

Cooperative Multi-Agent Reinforcement Learning with Hypergraph Convolution

Yunpeng Bai^{*,1,4}, Chen Gong^{*,2,3,4}, Bin Zhang^{*,1,4}, Guoliang Fan^{1,4}, , Xinwen Hou^{2,4}, Yu Liu^{2,4}

¹ Fusion Innovation Center, Institute of Automation, Chinese Academy of Sciences

² Comprehensive information system research Center, Institute of Automation, Chinese Academy of Sciences

³ School of Computing and Information Systems, Singapore Management University

⁴ School of Artificial Intelligence, University of Chinese Academy of Sciences

{baiyunpeng2020, gongchen2020, zhangbin2020, guoliang.fan, xinwen.hou, yu.liu}@ia.ac.cn

Abstract—Recent years have witnessed the great success of multi-agent systems (MAS). Value decomposition, which decomposes joint action values into individual action values, has been an important work in MAS. However, many value decomposition methods ignore the coordination among different agents, leading to the notorious “lazy agents” problem. To enhance the coordination in MAS, this paper proposes *HyperGraph CoNvolution MIX* (HGCN-MIX), a method that incorporates hypergraph convolution with value decomposition. HGCN-MIX models agents as well as their relationships as a hypergraph, where agents are nodes and hyperedges among nodes indicate that the corresponding agents can coordinate to achieve larger rewards. Then, it trains a hypergraph that can capture the collaborative relationships among agents. Leveraging the learned hypergraph to consider how other agents’ observations and actions affect their decisions, the agents in a MAS can better coordinate. We evaluate HGCN-MIX in the StarCraft II multi-agent challenge benchmark. The experimental results demonstrate that HGCN-MIX can train joint policies that outperform or achieve a similar level of performance as the current state-of-the-art techniques. We also observe that HGCN-MIX has an even more significant improvement of performance in the scenarios with a large amount of agents. Besides, we conduct additional analysis to emphasize that when the hypergraph learns more relationships, HGCN-MIX can train stronger joint policies.

Index Terms—Hypergraph Convolution, Multi-Agent Reinforcement Learning, Coordination and Control

I. INTRODUCTION

Deep Reinforcement Learning (DRL), which uses deep neural networks (DNN) to learn an optimal policy, has increasingly gained attention from researchers and achieved successful applications in various fields, including video games [1]–[4], robotics [5], etc.

In reality, many tasks require multiple agents, e.g. autonomous vehicles [6], resource allocation problems [7], etc. The single-agent RL algorithms ignore the coordination between agents, such as the independent Q learning (IQL) [8], which demonstrates poor performance in multi-agent systems (MAS). Therefore, enhancing coordination in MAS becomes an essential direction to improve the performance of multi-agent reinforcement learning (MARL). MARL algorithms

* THESE AUTHORS CONTRIBUTED EQUALLY.

THIS PROJECT WAS SUPPORTED BY NATIONAL DEFENSE FOUNDATION REINFORCEMENT FUND, NATIONAL DEFENSE BASIC SCIENTIFIC RESEARCH PROGRAM (JCKY2019203C029, JCKY2019-207B022) AND NATIONAL FOUNDATION FUND (Y2009-1B-02-353).

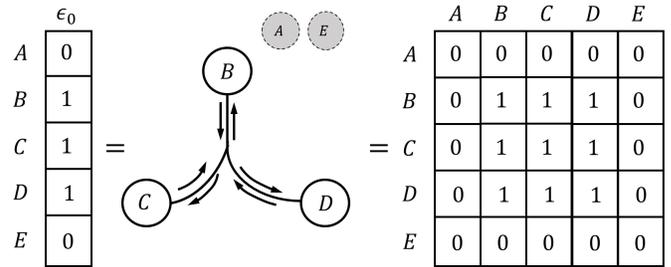


Fig. 1. As we shown in the mid figure, the hyperedge ϵ_0 in the left part connects nodes B, C, D and disconnects A and E . The mid figure also presents that how neighborhood information on the hyperedge ϵ_0 flows. And the right part shows a normal graph that equals to hyperedge ϵ_0 .

train a joint policy [9] to maximize the cumulative rewards in a Markov Decision Process (MDP). Due to the non-stationarity problem [10], it is hard to explain how the environments feedback to the agents a reward based on the agents’ actions. Value decomposition [11] is proposed based on the centralized training with decentralized execution (CTDE) paradigm [12], to study that how individual action values (introduced in Section III-A) affect joint action values. In CTDE, during the training phase, each agent acquires global information, but is not dependent on the global information during the execution phase [12]. It is noticed that the global information includes individual observation, actions, from other agents. Value-Decomposition Networks (VDN) [11] and QMIX [13], which are proposed on top of value decomposition (explained in Section III-B), do not consider cooperation between multi-agents. Insufficient cooperation causes the notorious issue, i.e., “lazy agents” [13] who take non-optimal actions due to the wrong credit assignment. For example, 6 of 10 agents have learned optimal or sub-optimal policies, but the rest four agents learn nothing.

To alleviate this “lazy agents” problem, several studies (CommNet [14], DGN [15]) utilize the observations or embedding of recurrent neural network (RNN) transmitted by neighbor agents to enhance communication. Through better communication, agents obtain more information about other agents, especially other agents’ observations and actions; thus, DGN [15] and CommNet [14] can significantly improve

prior state-of-the-art works. For each agent, collecting the information from other agents, i.e., observations, actions, is an important part to enhance the coordination of MAS. Graph Convolutional Network (GCN) achieves the integration by aggregating and propagating neighborhood information for those irregular or non-euclidean nature of graph-structured data [16]. Through GCN, we learn the embedding of each node or the whole graph based on the geometric structure. As shown in Figure 3(a), in most value decomposition methods, such as QMIX, QPLEX [17], individual observations and actions are transmitted to individual action values through RNN. We take individual action values as the embedding of each agents' observations and actions. Treating each agent as each node in a graph and performing GCN, the embedding of agents' observations and actions is attached with neighborhood information. However, it is difficult to represent the complex connections between a large of agents with a simple graph for these two reasons: (1) there is no prior knowledge for connections between agents to build the graph, and (2) connections between agents have a dynamic change due to the uncertainty of environments, observations, and global states. Hypergraphs [18], [19] can overcome these shortcomings. In hypergraphs, the value of hyperedge denotes the degree of connection between nodes. Specifically, the large weight represents a strong connection and vice versa. As presented in Figure 1, we list a simple example that equals a simple hyperedge to a regular graph. As for the uncertainty of environments, we use neural networks that take dynamic individual observations as input to learn the connections in our hypergraphs. Thus, hypergraphs are automatically adjusted by different observations.

