

Influence of Team Interactions on Multi-Robot Cooperation: A Relational Network Perspective

Yasin Findik^{1*}, Hamid Osooli^{1*}, Paul Robinette², Kshitij Jerath³, and S. Reza Ahmadzadeh¹

Abstract—Relational networks within a team play a critical role in the performance of many real-world multi-robot systems. To successfully accomplish tasks that require cooperation and coordination, different agents (e.g., robots) necessitate different priorities based on their positioning within the team. Yet, many of the existing multi-robot cooperation algorithms regard agents as interchangeable and lack a mechanism to guide the type of cooperation strategy the agents should exhibit. To account for the team structure in cooperative tasks, we propose a novel algorithm that uses a relational network comprising inter-agent relationships to prioritize certain agents over others. Through appropriate design of the team’s relational network, we can guide the cooperation strategy, resulting in the emergence of new behaviors that accomplish the specified task. We conducted six experiments in a multi-robot setting with a cooperative task. Our results demonstrate that the proposed method can effectively influence the type of solution that the algorithm converges to by specifying the relationships between the agents, making it a promising approach for tasks that require cooperation among agents with a specified team structure.

I. INTRODUCTION

Effective coordination and cooperation among multiple agents (e.g., robots) is imperative to accomplishing collective or individual objectives [1] in a range of applications, including autonomous driving [2], [3], search & rescue [4], [5], and logistics & transportation [6]. In recent years, deep learning techniques within multi-agent reinforcement learning (MARL) approaches have shown promising results in addressing the cooperation challenges in multi-agent environments [7], [8], [9]. One of the most prominent paradigms in MARL for cooperative tasks is the centralized training with decentralized execution (CTDE) [10].

State-of-the-art CTDE approaches developed in recent years are widely recognized for tackling coordination problems [11], [12], [13], [14], [15]. These CTDE approaches offered no control over which coordination strategy the algorithms yield, which is because of the absence of a mechanism for specifying the particular type of agent-agent relations. To address this issue, some algorithms borrowed and applied the idea of coordination graph [16] to share action values among the agents using deep learning [17]. This approach requires the algorithm to learn how to weigh and share action

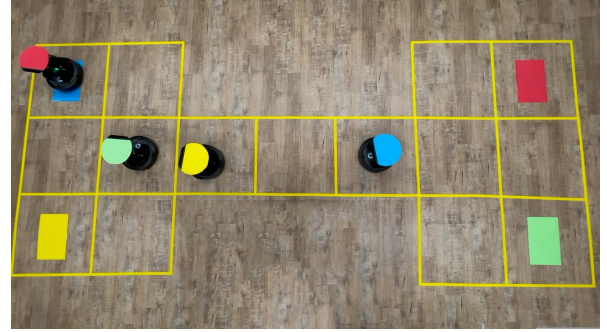


Fig. 1: Our experimental setup where four Turtlebot4 robots learn to navigate the *Switch* environment cooperatively so each one can reach the goal location matching their color.

values from an undirected graph, which does not allow for the influence of certain behaviors among others. We posit that one way to design cooperation strategies that account for team structure is to define a relational network represented with a directed graph, which captures the importance and priorities agents put on each other. We propose a novel CTDE-based algorithm that exerts influence on the type of cooperation strategy by utilizing a relational network that determines the relationships in the team. Our proposed algorithm utilizes this relational network to prioritize certain agents over others, enabling agents to uncover new behaviors that successfully accomplish the task while maintaining the team structure.

We conducted six experiments in both simulation and real-world settings, using our algorithm in a multi-robot task in which coordination among agents was influenced through different relational networks. We demonstrated that our proposed method can effectively steer the algorithm to converge to a specific team behavior governed by the relational network among agents. Unlike most methods, the proposed approach enables effective cooperation in a multi-agent team and allows for the emergence of different behaviors by leveraging relational networks.

II. RELATED WORK

Multi-agent reinforcement learning (MARL) in cooperative settings has become an active area of research in recent years, with various approaches being explored for successfully training agents to cooperate towards a shared goal. One such approach is fully decentralized learning, where each agent learns its own policy independently. In these approaches, the cooperative behavior emerges as the learned policies are employed in the environment. For example,

*indicates equal contribution

¹ PeARL Lab, Richard Miner School of Computer and Information Sciences, University of Massachusetts Lowell, MA, USA {yasin.findik, hamidosooli}@student.uml.edu, reza@cs.uml.edu

² Department of Electrical and Computer Engineering, University of Massachusetts Lowell, MA, USA paul_robinette@uml.edu

³ Department of Mechanical Engineering, University of Massachusetts Lowell, MA, USA kshitij.jerath@uml.edu

