasing representational complexity. The methods applied with LINDA are denoted as LINDA-VDN, LINDA-QMIX, and LINDA-QPLEX, respectively. To show the effectiveness of LINDA applied methods, We compare with RODE [40], which is the state-of-the-art method, CollaQ [46], which uses manual local information decomposition, and LIAM [25], which uses autoencoders for opponent modelling.

Because VDN suffers from structural constraints and limited representation complexity [27], it tends to fail in complex environments like SMAC. We evaluate LINDA-VDN in LBF, which is simpler for agents

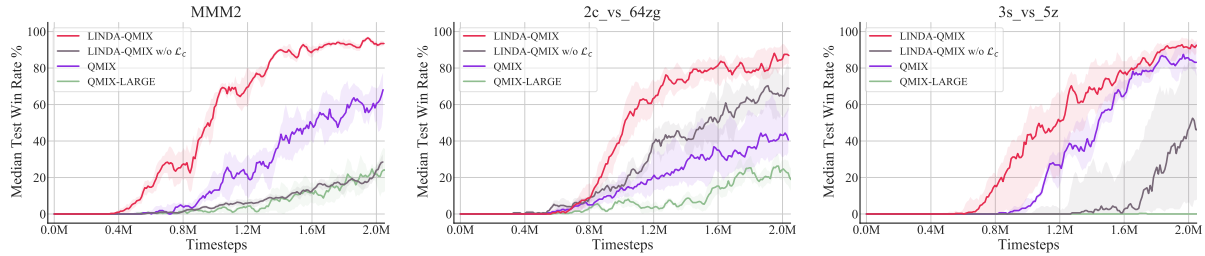


Figure 8 Ablation study for LINDA-QMIX and other ablation methods on three different SMAC maps.

to achieve goals. We configure three different grid sizes, 8×8 , 10×10 , and 16×16 , all with 2 agents and 1 food. For LINDA-QMIX and LINDA-QPLEX, we test them on multiple SMAC maps. For evaluation, we carry out each experiment with 5 random seeds, and the results are shown with a 95% confidence interval.

Figure 6 shows the learning curves of LINDA-VDN and VDN in different sizes of exploration space of the LBF environment. As the grid size grows, each agent needs to search a larger space to collaborate with others for food, and it takes more timesteps for agents to converge to a high reward. We show that the application of LINDA, i.e., LINDA-VDN, improves the performance of VDN. For smaller grid size, e.g. 8×8 , LINDA-VDN accelerates the learning speed and stabilizes the performance curve when converged. For larger grid sizes, the improvement is rather more significant.

In the SMAC environment, we test LINDA-QPLEX, LINDA-QMIX and other methods on six *super hard* maps, three *hard* maps and six *easy* maps. The details of the SMAC maps are described in Appendix C Table C1. From Figure 5, we can find LINDA-QMIX and LINDA-QPLEX outperform the vanilla QMIX and QPLEX, respectively, which indicates the effectiveness of LINDA for improving the coordination ability in complex sceneries. The superiority of QPLEX over other methods except for QPLEX-LINDA before 0.8M is for the improved network representation of QPLEX makes it good at some maps such as 3s5z_vs_3s6z. LIAM can solve the POMDP and improve the coordinate ability of QMIX in some way, and we find it is not competitive with LINDA. CollaQ only has a slight performance advantage over QMIX. We guess it is because CollaQ only uses local observation to capture the relationship between agents, which may not solve complex coordination problems.

We plot the averaged median test win rate across the six super hard maps in Figure 5. Compared with QMIX and QPLEX, the application of LINDA yields better learning performance and outperforms other methods. LINDA improves the collaboration for QMIX and QPLEX in all the six super hard maps, especially in tough scenarios. For example, in 6h_vs_8z, a map requiring strong coordination ability, all other methods failed, but LINDA-QPLEX achieved a high winning rate. In Appendix D, we further show the learning curves on the six easy scenarios and the hard maps. The LINDA framework still achieves slight performance improvement on most easy maps, where micro-tricks and cohesive collaboration are unnecessary. On the contrary, previous methods such as RODE perform worse than QMIX because the additional role learning module needs more samples to learn a successful strategy [40]. The overall experiments show that the LINDA framework is robust to both easy and tough scenarios and effectively enhances the learning efficiency, especially on challenging tasks.