As the analysis mentioned above, our contributions are summarized as follows, we proposed an end-to-end method termed HGCN-MIX, which combines value decomposition and hypergraph convolution (HGCN), aiming at enhancing the coordination between agents via HGCN. In HGCN-MIX, we build the hypergraphs generated via agents' individual observations without any prior knowledge. Besides, compared with the traditional fixed and binary hypergraph, in HGCN-MIX the hypergraphs are non-binary and dynamic adjust with the change of agent's observations. We perform our algorithms on 12 different scenarios in SMAC [20]. The experimental results show that our method outperforms or achieves a similar level of performance as the state-of-art MARL approaches in most scenarios. Especially in *MMM2*, and *27m_vs_30m*, HGCN-MIX obtains a higher 20% winning rate than QMIX. Besides, we conduct an ablation study to research the influence of using hypergraphs with different number of hyperedges. Ablation emphasizes that connections, represented as hyperedges, play an important role in training a strong joint policy.

II. RELATED WORKS

This section discusses several works in MARL and value decomposition methods briefly, and then introduces the developments of GCN and HGCN.

A. Multi-Agent Reinforcement Learning (MARL)

Various researches contribute to overcoming the multi-agent cooperation challenge. Independent Q learning (IQL) [8] extends the architecture of Deep Q-Learning Network (DQN) [1] to MAS. In mean field MARL, considering that each agent's effect on the overall population can become infinitesimal [21]. Although mean field MARL has demonstrated excellent performance in mean field games (MFGs), MFGs require homogeneous agents [22]. Recently, based on the centralized training with decentralized execution (CTDE) paradigm, some remarkable approaches such as value function decomposition [11], and multi-agent policy gradient methods [23] have been proposed. Compared with value-based methods, policy gradient methods can be easily applied to the environment with continuous actions space. In Multi-agent deep deterministic policy gradient (MADDPG) [12], the critic network takes overall individual observations and actions as inputs. As mentioned in [24], although policy gradient (PG) methods significantly outperform the QMIX in some continuous environments (multi-agent particle world environment (MPE) [12]), value-based methods achieve more sample-efficient than PG approaches in some discrete environments (SMAC).

B. Value Decomposition Methods for MARL

Value decomposition methods decompose the joint state-action value function into a single global reward. On the basis of CTDE, plenty of MARL algorithms are proposed, e.g., Value-Decomposition Networks (VDN) [11], QMIX [13], QPLEX [17]. VDN denotes the joint action value as the sum of individual action values. QMIX approximates the joint state-action value function monotonously. MMD-MIX [25] applies MMD-DQN [26] to collect randomness information in situations. Qtran [27] relaxes the constraints of joint action-value functions, but has a poor performance in some complex tasks. The most related to our work is GraphMIX [28], which combines graph neural networks (GNNs) [29] and value decomposition. However, GraphMIX only succeeds in several particular scenarios (i.e., *6h_vs_8z*).

C. Graph Convolutional Network (GCN) and Hypergraph Convolutional Network (HGCN)

GCN is proposed in [16] for classification tasks. GCN aims to capture neighborhood information for the irregular or non-euclidean nature of graph-structured data. GraphSAGE leverages the node feature information [30]. Only sampled neighbor nodes features are used for leveraging through aggregating functions. ChebNET [31] performs Fourier transform and converts spatial graph data into Fourier domain. Although GCN has been widely used in many areas, such as event detection [32] and knowledge graph [33], we select HGCN to aggregate instead of GCN. Replacing graphs with hypergraphs, Bai et al. deviate the two forms of convolution in HGCN [19]. [34] performs k-means clustering method to construct hypergraphs. The hypergraph attention mechanism has also been proposed in [19]. Hypergraphs in the hypergraph attention mechanism are also dynamic and non-binary. Reasons about our selection are elaborated in Section IV-A.

III. BACKGROUND

In this section, we briefly reviews about the basis of Dec-POMDP, value decomposition, and the definitions of hypergraph and hypergraph convolution.

A. Decentralized partially observable Markov decision process (Dec-POMDP)

Dec-POMDP [35] describes the fully cooperative multi-agent tasks. A Dec-POMDP can be defined by a tuple $G = \langle \mathcal{S}, \mathcal{U}, \mathcal{P}, \mathcal{Z}, r, \mathcal{O}, n, \gamma \rangle$. $s \in \mathcal{S}$ denotes the global state of the environment. Each agent $a \in \mathcal{A} := \{1, \dots, n\}$ takes an individual action $u_a \in \mathcal{U}$ at each timestep. \mathcal{U} is the action set. All the individual actions in the same timestep form a joint action $\mathbf{u} \in \mathcal{U} \equiv \mathcal{U}^n$. $\mathcal{P}(s'|s, \mathbf{u}) : \mathcal{S} \times \mathcal{U} \rightarrow \Delta(\mathcal{S})$ denotes the transition probability for any joint actions of both agents, where $\Delta(\cdot)$ is the simplex. All the agents share the same reward function $r(s, \mathbf{u}) : \mathcal{S} \times \mathcal{U} \rightarrow \mathbb{R}$.

Each agent gets its own local individual partial observation $z \in \mathcal{Z}$ according to the observation function $\mathcal{O}(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{Z}$. The previous of actions and observations for a agent makes up its own action-observation history $\tau_a \in T \equiv (\mathcal{Z} \times \mathcal{U})$. Policy of agent a is represented as: $\pi_a(u_a|\tau_a) : T \times \mathcal{U} \rightarrow [0, 1]$. The cumulative discounted reward in timestep t is defined as $R^t = \sum_{l=t}^{\infty} \gamma^{l-t} r_l$ [36], and we aim to maximize the expected reward $\mathbb{E}_{\tau \sim \pi} [R^0 | s_0]$, where τ is a trajectory and $\gamma \in [0, 1)$ is the discount factor.

B. Value decomposition

In Dec-PODMP, agents share the same joint reward, leading to the design of a unique reward function to train independent agents difficultly. As we have discussed in Section I, directly applying DQN for each agent without any coordination may lead to “lazy agents” [13].