Independent Q-learning [18] learns a separate action-value table for each agent via Q-learning [19]. Since tabular Q-learning is limited in tackling high-dimensional state and action space, the Independent Q-learning framework was later extended with function approximation [20]. However, independent learning in multi-agent settings suffers from the non-stationarity of the problem caused by other agents' actions from the perspective of one agent. Because the Markov property is no longer valid in non-stationary environments, the convergence of Q-learning based decentralized algorithms cannot be guaranteed [21]. Later, the introduction of fully decentralized algorithms with networked agents [22], which could communicate and share their observations and actions, resulted in a method that could tackle the non-stationarity problem. In frameworks with networked agents, an undirected graph governs the connections between agents. Similarly, we also leverage a graph to represent some form of rapport between the agents. However, unlike frameworks that focus on the use of a communication channel [22], our graph symbolizes and establishes relations between agents, effectively creating a relational network that directly affects the team behavior.

Another approach to MARL in cooperative settings is fully centralized learning, where all agents share a single controller and learn a joint policy or value function together [23]. However, this approach is computationally expensive and can become intractable since the observation and action space grow exponentially with the number of agents. To address the limitations of both fully centralized and fully decentralized learning in cooperative MARL settings, a novel paradigm known as centralized training with decentralized execution (CTDE) [10] has been proposed. In this approach, individual agents carry out their actions while a centralized mechanism integrates their strategies. Both policy-based and value-based methods have been successfully applied in implementing the CTDE paradigm. During training, policy-based methods such as Multi-Agent Deep Deterministic Policy Gradient (MADDPG) [14] and Multi-Agent Proximal Policy Optimization (MAPPO) [15] allow the critic to consider global observations of all agents, while value-based techniques such as Value Decomposition Networks (VDN) [11], QMIX [12], and QTRAN [13] improve upon Q-Learning [19] by incorporating a centralized function that calculates the joint Q-value from the individual action-values of each agent.

From fully decentralized paradigm to CTDE, previous approaches usually have focused on converging to an optimal solution for solving cooperation problems in MARL. However, the ability of a multi-robot team to arrive at an optimal or sub-optimal solution is strongly influenced by the underlying team interactions [24], [25]. Unlike [25], our approach does not rely on reward sharing among agents to alter their individual rewards based on the relationship. Instead, our method entails modifying the way agents contribute to the team reward, which is determined by the relational network. Ultimately, agents still receive the same individual reward provided by the environment.

To better understand these effects, we introduce the inte-

gration of the relational network technique, which effectively captures such team interactions and rewards, into the CTDE framework. We choose to explore and study this idea with VDN, as it is a simple yet effective CTDE approach for learning cooperative behaviors. However, our proposed notion of relational networks can be applied to other CTDE approaches as well.

III. METHODOLOGY

A. Markov Decision Process

We characterized Decentralized Markov Decision Process as a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, n, \gamma \rangle$ where $s \in \mathcal{S}$ indicates the true state of the environment, the joint set of individual actions and rewards are represented by $\mathcal{A} := \{a_1, a_2, \dots, a_n\}$, $\mathcal{R} := \{r_1, r_2, \dots, r_n\}$, respectively, $\mathcal{T}(s, A, s') : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto [1, 0]$ is the dynamics function defining the transition probability, n is the the number of agents, and $\gamma \in [0, 1)$ is the discount factor.

B. Value Decomposition Networks

In this section, we provide a brief background on the Value-Decomposition Networks (VDN) [11] upon which we built our proposed method. VDN falls under the category of CTDE approaches. It addresses the non-stationarity problem of decentralized learning effectively through centralized training and it addresses the scalability issue of centralized learning through decentralized execution.

VDN maintains a separate action-value function for each agent $i \in \{1, \dots, n\}$, denoted as $Q_i(s, a_i) = \mathbb{E}[G | S = s, A = a_i]$ and approximated as $\hat{Q}_i(s, a_i, w_i)$ where w_i is the weight vector and G is the return. VDN combines these individual Q_i values to obtain the central action value Q_{tot} , as

$$Q_{tot} = \sum_{i=1}^n \hat{Q}_i(s, a_i, w_i). \quad (1)$$

The ultimate goal is to minimize the Temporal Difference (TD) [26] error:

$$e_{TD} = r_{team} + \gamma \max_{u'} (Q_{tot}(s', u')) - Q_{tot}(s, u), \quad (2)$$

where u is the joint action of the agents and r_{team} is the sum of agents' rewards with uniform weights. The algorithm minimizes the TD error by backpropagating e_{TD}^2 to each agent's model. This process allows the agents to coordinate their actions towards maximizing the team reward, which is of utmost importance in most cooperative settings. The central aspect of CTDE paradigm becomes evident during training, where the agent networks are trained with a centralized Q_{tot} that is computed from the sum of individual Q_i values. Whereas, during execution, each agent's actions are determined by their own neural network, rendering the execution decentralized.