5.3 Ablation Study

To understand the superior performance of LINDA, we carry out ablation studies to verify the contribution of awareness learning. To this end, we add the LINDA structure to QMIX without the mutual information

loss $\mathcal{L}_c(\theta_c)$ during training and denote it as *LINDA-QMIX w/o $\mathcal{L}_c(\theta_c)$* . Besides, to test whether the superiority of our method comes from the increase in the number of parameters, we also test QMIX with a similar number of parameters with LIDA-QMIX and denote it as *QMIX-LARGE*. As shown in Figure 8, we compare LINDA-QMIX with LINDA-QMIX w/o $\mathcal{L}_c(\theta_c)$, QMIX-LARGE, and QMIX on three SMAC maps: *MMM2*, *2c_vs_64zg* and *3s_vs_5z*. The experimental results indicate that the mutual information loss plays a significant role in enhancing learning performance. Besides, it also proves that larger networks will not definitely bring performance improvement, and the superiority of LINDA does not come from the larger networks.

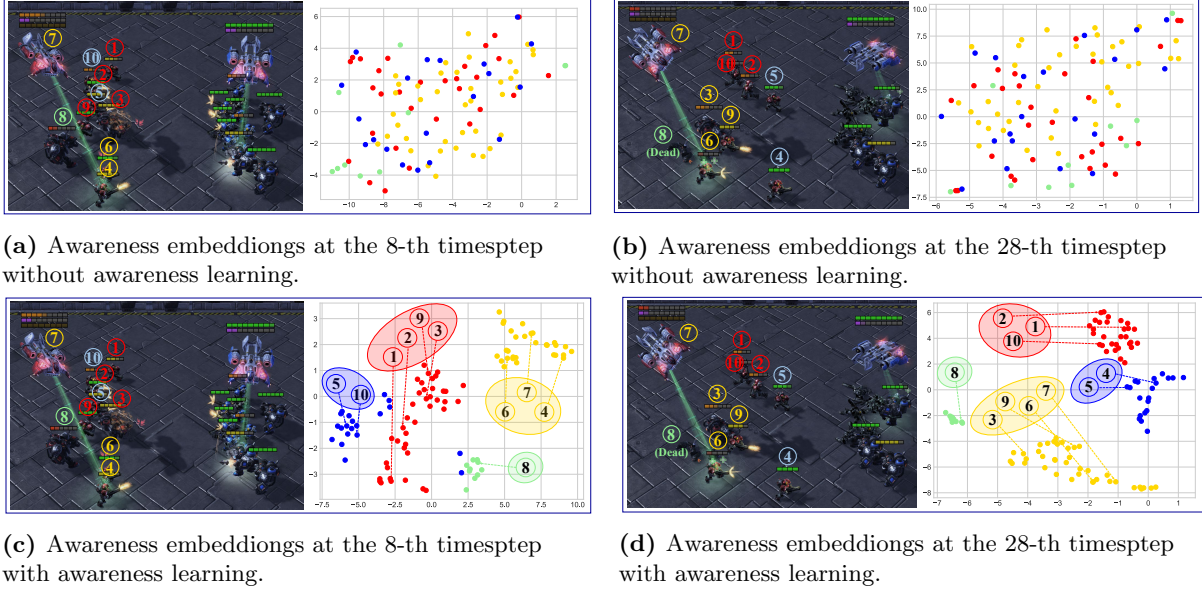


Figure 9 Visualization of the video frames and t-SNE projection of the representation space at the early (8-th) and the late (28-th) timesteps. Awareness embeddings of LINDA-QMIX w/o $\mathcal{L}_c(\theta_c)$ distributed are chaos in the representation space (see (a) and (b)). (c) and (d) reveals that awareness embeddings of LINDA-QMIX formed clusters and there exists correspondence between the groups formed in the game and in the awareness representation space. Besides, the groups dynamically changed at the early timestep (c) and the late timestep (d).

5.4 Awareness Embedding Representations

To further analyze the learned awareness in the representation space, we visualize the awareness embeddings on the SMAC map *MMM2*. *MMM2* is a heterogeneous scenario, in which 1 Medivac, 2 Marauders and 7 Marines face 1 Medivac, 3 Marauders and 8 Marines. In this task, different types of agents should cooperate well to fully exert the advantage of each unit type. We collect the awareness embeddings of all the agents at two different timesteps in one episode. Since there are 10 agents, each with 10 awareness embeddings for teammates, there are 100 embeddings in total at each timestep. We reduce the dimension of each awareness embedding by t-SNE [34] to show them in a 2-dimensional plane.