Based on the CTDE paradigm, value decomposition methods learn the decomposition of the joint value function. The core of CTDE is that a centralized critic (named “mixing network” in Figure 2) guides the agents during training. This critic only exists in the training phase and has the ability to access global states, individual observations, individual action values, and actions. Besides, during the training phase, each agent can acquire neighborhood information, including individual observation, actions from other agents. During the execution phase, the centralized critic is deleted from the structure, and agents are not allowed to access neighborhood information. There are several meaningful works proposed based on value decomposition, such as VDN, and QMIX. One of the most important equation in value decomposition is Individual-Global-Max (IGM) equation [27]:

$$\arg \max_{\mathbf{u}} Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u}) = \begin{pmatrix} \arg \max_{u_1} Q_1(\tau_1, u_1) \\ \vdots \\ \arg \max_{u_n} Q_n(\tau_n, u_n) \end{pmatrix}. \quad (1)$$

IGM assumes that the optimality of each agent is consistent with the global optimality. Almost all algorithms based on value decomposition satisfy the IGM condition and finally

achieve convergence. QMIX approximates the joint state-action value function monotonously:

$$\frac{\partial Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u})}{\partial Q_a(\tau_a, u_a)} \geq 0, \forall a \in 1, \dots, n, \quad (2)$$

where τ_a denotes action-observation history of agent a , and u_a denotes individual action of agent a . The joint value functions in Qtran is formulated as:

$$Q_{\text{tot}}(s, \mathbf{u}) = \phi(s, Q_0, (\tau_0, u_0), \dots, Q_n, (\tau_n, u_n)), \quad (3)$$

where ϕ varies in each method. In some complex tasks, Qtran has a poor performance. These methods (e.g., VDN, QMIX, etc.) pay more attention on decomposing the joint value function, but ignore the coordination between agents.

C. Definition of Hypergraph

A hypergraph [37] \mathcal{G} is consisted of a vertex set $\mathbf{V} = \{v_1, v_2, \dots, v_N\}$ (N represents the number of vertices) and the edge set \mathbf{E} with M hyperedges. The adjacency matrix of the hypergraph $H \subseteq \mathbb{R}^{N \times M}$ indicates the connection between a vertex and hyperedge. For a vertex $v_i \in \mathbf{V}$ and a hyperedge $\epsilon \in \mathbf{E}$, when v_i is connected by ϵ , $H_{i\epsilon} = 1$ exists, otherwise we have $H_{i\epsilon} = 0$. Besides, each hyperedge ϵ has a non-negative weight $w_{\epsilon\epsilon}$. All the $w_{\epsilon\epsilon}$ is form a diagonal matrix $\mathbf{W} \in \mathbb{R}^{M \times M}$. In a hypergraph, the vertex degree $\mathbf{D} \in \mathbb{R}^{N \times N}$ and the hyperedge degree $\mathbf{B} \in \mathbb{R}^{M \times M}$ are defined as $D_{ii} = \sum_{\epsilon=1}^M w_{\epsilon\epsilon} H_{i\epsilon}$, $B_{\epsilon\epsilon} = \sum_{i=1}^N H_{i\epsilon}$. It is noticed that \mathbf{D} and \mathbf{B} are diagonal matrices.

In an indirect graph, the degree of an edge is 2, presenting that each edge only connects 2 vertices. However, in a hypergraph, an edge may connect more than 2 vertices. Thus, the property of the hypergraph allows more information about the nodes included in the hyperedge. Furthermore, it is easy to incorporate prior knowledge of connections between agents by selecting various hyperedges.

D. Spectral convolution on hypergraph

Given a hypergraph $\mathcal{G} = (\mathbf{V}, \mathbf{E}, \boldsymbol{\Delta})$ with N vertices and M hyperedges, $\boldsymbol{\Delta} \in \mathbb{R}^{N \times N}$ is the hypergraph Laplacian positive semi-definite matrix [18]. Then, $\boldsymbol{\Delta}$ is factorized as $\boldsymbol{\Delta} = \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^T$, where $\boldsymbol{\Phi} = \text{diag}(\phi_1, \dots, \phi_n)$ is an orthonormal eigen vectors and $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is a diagonal matrix whose diagonal elements are the corresponding eigenvalues. The Fourier transform [16], [38] for a signal $\mathbf{x} = (x_1, \dots, x_n)$ is defined as $\hat{\mathbf{x}} = \boldsymbol{\Phi}^T \mathbf{x}$. Then, the spectral convolution of signal \mathbf{x} and filter g is defined as:

$$g * \mathbf{x} = \boldsymbol{\Phi} g(\boldsymbol{\Lambda}) \boldsymbol{\Phi}^T \mathbf{x}, \quad (4)$$

where $g(\boldsymbol{\Lambda}) = \text{diag}(g(\lambda_1), \dots, g(\lambda_n))$ is a function of the Fourier coefficients. However, the computing of our method incurs $\mathcal{O}(n^2)$ computational complexity. To solve this problem, Feng et al. [16], [18], [31] derive spectral convolution as

$$g * \mathbf{x} \approx \theta_0 \mathbf{x} - \theta_1 \mathbf{D}^{-\frac{1}{2}} \mathbf{H} \mathbf{W} \mathbf{B}^{-1} \mathbf{H}^T \mathbf{D}^{-\frac{1}{2}} \mathbf{x},$$

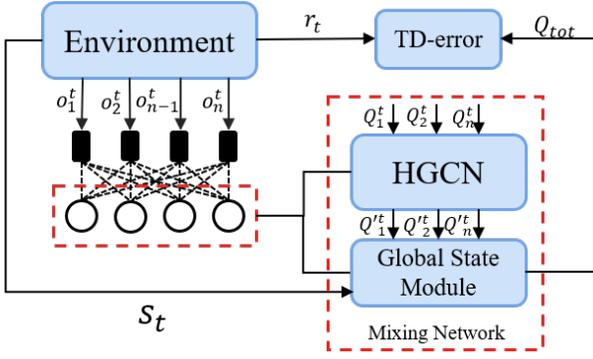


Fig. 2. The architecture of HGCN-MIX. At timestep t , agents first get their individual observations \mathbf{o}^t and output the individual action values \mathbf{Q} . Following Eq. (6), the hypergraph are generated with \mathbf{o}^t for HGNC. Then, HGNC takes individual action values \mathbf{Q} as the input and output the transformed action values \mathbf{Q}' . The “global state module” takes the global state s_t and \mathbf{Q}' , and outputs the Q_{tot} . The Q_{tot} and r_t are used to calculate the “TD-error”.