C. Graph-based Value Decomposition Network

Most MARL problems typically have multiple solutions with different levels of optimality. Although VDN and several other approaches do not provide a guarantee of

finding the global optimal solution for every possible MARL environment or task, in general, they aim to maximize team rewards by improving coordination and cooperation among agents, and converge to one of several possible solutions (potentially the global optimal solution). However, if there are multiple (optimal) solutions, the randomness of agent's exploration determines the solution to which the algorithm converges. To be able to influence the team behavior and have more control over the dominant solution, we introduce Graph-based VDN (G-VDN). Similar to VDN, G-VDN utilizes an individual action-value function Q_i approximator for each agent and a summation function that aggregates the distinct Q_i values. The difference between the two, however, lies in the approach they use to aggregate team rewards. In VDN, every agent contributes their individual rewards to the team, uniformly. On the other hand, G-VDN aggregates individual rewards by considering the relational network between agents in computing joint rewards, thereby shaping the converged behavior by taking inter-agent relationships into account.

Crucially, we modify the MDP to additionally include the relational network which is represented by a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$, where each agent $i \in \{1, \dots, n\}$ is a vertex v_i , \mathcal{E} is the set of directed edges e_{ij} directed from v_i to v_j , and the weight of the edges is represented by the matrix \mathcal{W} with elements $\omega_{ij} \in [0, 1]$ for each edge. The direction of the edges and the associated weights signify the relationship between the agents, with an edge e_{ij} directed from i to j indicating that agent i cares about or is vested in the outcomes for agent j . Based on the modified MDP represented as a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \mathcal{G}, n, \gamma \rangle$, we reformulate (2) by introducing the aggregated reward, such that r_{team} is the weighted sum of agents' reward, defined as follows:

$$r_{\text{team}} = \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{E}_i} \omega_{ij} r_j, \quad (3)$$

where \mathcal{E}_i denotes the set of vertex indices that have an edge directed from v_i , and r_j is the reward of the agent represented by v_j .

In VDN, the main goal of each agent i is to learn a policy that maximizes the overall team performance. Then, it is utilized to determine the best action for agent a_i for a specific state s . This can be expressed as follows:

$$a_i = \pi_i^{\text{VDN}}(s) = \max_a (Q_i(s, a, w_i)),$$

where the learned policy of agent i is denoted as π_i , and Q_i represents its action-value function. This notation is applied to all agents in the following form:

$$u = \pi^{\text{VDN}}(s) = \max_u (Q(s, u, w)),$$

where π represents a collection of policies, and u is the joint actions of the agents. The formulation of G-VDN can be stated as follows:

$$u = \pi^{\text{G-VDN}}(s|\mathcal{G}) = \max_u (Q(s, u, w|\mathcal{G})),$$

where the policy is conditioned by the relational network \mathcal{G} .

In cases where \mathcal{G} represents a self-interest relational network, it is worth mentioning that the policies for G-VDN and VDN become equivalent. This equivalence arises because, in VDN, each agent's reward is contributed equally to the overall team reward.

$$\pi^{\text{G-VDN}}(s|\mathcal{G}) \equiv \pi^{\text{VDN}}(s),$$

where $\mathcal{G} \in \{\text{a, d, h}\}$ networks depicted in Fig. 2. Overall, G-VDN allows for the shaping and influencing of cooperation strategies among agents, whereas VDN leads to agents randomly converging to one of multiple solutions without any control over the team behavior.

D. Training Phase

VDN employs a Deep Q-Network (DQN) [9] to represent each agent's Q-function, enabling it to handle high dimensional state and action spaces in multi-agent systems. Yet, training a DQN can be challenging due to instability and divergence caused by updating the Q-network parameters in each step, violating the independently and identically distributed (i.i.d) data points assumption. To overcome these challenges, Mnih et al. [9] also introduce several techniques such as experience replay and fixed Q-target networks, which have now become standard in many deep reinforcement learning algorithms.

Similar to VDN, we (i) maintain a replay memory to store the experiences, and randomly replay these experiences m times during each training step, (ii) utilize two DQNs for each agent: a prediction network, and a fixed target network which is the prediction network from a previous iteration. The training process of the networks entails three main components: generating episodes, updating the target network, and training the prediction network.