As shown in Figure 9a and 9b, the awareness embeddings generated by LINDA-QMIX w/o $\mathcal{L}_c(\theta_c)$, which is without awareness learning, distributed almost randomly in the representation space. Contrarily, with the proposed mutual information loss $\mathcal{L}_c(\theta_c)$, in Figure 9c and 9d, agents' awareness embeddings automatically form several clusters in the representation space. According to the positions of self-to-self awareness embeddings, we divide the agents into K groups $G = \{G_1, \dots, G_K\}$. We color the awareness embeddings by the group each agent belongs to. That is, for each group G_k , the representations $\{c_j^i \mid 1 \leq j \leq n, i \in G_k\}$ are painted the same color, where n is the number of agents. In the video frame of the same timestep, we see the correspondence between the agent groups formed in the game and in the awareness representation space. The agents in the same group tend to build similar awareness and achieve more cooperation. Besides, In the early timestep (Figure 9c) and the late timestep (Figure 9d), the formed groups dynamically changed and adapted according to the battle situation. As shown in Figure 9d, when the 7-th agent was dead, its awareness embeddings collapsed to a small region because the observation inputs of a dead agent are all zeros in the environment implementation. Other agents formed new groups and conducted new tactics for cooperation.

The phenomenon reveals the relationship between awareness representations and agents’ cooperation strategies. The aggregation of the learned awareness makes the group members share similar latent embeddings. Such consistency may encourage neural networks to output stable and consistent action value functions, thus reaching a consensus on the strategies among the group members and achieving better collaboration. Further, the pairwise relationship of the learned awareness implicitly forms groups among agents, which is similar to role learning [16, 31, 39, 40] in multi-agent systems. It indicates that pairwise awareness representations implicitly incorporate roles, and can be degraded into role-based methods by further processing the relationships of awareness among agents.

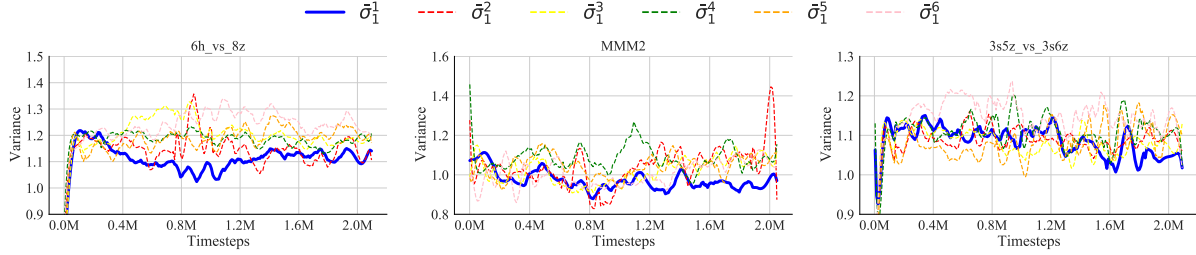


Figure 10 Visualization of the awareness variances. $\bar{\sigma}_1^i$ denotes the average of the awareness variance for agent 1. We take the first 6 agents in each SMAC map for concise visualization. The blue bold curve, representing the variance of the self-to-self awareness distribution, is nearly the lowest among all the curves.

5.5 Visualization of Awareness Distributions

We conduct visualization on the dynamic learning process of the awareness distributions. For each agent i , we visualize its average of the variance of the multivariate awareness distribution for agent 1, which is denoted as $\bar{\sigma}_1^i$. Figure 10 presents $\bar{\sigma}_1^i, 1 \leq i \leq 6$ (for the first 6 agents) over the course of training. The blue bold curve, which represents $\bar{\sigma}_1^1$, is nearly the lowest among all the curves during training. The variance of a self-to-self awareness distribution is relatively lower than the variance of other-to-self awareness distribution. It means that agents build more certain awareness for self than others, which is consistent with our intuition.

To further demonstrate whether a lower variance indicates a more precise awareness embedding, we visualize the variance and the difference of mean among the agents. To clearly show the relationship, we conduct visualization on a two-agent scenario, *2c_vs_64zg*. In Figure 11, the red and blue solid curve represents the average of the awareness variance, which is denoted as $\bar{\sigma}_1^1$ and $\bar{\sigma}_1^2$ respectively. The red dashed curve represents the L1-norm of the difference between the awareness mean, which is denoted as $\|\mu_1^2 - \mu_1^1\|$. We find that the peaks of $\|\mu_1^2 - \mu_1^1\|$ and $\bar{\sigma}_1^2$ are highly overlapping, especially after 1.2 millions of timesteps when the learning gradually converges. The peak of $\bar{\sigma}_1^2$ means that at that timestep, A2 (i.e., agent 2) is uncertain about the state of A1 probably because A1 is out of the view of A2. The peak of $\|\mu_1^2 - \mu_1^1\|$ means that A1 and A2 have a huge difference in their awareness for A1. The phenomenon of peak alignment indicates that the inference uncertainty is highly related to the proximity of awareness. The awareness variance may serve as an indicator for agents to estimate the confidence of awareness. It suggests that for further research, we can put different degrees of emphasis on different awareness according to the confidence. Higher-level modules such as the attention mechanism can be built based on awareness confidence.