where θ_0 and θ_1 are the parameters of filtering over all nodes. Further, with the utilization of a single parameter θ is defined as:

$$\begin{cases} \theta_1 = -\frac{1}{2}\theta \\ \theta_0 = \frac{1}{2}\theta \mathbf{D}^{-\frac{1}{2}} \mathbf{H} \mathbf{B}^{-1} \mathbf{H}^T \mathbf{D}^{-\frac{1}{2}} \end{cases}.$$

Then, the hypergraph convolution (HGNC) is derived as:

$$\mathbf{x}^{(l+1)} = \mathbf{D}^{-\frac{1}{2}} \mathbf{H} \mathbf{W} \mathbf{B}^{-1} \mathbf{H}^T \mathbf{D}^{-\frac{1}{2}} \mathbf{x}^{(l)} \mathbf{P}, \quad (5)$$

where \mathbf{P} denotes the weight matrix between the (l) -th and $(l+1)$ -th layer [19]. For each node in HGNC, more information can be accessed between those nodes connected by a common hyperedge. The greater the weights, the more information about observations, histories and actions is provided.

IV. METHOD

In this section, we first introduce the reasons that we select HGNC instead of GCN or hypergraph attention. Secondly, we demonstrate the difficulties of traditional hypergraphs and our method to construct hypergraphs. Then, we introduce our model, HGCN-MIX, an end-to-end method base on CTDE paradigm. Finally, we have a brief discussion about the loss function that follows the TD-error [36].

A. Why hypergraph convolution

1) *Why hypergraph*: As discussed in Section I, GCN has the ability to aggregate neighbourhood information, but compared with the graph, there are two main reasons to select the hypergraph: (1) there are no prior knowledge about the connectivity between nodes (which denote agents) to build a regular graph. (2) relationships between agents have a dynamic change due to the dynamic environments, observations, and global states. Therefore, we need multiple graphs to learn the different relationships between agents. Due to a regular graph is equal to a hyperedge, as we shown in Figure 1, it is easily for us to build the multiple hyperedges representing multiple regular graphs.

2) *Why not attention*: It is difficult to learn the connections between agents for value decomposition methods in hypergraph attention mechanism because of two disadvantages: (1) hypergraph attention mechanism takes the pairwise similarity as the weights of connections in a hypergraph. (2) hypergraph attention takes the encoding of nodes as input. Nodes with different pairwise similarity get low weights even if the connections between these nodes are strong. For example, the two cikissis in *2c_vs_64zg* may take opposite actions to coordinate. The encoding of their observations and actions may be huge different but connections between them are quite strong. We aim at finding more different connections between agents, so pairwise similarity is not suitable here. In value-decomposition methods, individual action values are essential factors for selecting actions. Besides, the learning of embedding takes additional computational resource.

B. Difficulties of traditional hypergraphs

In Section III-C, we elaborate the definition of hypergraphs. Traditional hypergraphs usually have three drawbacks: (1) prior knowledge is required to build hypergraph, but the prior knowledge is usually hard to construct; and (2) the hypergraph is fixed when it is established, which leads to fixed connections between agents; and (3) the hypergraph only include two value (i.e., 1 or 0), indicating whether the node is connected by the hyperedge. Based on the analysis mentioned above, we rethink that agents need different levels of collaboration in different situations, e.g., in *MMM2* agents need a high-level collaboration, but only low-level collaboration is needed for agent in *2s3z*. Surely, the traditional hypergraph can not solve these problems we mentioned.

Instead of building a fixed binary hypergraph beyond prior knowledge, we generate a dynamic hypergraph via neural networks which adjusts automatically according to the agents’ individual observations. We use real-value hypergraph to describe the connections between agents, where different values of hyperedges represent different levels of collaboration.

C. Building hypergraph

To describe the degrees of collaboration appropriately, we replace the original binary hypergraph with a real-value hypergraph. As mentioned in [19], non-binary hypergraphs also follow Eq. (5).

In HGCN-MIX, a linear function is utilized, which takes individual observation $z_i \in \mathcal{Z}, i \in \{1, \dots, n\}$ as input, and generates a real values vector $\mathbf{A}_i \in \mathbb{R}^m$, where m denotes the number of hyperedges. Then, under the ReLU operation, the negative weights are modified to 0. On the one hand, negative weights break the IGM equation (Eq. (1)). On the other hand, negative weights make the hypergraph irreversible. Each value in \mathbf{A}_i , denoted as $\omega_{i,j}, j \in \{1, \dots, m\}$, presents the weights for hyperedge ϵ_j . The entirety of $\omega_{i,j}$ consists of the first part H_1 of our hypergraph H :

$$H_1 = \text{ReLU}(\text{Linear}(\mathcal{Z})). \quad (6)$$

The above equation lead agents to collaborate with others in different observations.

Inspired by GCN, we add an n -dimensional identity matrix, which is denoted as $H_2 \subseteq \mathbb{R}^{n \times n}$, and n is the number of agents. H_2 makes agents concentrate on themselves and maintain track of their own information. It is noticed that compared with H_1 , a fixed value is impossibly fitness for all the environments, e.g. in $3m$, 1 is too large but too small for $MMM2$. Therefore, we multiply H_2 by the average of H_1 to make the H_2 having the same orders as the H_1 's. Thus, we formulate H_2 as:

$$H_2 = \mathbb{I}_n \times \text{avg}(H_1).$$

In some environments with weak cooperation or few account of agents, H_2 matrix paves the path for keeping their own information.

We denote the average of H_1 as μ . Then, joining the two H_1 and H_2 matrices together, we reformulate the final hypergraph $H \subseteq \mathbb{R}^{n \times (m+n)}$ as:

$$H = \begin{bmatrix} \omega_{1,1} & \cdots & \omega_{1,m} \\ \omega_{2,1} & \cdots & \omega_{2,m} \\ \vdots & \cdots & \vdots \\ \omega_{n,1} & \cdots & \omega_{n,m} \end{bmatrix} \begin{bmatrix} \mu & \cdots & 0 \\ 0 & \cdots & 0 \\ \vdots & \cdots & \vdots \\ 0 & \cdots & \mu \end{bmatrix}, \quad (7)$$

$\underbrace{\hspace{10em}}_{H_1: \text{Learning}} \quad \underbrace{\hspace{10em}}_{H_2: \text{One-hot}}$

where $\omega_{i,k}, i \in [1, n]$, and $k \in [1, m]$ is the weight of learned connection between agents, and n denotes the number of agents and m denotes the number of hyperedges, μ is the average value of H_1 .