Generating Episode: An episode is generated by using the prediction network of each agent with ε -greedy approach. A large number of transitions experienced by each agent during the episode are stored in the replay memory. Each transition is a tuple $\langle S, A, R, S' \rangle$, where S is a set of current states of the agents $\{s_1, s_2, \dots, s_n\}$, A contains the actions taken by each agent $\{a_1, a_2, \dots, a_n\}$, R are the rewards received by the agents $\{r_1, r_2, \dots, r_n\}$ and S' keeps the next state of the agents $\{s'_1, s'_2, \dots, s'_n\}$. The size of the memory is defined as a hyperparameter, and when it becomes full, the oldest transitions are replaced with new incoming ones.

Target Network Update: We use fixed target networks to increase the stability of agents' prediction network and prevent them from diverging. After a specified number of episodes have been generated, the weights of each agent's target network are updated with those of the corresponding agent's prediction network. The idea of incorporating target networks is akin to having a calculated plan for achieving success in a game, and persisting with it until a superior one emerges, rather than altering one's strategy after each move. By allowing the network to consider multiple moves instead of continuously adjusting its parameters, it has the

potential to develop a stronger model before being utilized for decision-making [9].

Prediction Network Training: The prediction networks are trained each time an episode is generated. To train the prediction networks, a batch of transitions of size b is sampled from memory. TD error is calculated using (2) for each instance in b and the team reward used in the TD error is computed by utilizing (3). Lastly, the model parameters are updated using an optimizer in the direction that reduces the TD error. This step that includes sampling batches from the memory and updating the networks is repeated m times per episode, where m is a hyperparameter.

IV. EXPERIMENTS

A. Environment

To evaluate the effectiveness of the proposed approach in influencing and steering agents' behaviors, we applied G-VDN to the *Switch* environment [27] in three scenarios with varying numbers of agents both in simulation and real-world. The accompanying video shows the execution of the experiments in real-world¹. As shown in Fig. 3(a), the *Switch* environment consists of a grid-world that features two areas connected through a narrow passageway. The *Switch* environment in its original size accommodates a maximum of four agents marked with different colors that are trying to reach their respective goals, marked in boxes of the same color as each agent. Each agent's goal cell is positioned on the opposite side of its initial location, posing a challenge given that only one agent can pass through the passageway at a time. Agents can take one of the five actions for each time step: move left, right, up, down, or stay still. Upon reaching their goal, they are awarded +5 points. The agents are penalized with -0.1 points for each time step spent away from their goal. The episode concludes once all agents have arrived at their respective goal locations or after a maximum of 50 time steps. To successfully reach their goal locations, the agents must cooperate to avoid blocking each other's paths while maximizing total team rewards. We chose the *Switch* environment because firstly, learning an optimal behavior requires cooperation among the agents. Secondly, the presence of the bridge allows for easily observable and measurable agents' behaviors and facilitates the evaluation of relational networks effects in different scenarios. And finally, since *Switch* was used in the original VDN paper, we can provide a clear comparison between our method and VDN.

B. Models and Hyperparameters

In our setup for the *Switch* environment, we configured our algorithm as follows: We used a Multi-Layer Perceptron (MLP) with two hidden layers, each consisting of 128 neurons and utilizing the ReLU activation function. For training each agent's prediction model, we performed $m = 10$ iterations per episode using batches of size $b = 32$, randomly sampled from a replay memory with a capacity of 50k time-steps. We employed the *Adam* optimizer with

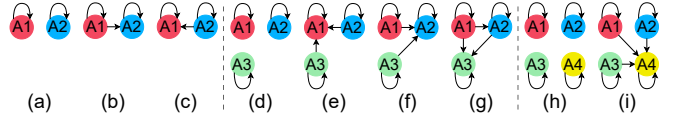


Fig. 2: Relational networks used in G-VDN. The edges represent whether an agent puts importance on another agent, including themselves. The relational networks in (b-c) were used in *two-agent* experiments in section IV-C.1; (e-g) were used in *three-agent* experiments in section IV-C.2; (i) was used in the *four-agent* experiment in section IV-C.3.

a learning rate of 0.001 and utilized the squared TD-error as the loss function. The weights of the target network were updated every 200 episodes with the weights of the prediction network. To encourage exploration, we employed the ϵ -greedy method, where ϵ linearly decreased over time. Finally, we set the discount factor (γ) to 0.99.

C. Results

As the environment can accommodate up to four agents, we have designed three different setups for our experiments, comprising of two, three, or four agents. This has resulted in six independent experiments. In all experiments, we compared our algorithm against VDN over the individual and collective rewards collected by the team.