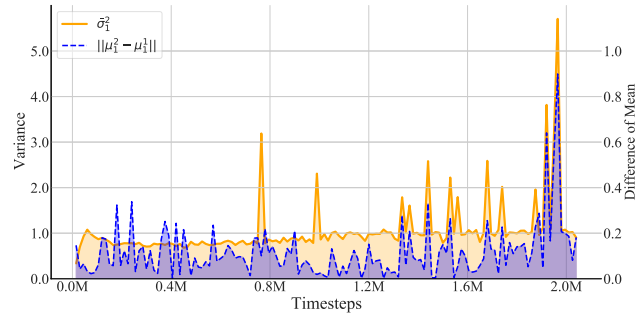


Figure 11 Visualization of awareness distributions on a 2-agent SMAC map *2c_vs_64zg*. The peaks of the curves $\|\mu_1^2 - \mu_1^1\|$ and $\bar{\sigma}_1^2$ are almost perfectly aligned.

6 Conclusion

We propose a novel framework, multi-agent Local Information Decomposition for Awareness of teammates (LINDA), which learns to build awareness for teammates on the local network to alleviate partial observability. We design an information-theoretic loss for awareness learning. LINDA is agnostic to specific algorithms and is flexibly applicable to existing MARL methods that follow the CTDE paradigm. We apply LINDA to three value-based MARL algorithms, and results show that LINDA makes a significant performance improvement, especially on super hard tasks in the SMAC benchmark environment. We also demonstrate the interpretability of the learned awareness distributions and show that LINDA forms awareness groups and promotes cooperation. Further research on building higher-level modules such as the attention mechanism based on awareness would be of interest.

Acknowledgements This work was supported by National Natural Science Foundation of China (Under Grant No. 61773198).

