
State-based Episodic Memory for Multi-Agent Reinforcement Learning

Xiao Ma

Department of Computer Science and Technology
Nanjing University
max@lamda.nju.edu.cn

Wu-Jun Li

National Key Laboratory for Novel Software Technology
Department of Computer Science and Technology
Nanjing University
liwujun@nju.edu.cn

Abstract

Multi-agent reinforcement learning (MARL) algorithms have made promising progress in recent years by leveraging the centralized training and decentralized execution (CTDE) paradigm. However, existing MARL algorithms still suffer from the sample inefficiency problem. In this paper, we propose a simple yet effective approach, called state-based episodic memory (SEM), to improve sample efficiency in MARL. SEM adopts episodic memory (EM) to supervise the centralized training procedure of CTDE in MARL. To the best of our knowledge, SEM is the first work to introduce EM into MARL. We can theoretically prove that, when using for MARL, SEM has lower space complexity and time complexity than state and action based EM (SAEM), which is originally proposed for single-agent reinforcement learning. Experimental results on StarCraft multi-agent challenge (SMAC) show that introducing episodic memory into MARL can improve sample efficiency and SEM can reduce storage cost and time cost compared with SAEM.

1 Introduction

Reinforcement learning (RL) has achieved promising success in a variety of challenging domains, including game playing [10] and robotics [7]. There are multiple agents to collaboratively make sequential decisions in many real-world applications, such as autonomous driving [17], intelligent robotic control [5], and game playing [1, 16]. Therefore, multi-agent reinforcement learning (MARL) has attracted much attention. In MARL, each agent has its partial observations and chooses its actions in a shared environment with other agents' interactions. This complex interaction model causes many challenges, such as instability, sample inefficiency, the moving target problem (non-stationarity), the exponential growth of the agents' action space.

Early MARL methods, such as independent Q-learning [20], adopt the decentralized policy for training. In these methods, each agent takes action independently and treats other agents as part of the environment. Since other agents' policies will change, leading to a non-stationary environment, these decentralized policy based methods suffer from sample inefficiency and instability problems. To solve these problems caused by decentralized policies, researchers recently propose to use the paradigm called centralized training and decentralized execution (CTDE) [12]. Benefiting from CTDE, value-decomposed MARL algorithms [9, 19, 15, 3, 14, 22] have been proposed in recent years. During *training* phase of these algorithms, a joint action-value function, including all agents'

individual action-value functions, is learned using additional information such as global states, actions, or rewards. Each agent can act independently based on its local observation and its individual action-value function without any additional information in the *execution* phase. Although value-decomposed MARL algorithms have achieved promising success, improving sample efficiency is still a critical problem for MARL. Episodic memory (EM), which can record the best experiences, has been applied in single-agent RL to improve sample efficiency [13, 2, 6, 8]. However, to our best knowledge, EM has not been introduced into MARL.

In this paper, we make the first attempt to introduce EM into MARL and propose a novel method, called state-based episodic memory (SEM), to improve sample efficiency for MARL. The contributions of this work are briefly outlined as follows:

- SEM is the first work to introduce EM into MARL to improve sample efficiency.
- SEM establishes only one lookup table to record global states and their corresponding highest discounted return among all executed joint actions. SEM replays the highest discounted return from the table as the EM target to supervise the centralized training in the paradigm of CTDE.
- When using for MARL, SEM can be theoretically proved to have lower space complexity and time complexity than state and action based EM (SAEM), which is originally proposed for single-agent RL [13, 2, 6, 8].
- Experimental results on StarCraft multi-agent challenge (SMAC) show that introducing episodic memory into MARL can improve sample efficiency and SEM can reduce storage cost and time cost compared with SAEM.