In all the figures that depict the results of the experiments (i.e., Figs. 3, 4, and 5), the shaded regions display the average training reward over 10 runs, while the solid lines represent the average test rewards of the agents. The test rewards are determined by interrupting the training every 50 episodes and evaluating the individual rewards of the agents using the greedy strategy. In Figs. 3, 4, and 5, collective reward depicts the sum of individual rewards that is achieved by the agents in the environment. Note that collective reward is simply the addition of the individual agents' rewards (i.e., the team reward maximized in VDN), which is different from the team reward which is used for training G-VDN.

1) *Two-Agent Scenario:* This scenario consists of two agents represented by the red and blue circles and their respective goal locations marked with the same colors as shown in Fig. 3(a). To successfully reach their respective goal locations, one agent must learn to wait until the other agent passes the bridge. Therefore, there are two optimal policies for the agents to follow: either the red agent passes the bridge first or the blue agent does. In each run, VDN converges to one of the two optimal policies that the algorithm has experienced first based on the randomness in the exploration phase. As seen in Fig. 3(b) and Table I, the individual rewards for both agents are nearly identical as a result of averaging their rewards over 10 runs. In other words, in some runs, the red agent learns to cross the bridge first, leading to the maximization of its individual reward, while in others, the blue agent crosses the bridge and maximizes its own reward. Next, we designed two relational networks for G-VDN to better understand the influence of team interactions on optimal solutions. Specifically, in one relational network, the red agent assigns importance to the blue agent

¹Accompanying video: <https://youtu.be/8m-rSsdloXU>

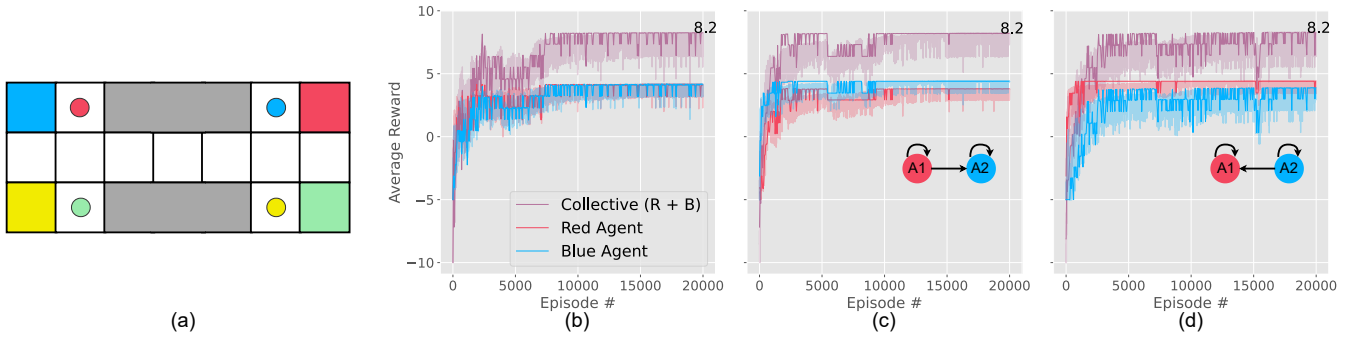


Fig. 3: (a) *Switch* environment with four agents. In the three-agent setup, the yellow agent and its goal location were removed. For the two-agent setup, the green agent and its goal location were in addition removed. Two-agent experiment results: (b) VDN, (c-d) G-VDN with relational networks in Fig. 2(b) and Fig. 2(c), respectively.

(Fig. 2(b)), while in the other relational network, the blue agent places importance on the red agent (Fig. 2(c)). The directed edges indicate whether an agent places importance on another agent, which dictate whether an agent yields to the other agent and lets them pass the bridge first for this environment. The results in Fig. 3(c) and Table I show that when training the network of agents with G-VDN and using the relationships shown in Fig. 2(b), the blue agent accumulates higher individual rewards (4.40 ± 0.00) as it is the first to use the bridge to reach its goal than the red agent (3.80 ± 0.00). Conversely, when Fig. 2(c) was applied, the situation is reversed with the red agent having a higher total reward (4.39 ± 0.02), as depicted in Fig. 3(d). Comparing the collective rewards from these three experiments reveals that they are all the same as the graph steers which *optimal* solution the algorithm converges to.