References

- 1 Stefano V Albrecht and Peter Stone. Reasoning about hypothetical agent behaviours and their parameters. In *the 16th Conference on Autonomous Agents and MultiAgent Systems*, pages 547–555, 2017.
- 2 Stefano V Albrecht and Peter Stone. Autonomous agents modelling other agents: A comprehensive survey and open problems. *Artificial Intelligence*, 258:66–95, 2018.
- 3 Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. In *the 5th International Conference on Learning Representations*, 2017.
- 4 Yongcan Cao, Wenwu Yu, Wei Ren, and Guanrong Chen. An overview of recent progress in the study of distributed multi-agent coordination. *IEEE Transactions on Industrial Informatics*, 9(1):427–438, 2012.
- 5 Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1724–1734, 2014.
- 6 Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. Shared experience actor-critic for multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, 2020.
- 7 Haoyan Cui and Zhen Zhang. A cooperative multi-agent reinforcement learning method based on coordination degree. *IEEE Access*, 9:123805–123814, 2021.
- 8 Jakob N. Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *the 32nd AAAI Conference on Artificial Intelligence*, pages 2974–2982, 2018.
- 9 Sven Gronauer and Klaus Diepold. Multi-agent deep reinforcement learning: a survey. *Artificial Intelligence Review*, pages 1–49, 2021.
- 10 He He and Jordan L. Boyd-Graber. Opponent modeling in deep reinforcement learning. In *the 33rd International Conference on Machine Learning*, pages 1804–1813, 2016.
- 11 Zhang-Wei Hong, Shih-Yang Su, Tzu-Yun Shann, Yi-Hsiang Chang, and Chun-Yi Lee. A deep policy inference q-network for multi-agent systems. In *the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 1388–1396, 2018.
- 12 Siyi Hu, Fengda Zhu, Xiaojun Chang, and Xiaodan Liang. Updet: Universal multi-agent reinforcement learning via policy decoupling with transformers. *arXiv preprint arXiv:2101.08001*, 2021.
- 13 Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In *the 36th International Conference on Machine Learning*, pages 2961–2970, 2019.
- 14 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 15 Landon Kraemer and Bikramjit Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190:82–94, 2016.
- 16 Kemas M Lhaksmana, Yohei Murakami, and Toru Ishida. Role-based modeling for designing agent behavior in self-organizing multi-agent systems. *International Journal of Software Engineering and Knowledge Engineering*, 28(01):79–96, 2018.
- 17 Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*, pages 6379–6390, 2017.
- 18 Xueguang Lyu, Yuchen Xiao, Brett Daley, and Christopher Amato. the contrasting centralized and decentralized critics in multi-agent reinforcement learning. In *the 20th International Conference on Autonomous Agents and Multiagent Systems*, pages 844–852, 2021.
- 19 Ann Nowé, Peter Vrancx, and Yann-Michaël De Hauwere. Game theory and multi-agent reinforcement learning. *Reinforcement Learning*, 12:441–470, 2012.
- 20 Frans A. Oliehoek and Christopher Amato. *A Concise Introduction to Decentralized POMDPs*. Springer, 2016.
- 21 Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. Optimal and approximate q-value functions for decentralized pomdps. *Journal of Artificial Intelligence Research*, 32:289–353, 2008.
- 22 Afshin OroojlooyJadid and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement learning. *arXiv preprint arXiv:1908.03963*, 2019.
- 23 Praveen Palanisamy. Multi-agent connected autonomous driving using deep reinforcement learning. In *International Joint Conference on Neural Networks*, pages 1–7, 2020.
- 24 Georgios Papoudakis and Stefano V Albrecht. Variational autoencoders for opponent modeling in multi-agent systems. *arXiv preprint arXiv:2001.10829*, 2020.
- 25 Georgios Papoudakis, Filippos Christianos, and Stefano V Albrecht. Local information opponent modelling using variational autoencoders. *arXiv preprint arXiv:2006.09447*, 2020.
- 26 Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted QMIX: expanding monotonic value function factorisation for deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, 2020.
- 27 Tabish Rashid, Mikayel Samvelyan, Christian Schröder de Witt, Gregory Farquhar, Jakob N. Foerster, and Shimon Whiteson. QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning. In *the 35th International Conference on Machine Learning*, pages 4292–4301, 2018.

- 28 Mikayel Samvelyan, Tabish Rashid, Christian Schröder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob N. Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. In *the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pages 2186–2188, 2019.
- 29 Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Hostallero, and Yung Yi. QTRAN: learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *the 36th International Conference on Machine Learning*, pages 5887–5896, 2019.
- 30 Peter Stone, Gal A Kaminka, Sarit Kraus, and Jeffrey S Rosenschein. Ad hoc autonomous agent teams: Collaboration without pre-coordination. In *Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010.
- 31 Peter Stone and Manuela Veloso. Task decomposition, dynamic role assignment, and low-bandwidth communication for real-time strategic teamwork. *Artificial Intelligence*, 110(2):241–273, 1999.
- 32 Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Flores Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 2085–2087, 2018.
- 33 Karl Tuyls and Gerhard Weiss. Multiagent learning: Basics, challenges, and prospects. *Ai Magazine*, 33(3):41–41, 2012.
- 34 Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(11):2579–2605, 2008.
- 35 Oriol Vinyals, Timo Ewalds, Sergey Bartunov, P. Georgiev, A. S. Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, J. Agapiou, Julian Schrittwieser, John Quan, Stephen Gaffney, S. Petersen, K. Simonyan, T. Schaul, H. V. Hasselt, D. Silver, T. Lillicrap, Kevin Calderone, Paul Keet, Anthony Brunasso, D. Lawrence, Anders Ekelmo, J. Repp, and Rodney Tsing. Starcraft ii: A new challenge for reinforcement learning. arXiv preprint arXiv:/1708.04782, 2017.
- 36 Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. QPLEX: duplex dueling multi-agent q-learning. In *the 9th International Conference on Learning Representations*, 2021.
- 37 Jianhong Wang, Jinxin Wang, Yuan Zhang, Yunjie Gu, and Tae-Kyun Kim. Shaq: Incorporating shapley value theory into q-learning for multi-agent reinforcement learning. arXiv preprint arXiv:2105.15013, 2021.
- 38 Jianhong Wang, Yuan Zhang, Tae-Kyun Kim, and Yunjie Gu. Shapley q-value: A local reward approach to solve global reward games. In *the 34th AAAI Conference on Artificial Intelligence*, pages 7285–7292, 2020.
- 39 Tonghan Wang, Heng Dong, Victor R. Lesser, and Chongjie Zhang. ROMA: multi-agent reinforcement learning with emergent roles. In *the 37th International Conference on Machine Learning*, pages 9876–9886, 2020.
- 40 Tonghan Wang, Tarun Gupta, Anuj Mahajan, Bei Peng, Shimon Whiteson, and Chongjie Zhang. RODE: learning roles to decompose multi-agent tasks. In *the 9th International Conference on Learning Representations*, 2021.
- 41 Tonghan Wang, Jianhao Wang, Chongyi Zheng, and Chongjie Zhang. Learning nearly decomposable value functions via communication minimization. In *the 8th International Conference on Learning Representations*, 2020.
- 42 Tonghan Wang, Liang Zeng, Weijun Dong, Qianlan Yang, Yang Yu, and Chongjie Zhang. Context-aware sparse deep coordination graphs. arXiv preprint arXiv:2106.02886, 2021.
- 43 Weixun Wang, Tianpei Yang, Yong Liu, Jianye Hao, Xiaotian Hao, Yujing Hu, Yingfeng Chen, Changjie Fan, and Yang Gao. Action semantics network: Considering the effects of actions in multiagent systems. In *the 8th International Conference on Learning Representations*, 2020.
- 44 Bin Wu. Hierarchical macro strategy model for MOBA game AI. In *the 33rd AAAI Conference on Artificial Intelligence*, pages 1206–1213, 2019.
- 45 Annie Xie, Dylan P. Losey, Ryan Tolsma, Chelsea Finn, and Dorsa Sadigh. Learning latent representations to influence multi-agent interaction. In *the 4th Conference on Robot Learning*, pages 575–588, 2020.
- 46 Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt Keutzer, Joseph E Gonzalez, and Yuandong Tian. Multi-agent collaboration via reward attribution decomposition. arXiv preprint arXiv:2010.08531, 2020.
- 47 Ming Zhou, Jun Luo, Julian Vilella, Yaodong Yang, David Rusu, Jiayu Miao, Weinan Zhang, Montgomery Alban, Iman Fadakar, Zheng Chen, et al. Smarts: Scalable multi-agent reinforcement learning training school for autonomous driving. arXiv preprint arXiv:2010.09776, 2020.