D. Mixing network

The structure of HGCN-MIX is presented in Figure 2. As shown in Figure 3(a), each agent applies a DRQN [39] to learn individual action value $Q_a \in \mathcal{Q}$, where $a \in \{1, \dots, n\}$. Then, the hypergraph convolution performs on \mathcal{Q} . In convolution, we set \mathbf{P} to \mathbb{I}_n , and \mathbf{W} to the learned parameters. To keep Eq. (1), we use the absolute value of \mathbf{W} for HCGN,

$$\mathbf{Q}' = \text{HGCN}(\text{HGCN}(\mathbf{Q}, \text{abs}(\mathbf{W}_1)), \text{abs}(\mathbf{W}_2)),$$

where $\mathbf{W}_i, i \in \mathbb{R}$ represents the weight matrix in the i -th layer and $\mathbf{Q}' \in \mathbb{R}^n$ denotes the modified action values. Each action value Q_a appends the information contained in all the Q_a by our method, so that agents can strengthen their collaboration. Then, a QMIX module attempts to add global state information to \mathbf{Q}' with \mathbf{Q}' as an input. As shown in Figure 3(b), Q_{tot} represents the joint action value, which can be formulated as:

$$\begin{aligned} \mathbf{Q}'' &= \text{elu}(\text{abs}(\text{MLP}(s_t)) \odot \mathbf{Q}' + \text{MLP}(s_t)) \\ Q_{tot} &= \text{sum}(\text{abs}(\text{MLP}(s_t)) \odot \mathbf{Q}'') + \text{MLP}(s_t) \end{aligned} \quad (8)$$

E. Loss function

Our loss function follows the TD-error:

$$\mathcal{L}(\theta) = \frac{1}{2} (y_{tot} - Q_{tot}(\boldsymbol{\tau}, \mathbf{u}|\theta))^2, \quad (9)$$

where the joint action-value function is parameterized by θ , and $y_{tot} = r + \gamma \max_{\mathbf{u}'} Q_{tot}(\boldsymbol{\tau}', \mathbf{u}'|\theta^-)$. The target network is parameterized by θ .

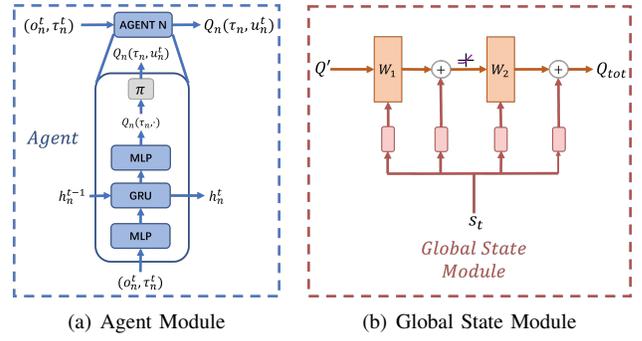


Fig. 3. Figure 3(a) shows the structure of the agent module. At timestep t , agent n take individual observation o_n^t and history τ_n^t as input, and output $\mathbf{Q}(\tau_n, \cdot)$ the values of legal actions allowed in the environment. Then, individual Q value $Q_n(\tau_n, u_n^t)$ is generated through the ϵ -greedy. Figure 3(b) shows the process to add the global state information. At timestep t , by applying four different MLPs and activation functions, global state s_t is transformed into four parts. Following Eq. (8), two of four parts have a vector inner product with transformed action values \mathbf{Q}' . The rest two parts are added to \mathbf{Q}' . Then, the joint action value Q_{tot} is generated.

V. EXPERIMENTS

In this section, firstly, we introduce settings in our experiments. Secondly, we conduct experiments on SMAC, which is commonly adopted for MARL algorithms. To demonstrate the validity of HGCN-MIX, we compare it with several well-known baseline algorithms: QMIX, QTRAN, VDN, RODE, and ROMA. Thirdly, we perform HGCN-MIX with a different number of hyperedges to explore the influence of hyperedges. Lastly, we conduct an ablation to explore the effectiveness of the learning part.

A. Settings

We conduct the experiments on SMAC [20] platform, which are commonly adopted to evaluate MARL algorithms. Our implementation is based on the framework for deep MARL algorithms, i.e., Pymarl [20]. Pymarl uses SMAC as the game environment and includes the implementations of QMIX, QTRAN, and VDN. The version of the Starcraft II is 4.6.2 (B69232) in our experiments¹. We select several well-known baseline algorithms include RODE [40], ROMA [41], QMIX, QTRAN [27], and VDN, as our baseline methods. Our experiments are carried out on Nvidia GeForce RTX 3090 graphics cards and Intel (R) Xeon (R) Platinum 8280 CPU. The policies are trained for 2 million timesteps. In HGCN-MIX, the number of learning hyperedge is set to 32. Other hyperparameters and the code of HGCN-MIX follow this link². We record the maximum median winning rates during the 2 million timesteps, as we shown in Table I.

B. Validation

In Figure 4, the solid line is the median winning rates of the 5 random seeds, and the shadow denotes the 25 – 75% percentiles of the winning rates. The experimental results

¹Performance is **NOT** usually comparable between versions [20].