2) Three-Agent Scenario: In this scenario, the *Switch* environment used with three agents: the red and green agents are on the left and the blue agent is on the right side of the bridge as shown in Fig. 3(a). The optimal policy, minimizing waiting time, is for the pair of agents situated on the left to cross the bridge initially, as they are heading in the same direction and can proceed consecutively. Yet, because either agent can go first (red or green), this results in two optimal policies for crossing the bridge, either red-green-blue or green-red-blue. Fig. 4(a) and Table I depict the results obtained with VDN, where such preferences of order cannot be specified. The results show that the average individual rewards of the red (4.38 ± 0.04) and green (4.35 ± 0.02) agents are similar to each other, as they randomly alternate taking the lead in crossing the bridge for the 10 runs, and reach higher scores than those of the blue agent (3.66 ± 0.04). However, G-VDN can influence the order using the relational network depicted in Fig. 2(e), where the red agent has the highest importance. As demonstrated in Fig. 4(b) and Table I, the order of bridge usage becomes consistent as red-green-blue with the individual rewards of (4.40 ± 0.00), (4.30 ± 0.00), and (3.70 ± 0.00), respectively.

We have demonstrated that the type of *optimal* solution the algorithm converges to can be guided via relational networks. G-VDN, in addition, enables the agents to discover new be-

haviors to converge to based on the inter-agent relationship. For instance, the relationships in Fig. 2(f) show that the red and green agents prioritize the blue agent, granting it the first turn. The second turn is then randomly assigned between the red and green, as their importance weights are equal. The results of training the agents with G-VDN using this relational graph are shown in Fig. 4(c), which indicates that the blue agent's (4.40 ± 0.00) individual reward is higher than the others' and the red (3.75 ± 0.03) and green (3.76 ± 0.07) agents' rewards are almost the same due to alternating turns. To examine the extent to which the agents follow the relational network and prioritize defined cooperation strategy instead of solely maximizing the collective and individual rewards, we define a new relational graph shown in Fig. 2(g). The graph elicits the following sequence of

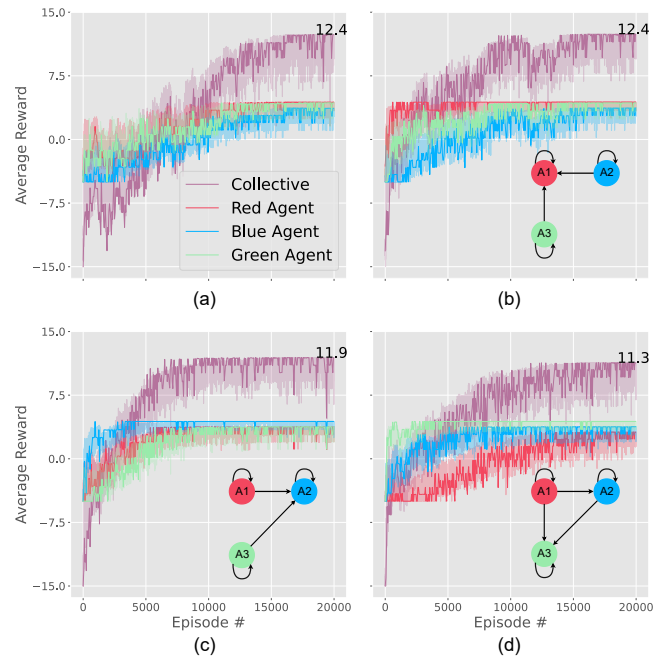


Fig. 4: Three-agent experiment results. (a) VDN, (b-d) G-VDN with relational networks in Fig. 2(e), Fig. 2(f), and Fig. 2(g), respectively.

TABLE I: Average reward with 95% confidence intervals for ten runs after training completed.

	Two-Agent Scenario			Three-Agent Scenario				Four-Agent Scenario	
	VDN	G-VDN	G-VDN	VDN	G-VDN	G-VDN	G-VDN	VDN	G-VDN
Relational Network	N/A	Fig. 2(b)	Fig. 2(c)	N/A	Fig. 2(e)	Fig. 2(f)	Fig. 2(g)	N/A	Fig. 2(i)
Red Agent	4.10±0.19	3.80±0.00	4.39±0.02	4.38±0.04	4.40±0.00	3.75±0.03	3.13±0.03	3.70±0.34	3.43±0.29
Blue Agent	4.15±0.15	4.40±0.00	3.86±0.06	3.66±0.04	3.70±0.00	4.40±0.00	3.78±0.02	3.78±0.25	3.61±0.17
Green Agent	—	—	—	4.35±0.02	4.30±0.00	3.76±0.07	4.40±0.00	3.77±0.36	3.54±0.23
Yellow Agent	—	—	—	—	—	—	—	3.80±0.30	4.34±0.09

crossing the bridge: green(left)–blue(right)–red(left), causing the red agent to wait the alternating turn instead of crossing the bridge a step after the green agent. Figure 4(d) and Table I show that the agents followed the defined relationships, demonstrating their ability to prioritize relationships over collective reward, even at a cost. It should be noted that, among the experiments in this section, the lowest collective reward is obtained with this relational graph as it requires the red agent to sacrifice more of its reward by taking the final turn.