Appendix A Mathematical Derivation

We maximize the mutual information between agent i 's awareness for agent j \mathbf{c}_j^i and agent j 's trajectory τ^j conditioned on agent i 's local trajectory τ^i :

$$\begin{aligned}
 I(\mathbf{c}_j^i; \tau^j | \tau^i) &= \mathbb{E}_{\tau, \mathbf{c}_j^i} \left[\log \frac{p(\mathbf{c}_j^i | \tau^i, \tau^j)}{p(\mathbf{c}_j^i | \tau^i)} \right] \\
 &= \mathbb{E}_{\tau, \mathbf{c}_j^i} \left[\log \frac{q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j)}{p(\mathbf{c}_j^i | \tau^i)} \right] + \mathbb{E}_\tau \left[D_{KL} \left(p(\mathbf{c}_j^i | \tau^i, \tau^j) \parallel q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j) \right) \right] \\
 &\geq \mathbb{E}_{\tau, \mathbf{c}_j^i} \left[\log \frac{q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j)}{p(\mathbf{c}_j^i | \tau^i)} \right] \\
 &= \mathbb{E}_{\tau, \mathbf{c}_j^i} \left[\log q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j) \right] - \mathbb{E}_{\tau, \mathbf{c}_j^i} \left[\log p(\mathbf{c}_j^i | \tau^i) \right] \\
 &= \mathbb{E}_{\tau, \mathbf{c}_j^i} \left[\log q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j) \right] + \mathbb{E}_{\tau^i} \left[H(\mathbf{c}_j^i | \tau^i) \right] \\
 &= \mathbb{E}_\tau \left[\int p(\mathbf{c}_j^i | \tau^i, \tau^j) \log q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j) d\mathbf{c}_j^i \right] + \mathbb{E}_{\tau^i} \left[H(\mathbf{c}_j^i | \tau^i) \right].
 \end{aligned}$$