²<https://github.com/cugbbaiyun/HGCN-MIX>

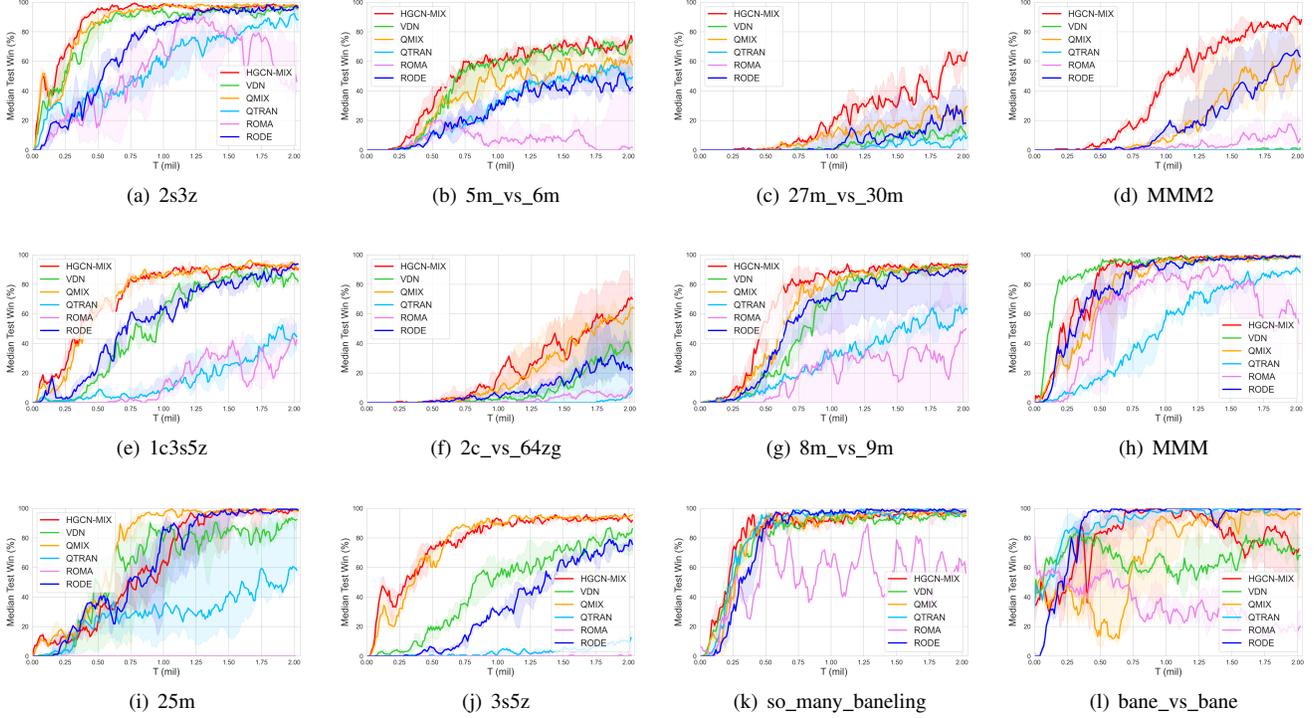


Fig. 4. Overall results in different scenarios. Exepect for *MMM2*, *MMM*, *so_many_baneling*, and *bane_vs_bane*, each digit represents the number of units and each letter corresponds to a type of unit (i.e., *27m* means 27 Marines) in scenario name. There are two super hard scenarios: *MMM2*, *27m_vs_30m* and three hard scenarios: *2c_vs_64zg*, *bane_vs_bane*, and *5m_vs_6m*. The rest scenarios are all easy scenarios. Notice that in those scenarios with special names, *MMM2*, *MMM*, *so_many_baneling*, and *bane_vs_bane* are all heterogeneous. *MMM2*, *so_many_baneling* are asymmetric.

illustrate that HGCN-MIX outperforms the QMIX approach significantly. In *27m_vs_30m*, and *MMM2*, HGCN-MIX achieves the 20% higher winning rate than baseline (i.e., QMIX). In *2s3z*, *5m_vs_6m*, and *8m_vs_9m*, the winning rates trend of HGCN-MIX has plateaued faster than baseline, however, in *25m*, *3s5z*, and *MMM*, we observe the opposite results. The learning of hypergraphs in HGCN-MIX leads to extra computational cost. In the rest 3 out of 7 easy scenarios, e.g., *1c3s5z*, *3s5z*, etc., HGCN-MIX reaches a similar performance as baseline methods.

It is noticed that in the hard scenario *2c_vs_64zg*, the curve of winning rate in HGCN-MIX convergent unstably. There are two reasons for this phenomenon: 1) The number of ally units is 2, due to the mechanism of the game, and the two cikissis usually take the contrast actions in most scenarios. For example, a cikissi goes down and the other one always goes up. It is difficult for HGCN to learn the coordination between these two cikissis. 2) The count of enemy units is 64. A number of enemy units lead to a large action space, so that we can hardly use the 1-dimensional Q value to represent individual observations and actions.

It can not be ignored that matrix multiplication and matrix inversion in HGCN-MIX bring the extra computational cost compared with QMIX. Considering the significant improvement in winning rates, we believe that this amount of extra cost is acceptable.

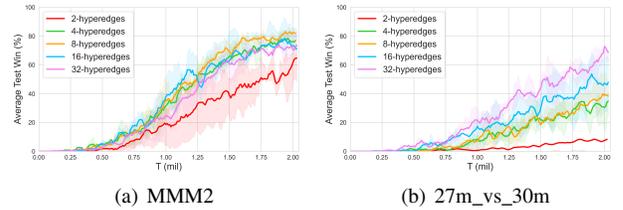


Fig. 5. Average winning rates of agent in two different scenarios. We select the two super hard scenarios: *MMM2* and *27m_vs_30m*. Both *27m_vs_30m* and *MMM2* are asymmetric and have large counts of agents (10 units in *MMM2* and 27 units in *27m_vs_30m*). Besides, *MMM2* is heterogeneous but *27m_vs_30m* is homogeneous.

C. Additional experiments

We perform this additional study to examine how the number of hyperedges influencing the performance of joint policies. We set the counts of hyperedges to 2, 4, 8, 16, and 32, respectively. All the hypergraphs learn without one-hot matrices. We test them on the two representative super hard scenarios: *27m_vs_30m* and *MMM2*, and collect the average winning rates as we shown in Figure 5.

The result in *27m_vs_30m* shows that, with the counts of hyperedges increasing, HGCN-MIX can train a stronger policy. Rather than improving the winning rates, the few hyperedges limit the cooperation between different agents.

TABLE I
PART OF THE MAXIMUM MEDIAN WIN RATIOS DURING TRAINING.

Algorithms	Maximum Median Win Ratios(%)					
	2s3z	56m	27m	3M2	1c3s5z	2c64zg
HGCN-MIX	100	81	69	94	97	71
QMIX	100	69	31	78	100	69
VDN	100	81	22	3	94	44
QTRAN	94	59	13	0	56	9
RODE	100	56	38	72	97	38
ROMA	94	25	0	19	50	13
Algorithms	89m	3M	25m	3s5z	somany	bane
HGCN-MIX	97	100	100	100	100	100
QMIX	97	100	100	100	100	100
VDN	97	100	97	91	100	94
QTRAN	69	91	62	16	100	100
RODE	94	100	100	84	100	100
ROMA	53	97	0	0	94	81

The bold underlined winning rate in each column denotes the best performance over 6 methods in each scenario. Scenario names are aggregated to save space. For example, 89m refers to *8m_vs_9m* and 3M2 refers to *MMM2*, bane refers to *bane_vs_bane*, somany refers to *so_many_baneling*.