3) *Four-Agent Scenario*: The final experiment involves four agents in the *Switch* environment as shown in Fig. 3(a) with equal numbers of agents each side. Applying VDN does not guarantee a deterministic solution for which side crosses the bridge first or which agent on the side takes the first turn. Instead, it leads to a stochastic selection process for these decisions, resulting in the learned team behavior varying from one run to the next. Our results depicted in Fig. 5(a) confirm that according to VDN, the individual rewards are very close (similar) due to the randomness caused by having multiple solutions that equally place the importance on all the agents without specifying any relationships (See Table I).

Yet, when using G-VDN, the relational network shown in Fig. 2(i) influences the order in which the bridge should be navigated. Since all agents place the importance to the yellow agent, the right side learns to take the first turn and the yellow agent is the first to use the bridge. This behavior is illustrated in Fig. 5(b). Similar to the three-agent scenario, we expect that the blue agent crosses the bridge after the yellow agent to compensate for the negative reward of remaining in the environment. However, due to the increased complexity of the team network and small amount of negative reward, this behavior has been observed in 80% of the runs (8/10 runs).

In the other 20%, the bridge was crossed by either of the three agents randomly. As shown in Fig. 5(b) and Table I, this behavior causes the individual rewards for the blue, red, and green agent to be very similar, while the yellow agent to attain the highest reward (4.34 ± 0.09). Between the green and the red agent, the order of using the bridge was determined randomly (since they have the same weight), resulting in similar individual rewards.

D. Real-World Experiments

To validate the effectiveness of our proposed algorithm in a real-world environment, we replicated the *Switch* grid-world in the lab as illustrated in Fig. 1. We employed four Turtlebot4 mobile robots marked by different colors, identical to our simulation experiments. Robots communicate through the Robotic Operating System (ROS) and a separate identification parameters has been assigned to each robot. We used a graph-based SLAM algorithm [28] to map the environment and localize the robots. Based on this map, we determined the coordinate of the cells in the *Switch* environment. In each time-step, the robots receive commands generated by (the trained) G-VDN through ROS to either stay in their current cell or move to a neighboring cell. The accompanying video shows the execution of the tasks.

V. CONCLUSION AND FUTURE WORK

We have proposed a novel algorithm that utilizes a relational network to influence agents' behavior within a team, leading to new cooperative behaviors. Given a team structure, our algorithm discovers and guides cooperation strategies that are driven by agents' relationships. As our results demonstrate, our algorithm is effective in influencing team behavior according to the specified relational network even in real-world multi-robot applications that require team cooperation.

Our proposed approach also has a few inherent limitations that open up novel research challenges for future work. Although in this paper, we consider the team structure to be given, methods for discovering the team structure could be investigated. For instance a high-level cost (e.g., observed behavior and capability differences of agents) can be used to find an optimal structure. Another challenging aspect of the proposed approach is that the relational network becomes denser as the number of agents increases. Thus, it becomes essential to develop methods for dealing with the tedious and eventually impractical task of adjusting relational weights.

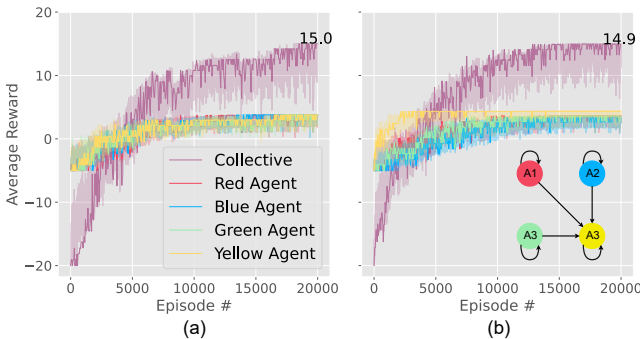


Fig. 5: Four-agent experiment results. (a) VDN, (b) G-VDN with relational network in Fig. 2(i).

ACKNOWLEDGEMENT

This work is supported in part by NSF (IIS-2112633) and the Army Research Lab (W911NF20-2-0089).