The awareness encoder of agent i is conditioned on the local trajectory τ^i . Thus, given τ^i , the awareness distribution $p(\mathbf{c}_j^i)$ is independent of τ^j . And we have

$$I(\mathbf{c}_j^i; \tau^j | \tau^i) \geq -\mathbb{E}_\tau \left[\mathcal{CE} \left[p(\mathbf{c}_j^i | \tau^i) \parallel q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j) \right] \right] + \mathbb{E}_{\tau^i} \left[H(\mathbf{c}_j^i | \tau^i) \right], \quad (\text{A1})$$

where \mathcal{CE} is the Cross-Entropy operator. We use a replay buffer \mathcal{D} in practice. For agent i 's n awareness distributions $\mathbf{c}_1^i, \dots, \mathbf{c}_n^i$, we can derive the minimization objective:

$$\mathcal{L}_c(\theta_c^i) = \sum_{j=1}^n \mathbb{E}_{\tau \sim \mathcal{D}} \left[D_{KL} \left[p(\mathbf{c}_j^i | \tau^i) \parallel q_\xi(\mathbf{c}_j^i | \tau^i, \tau^j) \right] \right]. \quad (\text{A2})$$

Appendix B Architecture, Hyperparameters, and Infrastructure

In this paper, we base our framework on three value factorization based methods, VDN [32], QMIX [27], and QPLEX [36]. Following the CTDE paradigm, each agent has an individual neural network to approximate its local utility, which is fed into a mixing network to estimate the global action value during training.

We base our implementations of LINDA-VDN, LINDA-QMIX and LINDA-QPLEX in the PyMARL framework and use its default mixing network structure and the same hyper-parameter setting with VDN [†], QMIX [§] and QPLEX [¶] from the original paper, respectively. The architecture of all agent networks is a DRQN with a recurrent layer comprised of a GRU with a 64-dimensional hidden state, with a fully-connected layer before and after, VDN only sum up all the DRON outputs. QMIX's mixing network consists of a single hidden layer of 32 units, utilizing an ELU non-linearity. The hyper-networks are then sized to produce weights of appropriate size. The mixing network of QPLEX is more complex, which consists of two main components as follows: (i) an Individual Action-Value Function for each agent, and (ii) a Duplex Dueling component that composes individual action-value functions into a joint action-value function under the advantage-based IGM constraint. The awareness encoder adopts a 64-dimensional hidden layer with LeakyReLU activation and outputs a 3-dimensional multivariate Gaussian distribution for each agent. The posterior estimator q_ξ also uses a 64-dimensional hidden layer with LeakyReLU activation. The awareness embeddings are sampled from the corresponding Gaussian distributions and concatenated with the trajectory to be fed into the local utility network.

Appendix C The SMAC Environment

We test the algorithms on different scenarios in SMAC. The detailed configurations of the scenarios are shown in Table C1.

Appendix D Additional Experimental Results

Performance on more maps. We mainly benchmark our method on the StarCraft II unit micromangement tasks. To test the generation of LINDA, we evaluate LINDA-QMIX and LINDA-QPLEX on other easy and hard maps. The additional results are shown in Figure D2 and D1. We find that in easy scenarios where micro-tricks and cohesive collaboration are unnecessary, the application of LINDA still brings slight performance improvement. On the contrary, previous methods such as RODE perform worse than QMIX because the additional role learning module needs more samples to learn a successful strategy [40]. The experimental results show that LINDA is robust to both easy and hard tasks.

[†]VDN code: <https://github.com/oxwhirl/pymarl>

[§]QMIX code: <https://github.com/oxwhirl/pymarl>

[¶]QPLEX code: <https://github.com/wjh720/QPLEX>

Table C1 SMAC challenges

Map Name	Ally Units	Enemy Units	Type	Challenge
2s_vs_1sc	2 Stalkers	1 Spine, 1 Crawler	Asymmetric, Heterogeneous	Easy
2s3z	2 Stalkers, 3 Zealots	2 Stalkers, 3 Zealots	Symmetric, Heterogeneous	Easy
3s5z	3 Stalkers, 5 Zealots	3 Stalkers, 5 Zealots	Symmetric, Heterogeneous	Easy
1c3s5z	1 Colossus, 3 Stalkers, 5 Zealots	1 Colossus, 3 Stalkers , 5 Zealots	Symmetric, Heterogeneous	Easy
MMM	1 Medivac, 2 Marauders, 7 Marines	1 Medivac, 2 Marauders , 7 Marines	Asymmetric, Heterogeneous	Easy
10m_vs_11m	10 Marines	11 Marines	Asymmetric, Homogeneous	Easy
5m_vs_6m	5 Marines	6 Marines	Asymmetric, Homogeneous	Hard
3s_vs_5z	3 Stalkers	5 Zealots	Asymmetric, Homogeneous	Hard
2c_vs_64zg	2 Colossi	64 Zerglings	Asymmetric, Homogeneous	Hard
15m_vs_17m	15 Marines	17 Marines	Asymmetric, Homogeneous	Super Hard
6h_vs_8z	6 Hydralisks	8 Zealots	Asymmetric, Homogeneous	Super Hard
3s5z_vs_3s6z	3 Stalkers, 5 Zealots	3 Stalkers , 6 Zealots	Asymmetric, Heterogeneous	Super Hard
1c3s5z_vs_1c3s6z	1 Colossus, 3 Stalkers, 5 Zealots	1 Colossus, 3 Stalkers , 6 Zealots	Asymmetric, Heterogeneous	Super Hard
MMM2	1 Medivac, 2 Marauders, 7 Marines	1 Medivac, 2 Marauders , 8 Marines	Asymmetric, Heterogeneous	Super Hard
MMM3	1 Medivac, 2 Marauders, 7 Marines	1 Medivac, 2 Marauders , 9 Marines	Asymmetric, Heterogeneous	Super Hard