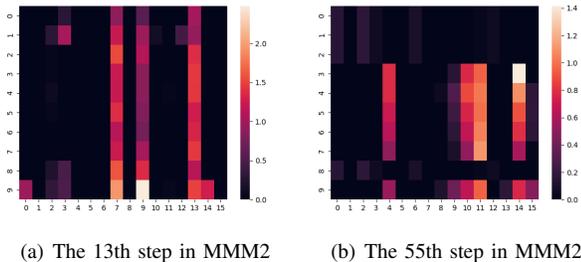


Fig. 6. Hypergraphs during training. The two figure show hypergraphs in different timesteps in the same episode. The scenario of the two figures is *MMM2*. Hypergraphs are generated as we have talked about in Section IV-C. In Figure 6(a) and 6(b), line represents agents connected with different hyperedges. Rows show weights in hyperedges, such as row 7 in Figure 6(a) represent the hyperedge ϵ_7 .

However, in *MMM2*, HGCN-MIX with few hyperedges also works well. With greater than 4 hyperedges, HGCN-MIX has reached the similar performance as baseline methods. To explore the reason that why HGCN-MIX works well in *MMM2* with few hyperedges, we record two hypergraphs in different timesteps and the same episode in *MMM2*. We visualize the hypergraphs and present in Figure 6.

How to read Figure 6. The same lines, such as 0, 1, 2, and 8 in Figure 6(b), represent these agents that have been slain, where the health of the agent is set to 0. Individual

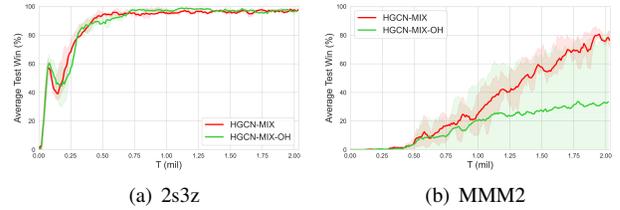


Fig. 7. Average test winning rates of HGCN-MIX and HGCN-MIX-OH in *2s3z* and *MMM2*. We select one easy scenario *2s3z* and one super hard scenario *MMM2* to explore the influence of the global state module and learning part in the hypergraph.

observations in dead agents are set to the same fixed mask values (i.e., -1). Therefore, the weights generated by neural networks have the same weights. The rows include many zeros introducing the hyperedges equipped with less cooperation. Besides, we can pay attention to the weights in the same hyperedges rather than different hyperedges. The diagonal matrix W (defined in Section III-C) balances the relative relationships between different hyperedges.

From Figure 6, we find that there are only 4 or 5 hyperedges learned different connections in *MMM2*. This result illustrates that only 4 or 5 essential connections exist between agents in *MMM2*. When the number of hyperedges is greater than 4, HGCN-MIX achieves a similar level of performance with the different number of hyperedges.

D. Ablation

In ablation, we conduct two algorithms (i.e., HGCN-MIX and HGCN-MIX-OH) on *MMM2* and *2s3z* to emphasize the influence of the learning part in Eq. (7). Hypergraph in HGCN-MIX keeps both one-hot part and learning part, and HGCN-MIX-OH only keeps the one-hot part in Eq. (7). In HGCN-MIX-OH, weights in one-hot part are set to 1. Following Eq. (5), we can derivate that when hypergraph is set as an identity matrix, no matter what W is, the result of HGCN satisfies: $x^{(l+1)} = x^l$. In this way, HGCN has no effectiveness. In other words, only the global state module works in the mixing network. We conduct HGCN-MIX and HGCN-MIX-OH on the super hard scenario *MMM2* and easy scenario *2s3z*, and record the average winning rate of the policy. It is noticed that the hyperparameters of training are the same as that of in HGCN-MIX approach.

As shown in Figure 7, in *2s3z*, both HGCN-MIX and HGCN-MIX-OH reach almost 100% winning rates, but in *MMM2* only HGCN-MIX achieves a high winning rate. HGCN-MIX-OH only works well in some easy scenarios and fails in super hard scenarios. This demonstrates that the global state module is not the most important point for the improvement of performance. The key to enhancing the coordination lies in the learning part.

VI. CONCLUSIONS AND FUTURE WORK

In this paper, we propose a novel MARL algorithm termed *HyperGraph CoNvolution MIX (HGCN-MIX)*, which combines hypergraph convolution with value decomposition to alleviate the notorious “lazy agent” issue. Besides, our approach

also overcomes the problem that decomposes joint state-action value functions without any prior knowledge about the connection between nodes (or agents). By using hypergraph convolution, for each agent, the aggregations of neighborhood information about observations and actions are easily achieved. Based on these information, agents may take better actions to enhance coordination. We perform HGCN-MIX on SMAC, a well-known benchmark for MARL. Experiments results show that HGCN-MIX outperforms in some scenarios with a large account of agents. We also conduct ablation experiments to illustrate the connections represented as hyperedges between different agents using neural networks. The learning part in HGCN-MIX plays an important role in HGCN-MIX. For the sake of HGCN-MIX only transforming the action values, it paves the path for the wide applications in many MARL algorithms based on value decomposition.

Since HGCN-MIX needs to learn the relationship of coordination between agents, the winning rates trend of HGCN-MIX has plateaued slower than baseline in some scenarios. The problem of speeding up the process of learning the coordination between agents will be studied in future work.