REFERENCES

- [1] L. Busoniu, R. Babuska, and B. De Schutter, "A comprehensive survey of multiagent reinforcement learning," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 38, no. 2, pp. 156–172, 2008.
- [2] G. Bresson, Z. Alsayed, L. Yu, and S. Glaser, "Simultaneous localization and mapping: A survey of current trends in autonomous driving," *IEEE Transactions on Intelligent Vehicles*, vol. 2, no. 3, pp. 194–220, 2017.
- [3] S. Grigorescu, B. Trasnea, T. Cocias, and G. Macesanu, "A survey of deep learning techniques for autonomous driving," *Journal of Field Robotics*, vol. 37, no. 3, pp. 362–386, 2020.
- [4] D. Calisi, A. Farinelli, L. Iocchi, and D. Nardi, "Autonomous navigation and exploration in a rescue environment," in *IEEE International Safety, Security and Rescue Robotics, Workshop, 2005*. IEEE, 2005, pp. 54–59.
- [5] A. Kleiner, J. Prediger, and B. Nebel, "Rfid technology-based exploration and slam for search and rescue," in *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2006, pp. 4054–4059.
- [6] V. G. Díaz, J. C.-W. Lin, and J. A. M. Molinera, "Editorial on "recent advances in logistics transportation with autonomous systems"," *Soft Computing*, vol. 25, no. 18, pp. 11 897–11 898, 2021.
- [7] M. Hüttenrauch, S. Adrian, G. Neumann, *et al.*, "Deep reinforcement learning for swarm systems," *Journal of Machine Learning Research*, vol. 20, no. 54, pp. 1–31, 2019.
- [8] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334–1373, 2016.
- [9] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, *et al.*, "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [10] F. A. Oliehoek, M. T. Spaan, and N. Vlassis, "Optimal and approximate q-value functions for decentralized pomdps," *Journal of Artificial Intelligence Research*, vol. 32, pp. 289–353, 2008.
- [11] 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," *arXiv preprint arXiv:1706.05296*, 2017.
- [12] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, "Monotonic value function factorisation for deep multi-agent reinforcement learning," *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 7234–7284, 2020.
- [13] 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 *International conference on machine learning*. PMLR, 2019, pp. 5887–5896.
- [14] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Advances in neural information processing systems*, vol. 30, 2017.
- [15] C. Yu, A. Velu, E. Vinitzky, Y. Wang, A. Bayen, and Y. Wu, "The surprising effectiveness of ppo in cooperative, multi-agent games," *arXiv preprint arXiv:2103.01955*, 2021.
- [16] C. Guestrin, D. Koller, and R. Parr, "Multiagent planning with factored mdps," *Advances in neural information processing systems*, vol. 14, 2001.
- [17] W. Böhmer, V. Kurin, and S. Whiteson, "Deep coordination graphs," in *International Conference on Machine Learning*. PMLR, 2020, pp. 980–991.
- [18] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *Proceedings of the tenth international conference on machine learning*, 1993, pp. 330–337.
- [19] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, pp. 279–292, 1992.
- [20] A. Tampuu, T. Matiisen, D. Kodelja, I. Kuzovkin, K. Korjus, J. Aru, J. Aru, and R. Vicente, "Multiagent cooperation and competition with deep reinforcement learning," *PloS one*, vol. 12, no. 4, p. e0172395, 2017.
- [21] P. Hernandez-Leal, M. Kaisers, T. Baarslag, and E. M. De Cote, "A survey of learning in multiagent environments: Dealing with non-stationarity," *arXiv preprint arXiv:1707.09183*, 2017.
- [22] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Basar, "Fully decentralized multi-agent reinforcement learning with networked agents," in *International Conference on Machine Learning*. PMLR, 2018, pp. 5872–5881.
- [23] C. Claus and C. Boutilier, "The dynamics of reinforcement learning in cooperative multiagent systems," *AAAI/IAAI*, vol. 1998, no. 746-752, p. 2, 1998.
- [24] Y. Findik, P. Robinette, K. Jerath, and S. R. Ahmadzadeh, "Impact of relational networks in multi-agent learning: A value-based factorization view," in *Proc. 62nd IEEE Conference on Decision and Control (CDC)*, 2023, p. 8.
- [25] H. Haeri, R. Ahmadzadeh, and K. Jerath, "Reward-sharing relational networks in multi-agent reinforcement learning as a framework for emergent behavior," *arXiv preprint arXiv:2207.05886*, 2022.
- [26] R. S. Sutton, "Learning to predict by the methods of temporal differences," *Machine learning*, vol. 3, pp. 9–44, 1988.
- [27] A. Koul, "ma-gym: Collection of multi-agent environments based on openai gym." <https://github.com/koulanurag/ma-gym>, 2019.
- [28] S. Macenski and I. Jambrecic, "Slam toolbox: Slam for the dynamic world," *Journal of Open Source Software*, vol. 6, no. 61, p. 2783, 2021.