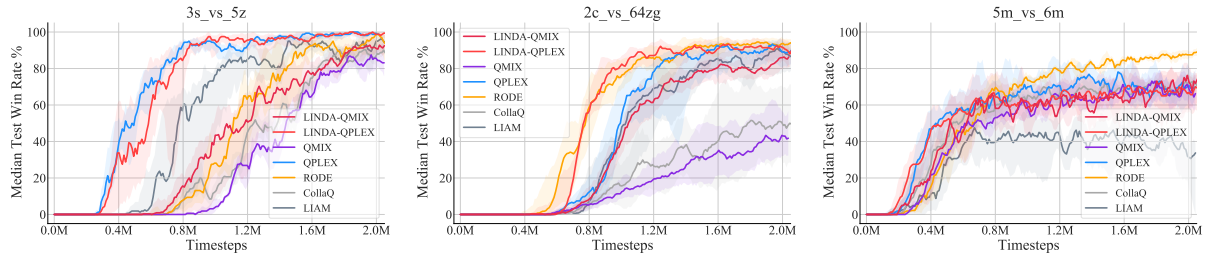


Figure D1 Test win rate for LINDA-QPLEX, LINDA-QMIX, and other baselines on three hard SMAC maps.

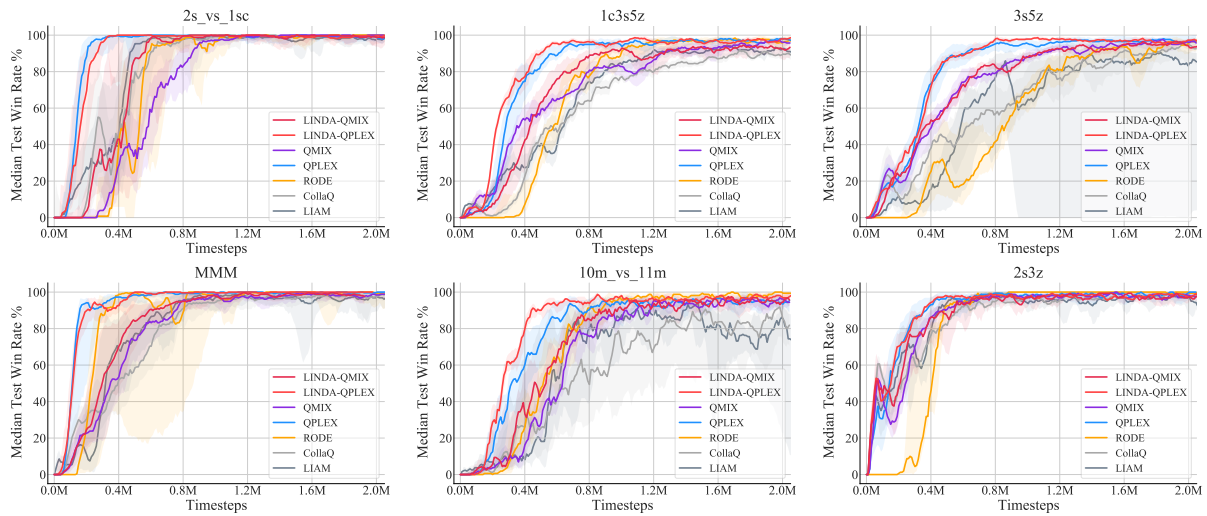


Figure D2 Test win rate for LINDA-QPLEX, LINDA-QMIX, and other baselines on six easy SMAC maps.