REFERENCES

- [1] V. Mnih, K. Kavukcuoglu, D. Silver *et al.*, “Human-level control through deep reinforcement learning,” *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [2] D. Silver, A. Huang, C. J. Maddison *et al.*, “Mastering the game of go with deep neural networks and tree search,” *nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [3] C. Gong, Y. Bai, X. Hou, and X. Ji, “Stable training of bellman error in reinforcement learning,” in *International Conference on Neural Information Processing*. Springer, 2020, pp. 439–448.
- [4] C. Gong, Q. He, Y. Bai, X. Hou, G. Fan, and Y. Liu, “Wide-sense stationary policy optimization with bellman residual on video games,” in *2021 IEEE International Conference on Multimedia and Expo, ICME 2021, Shenzhen, China, July 5-9, 2021*. IEEE, 2021, pp. 1–6.
- [5] S. Levine, C. Finn, T. Darrell *et al.*, “End-to-end training of deep visuomotor policies,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334–1373, 2016.
- [6] M. Rahmati, M. Nadeem, V. Sadhu *et al.*, “Uw-marl: Multi-agent reinforcement learning for underwater adaptive sampling using autonomous vehicles,” in *Proceedings of the International Conference on Underwater Networks & Systems*, 2019, pp. 1–5.
- [7] J. Cui, Y. Liu, and A. Nallanathan, “Multi-agent reinforcement learning-based resource allocation for uav networks,” *IEEE Transactions on Wireless Communications*, vol. 19, no. 2, pp. 729–743, 2019.
- [8] A. Tampuu, T. Matiisen, D. Kodolja *et al.*, “Multiagent cooperation and competition with deep reinforcement learning,” *PloS one*, vol. 12, no. 4, p. e0172395, 2017.
- [9] T. Zhang, Y. Li, C. Wang *et al.*, “Fop: Factorizing optimal joint policy of maximum-entropy multi-agent reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 491–12 500.
- [10] G. Papoudakis, F. Christianos, A. Rahman *et al.*, “Dealing with non-stationarity in multi-agent deep reinforcement learning,” *arXiv preprint arXiv:1906.04737*, 2019.
- [11] P. Sunehag, G. Lever, A. Gruslys *et al.*, “Value-decomposition networks for cooperative multi-agent learning,” *ArXiv*, vol. abs/1706.05296, 2018.
- [12] R. Lowe, Y. Wu, A. Tamar *et al.*, “Multi-agent actor-critic for mixed cooperative-competitive environments,” *arXiv preprint arXiv:1706.02275*, 2017.
- [13] T. Rashid, M. Samvelyan, C. Schroeder *et al.*, “Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 4295–4304.
- [14] S. Sukhbaatar, R. Fergus *et al.*, “Learning multiagent communication with backpropagation,” *Advances in neural information processing systems*, vol. 29, pp. 2244–2252, 2016.
- [15] J. Jiang, C. Dun, T. Huang *et al.*, “Graph convolutional reinforcement learning,” *arXiv preprint arXiv:1810.09202*, 2018.
- [16] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” *arXiv preprint arXiv:1609.02907*, 2016.
- [17] J. Wang, Z. Ren, T. Liu *et al.*, “Qplex: Duplex dueling multi-agent q-learning,” *arXiv preprint arXiv:2008.01062*, 2020.
- [18] Y. Feng, H. You, Z. Zhang *et al.*, “Hypergraph neural networks,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 3558–3565.
- [19] S. Bai, F. Zhang, and P. H. Torr, “Hypergraph convolution and hypergraph attention,” *Pattern Recognition*, vol. 110, p. 107637, 2021.
- [20] M. Samvelyan, T. Rashid, C. S. De Witt *et al.*, “The starcraft multi-agent challenge,” *arXiv preprint arXiv:1902.04043*, 2019.
- [21] Y. Yang, R. Luo, M. Li *et al.*, “Mean field multi-agent reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 5571–5580.
- [22] J.-M. Lasry and P.-L. Lions, “Mean field games,” *Japanese journal of mathematics*, vol. 2, no. 1, pp. 229–260, 2007.
- [23] C. Yu, A. Velu, E. Vinitzky *et al.*, “The surprising effectiveness of mappo in cooperative, multi-agent games,” *arXiv preprint arXiv:2103.01955*, 2021.
- [24] G. Papoudakis, F. Christianos, L. Schäfer *et al.*, “Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks,” in *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS)*, 2021. [Online]. Available: <http://arxiv.org/abs/2006.07869>
- [25] Z. Xu, D. Li, Y. Bai *et al.*, “Mmd-mix: Value function factorisation with maximum mean discrepancy for cooperative multi-agent reinforcement learning,” *arXiv preprint arXiv:2106.11652*, 2021.
- [26] T. Tang Nguyen, S. Gupta, and S. Venkatesh, “Distributional reinforcement learning with maximum mean discrepancy,” *arXiv e-prints*, pp. arXiv–2007, 2020.
- [27] K. Son, D. Kim, W. J. Kang *et al.*, “Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 5887–5896.
- [28] N. Naderialzadeh, F. H. Hung, S. Soleyman *et al.*, “Graph convolutional value decomposition in multi-agent reinforcement learning,” *arXiv preprint arXiv:2010.04740*, 2020.
- [29] F. Scarselli, M. Gori, A. C. Tsoi *et al.*, “The graph neural network model,” *IEEE transactions on neural networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [30] W. L. Hamilton, R. Ying *et al.*, “Inductive representation learning on large graphs,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 1025–1035.
- [31] M. Defferrard, X. Bresson, and P. Vandergheynst, “Convolutional neural networks on graphs with fast localized spectral filtering,” *Advances in neural information processing systems*, vol. 29, pp. 3844–3852, 2016.
- [32] T. H. Nguyen and R. Grishman, “Graph convolutional networks with argument-aware pooling for event detection,” in *Thirty-second AAAI conference on artificial intelligence*, 2018.
- [33] Z. Wang, Q. Lv, X. Lan *et al.*, “Cross-lingual knowledge graph alignment via graph convolutional networks,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 349–357.
- [34] J. Jiang, Y. Wei, Y. Feng *et al.*, “Dynamic hypergraph neural networks,” in *IJCAI*, 2019, pp. 2635–2641.
- [35] F. A. Oliehoek and C. Amato, *A concise introduction to decentralized POMDPs*. Springer, 2016.
- [36] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [37] C. Berge, *Hypergraphs: combinatorics of finite sets*. Elsevier, 1984, vol. 45.
- [38] R. N. Bracewell and R. N. Bracewell, *The Fourier transform and its applications*. McGraw-Hill New York, 1986, vol. 31999.
- [39] M. Hausknecht and P. Stone, “Deep recurrent q-learning for partially observable mdps,” in *2015 aaii fall symposium series*, 2015.
- [40] T. Wang, T. Gupta, A. Mahajan *et al.*, “Rode: Learning roles to decompose multi-agent tasks,” *arXiv preprint arXiv:2010.01523*, 2020.
- [41] T. Wang, H. Dong, V. Lesser *et al.*, “Roma: Multi-agent reinforcement learning with emergent roles,” *arXiv preprint arXiv:2003.08039*, 2020.