2 Notation

In this paper, we adopt similar notations as those in [15]. More specifically, we describe the fully cooperative multi-agent task as a decentralized partially observable Markov decision process (Dec-POMDP) [11]. A Dec-POMDP is defined by a tuple $G = \langle S, U, P, r, Z, O, n, \gamma \rangle$. Here, $S \in \mathbb{R}^F$ denotes the global state of the environment and U denotes the action space of each agent $a \in A \equiv \{1, \dots, n\}$. At each time step, each agent a chooses an action $u^a \in U$, forming a joint action $\mathbf{u} \in \mathbf{U} \equiv U^n$. It results in a transition on the environment according to state transition function $P(s'|s, \mathbf{u}) : S \times \mathbf{U} \times S \rightarrow [0, 1]$. All agents receive the same reward according to the same reward function $r(s, \mathbf{u}) : S \times \mathbf{U} \rightarrow \mathbb{R}$ and $\gamma \in [0, 1]$ is a discount factor. Dec-POMDPs consider partially observable scenarios in which each agent has individual observation $o \in O$ according to observation function $Z(s, a) : S \times A \rightarrow O$. Each agent has an observation-action history $\tau^a \in \mathcal{T} \equiv (O \times U)^*$ and chooses its action based on a stochastic policy $\pi^a(u^a | \tau^a) : \mathcal{T} \times U \rightarrow [0, 1]$. $\boldsymbol{\tau} \in \mathbf{T} \equiv \mathcal{T}^n$ is the joint action-observation history. The joint action-value function of the joint policy π is defined as: $Q_{tot}^\pi(\boldsymbol{\tau}_t, \mathbf{u}_t) = \mathbb{E}_{\boldsymbol{\tau}_{t+1:\infty}, \mathbf{u}_{t+1:\infty}} [R_t | \boldsymbol{\tau}_t, \mathbf{u}_t]$, where $R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$ is the discounted return.

3 Related Work

3.1 EM for Single-Agent RL

In single-agent RL, deep reinforcement learning algorithms need to take millions of interactions with the environments to attain human-level performance. However, humans can quickly exploit the high reward after the first discovery by using the hippocampus, which can record EM. Motivated by the hippocampus' ability, researchers proposed to use EM to achieve fast learning and improve sample efficiency for single-agent RL. Model-free episodic control (MFEC) [2] and neural episodic control (NEC) [13] try to use lookup tables to record the EM and retrieve useful values from lookup tables for action selection. MFEC and NEC can be seen as tabular RL methods, which lack good generalization compared with deep neural network-based RL methods. Episodic memory deep q-network (EMDQN) [8] combines EM with deep q-network (DQN) to achieve good generalization and improve sample efficiency by accelerating the training process of DQN. EM has been widely applied in single-agent RL, but it has not been introduced into MARL. Furthermore, all existing EM-based single-agent RL methods adopt state and action based EM (SAEM). SAEM will face many difficulties when SAEM is applied into MARL, which will be detailedly described in Section 4.

3.2 Value-based MARL Algorithms

The most common and straightforward value-based approach in the multi-agent setting is to break down a multi-agent learning problem into multiple independent single-agent learning problems, which is called decentralized value-based methods. One of the representative methods is independent Q-learning (IQL) [20]. Decentralized value-based methods benefit from scalability because each agent can make decisions based on a decentralized policy. However, each agent has to treat other agents as a part of the environment. The other agents' policies will change during the training procedure, making the environment not stationary. Therefore, these decentralized value-based methods suffer from sample inefficiency and instability problems.

The sample inefficiency and instability problems caused by decentralized policies can be alleviated by adding centralized training into MARL, resulting in a MARL paradigm called centralized training with decentralized execution (CTDE) [12]. Benefiting from the paradigm of CTDE, researchers have proposed some approaches, called value-decomposed MARL algorithms, that learn a centralized but decomposed Q value function to improve agent learning performance in recent years [19, 15, 18, 14, 22]. During execution, each agent still acts independently. In the value-decomposed MARL algorithms, individual-global-max (IGM) is a crucial principle, ensuring that the optimal joint action across agents in the joint action-value $Q_{tot}(\boldsymbol{\tau}, \mathbf{u}; \theta)$ and the collection of all optimal individual action in the $[Q_a(\boldsymbol{\tau}^a, \mathbf{u}^a; \theta)]_{a=1}^n$ are consistent. The parameters θ are learned by minimizing the following expected TD error:

$$L(\theta) = \sum_{b=1}^B \sum_{t=1}^T (Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b; \theta) - y_t^b)^2, \quad (1)$$

where $y_t^b = r_t^b + \gamma \max_{\mathbf{u}} Q_{tot}(\boldsymbol{\tau}_{t+1}^b, \mathbf{u}; \theta^-)$ and θ^- are the parameters of a target network that are periodically copied from θ . VDN [19] and QMIX [15] respectively propose two decomposition structures, additivity and monotonicity, which are sufficient conditions for the IGM. WQMIX [14] tries to use a weighted projection that places more importance on better joint actions to overcome the limitation of QMIX. QTRAN [18] tries to realize the entire IGM function class by using extra soft regularizations, which actually loses the IGM guarantee. QPLEX [22] uses a duplex dueling architecture and provides a guaranteed IGM consistency. Although these value-decomposed MARL algorithms can outperform decentralized value-based methods, they still suffer from sample inefficiency and instability problems.

4 State and Action based Episodic Memory for MARL

The EM for single-agent RL [13, 2, 6, 8] is defined on state and action. This is why we call it state and action based EM (SAEM) in this paper. To the best of our knowledge, EM, including SAEM, has not been introduced into MARL before. In this section, we will extend the application of SAEM from single-agent RL to MARL, by replacing the action in single-agent RL with the joint action in MARL.

More specifically, we establish a lookup table for each joint action and denote this lookup table as $Q^{\text{SA}}(s, \mathbf{u})$, where s is the global state and \mathbf{u} is the joint action. Each entry in the table $Q^{\text{SA}}(s, \mathbf{u})$ records the highest return ever obtained by taking joint action \mathbf{u} in state s . At the end of the episode, we store the episode $(\mathbf{o}_1, \mathbf{u}_1, s_1, r_1, \dots, \mathbf{o}_T, \mathbf{u}_T, s_T, r_T)$ into the replay buffer H and store $(s_t, \mathbf{u}_t, R(s_t))$ into a set M of size $|M|$, where $R(s_t) = \sum_{i=t}^T \gamma^{i-t} r_i$ is the discounted return received after taking joint action \mathbf{u}_t in state s_t . Q^{SA} is updated when the set M is filled:

$$Q^{\text{SA}}(s, \mathbf{u}) \leftarrow \begin{cases} R_{\text{Max}}(s, \mathbf{u}), & \text{if } (s, \mathbf{u}) \notin Q^{\text{SA}}; \\ \max \{Q^{\text{SA}}(s, \mathbf{u}), R_{\text{Max}}(s, \mathbf{u})\}, & \text{otherwise,} \end{cases} \quad (2)$$

where $R_{\text{Max}}(s, \mathbf{u}) = \max_k \{R_k(s, \mathbf{u})\}$ and $(s, R_k(s, \mathbf{u})) \in M, k \in \{1, 2, \dots, K\}$. For any state and joint action, we update their values with the highest discounted return. After updating $Q^{\text{SA}}(s, \mathbf{u})$, the set M is made empty. The value of each entry in the table Q^{SA} is updated increasingly and the number of entries also increases during training. We limit the maximum size of the table Q^{SA} and remove the least frequently assessed entry when Q^{SA} is filled.

In the centralized training phase, the loss function of SAEM is defined as:

$$L(\theta) = \sum_{b=1}^B \sum_{t=1}^T (1 - \lambda) (Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b; \theta) - y_t^b)^2 + \lambda \left(Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b; \theta) - E_{s_t}^{b, \mathbf{u}_t^b} \right)^2, \quad (3)$$

where $y_t^b = r_t^b + \gamma \max_{\mathbf{u}} Q_{tot}(\boldsymbol{\tau}_{t+1}^b, \mathbf{u}; \theta^-)$ is the vanilla target in the value-decomposed MARL algorithms and $E_{s_t}^{b, \mathbf{u}_t^b} = Q^{SA}(s_t^b, \mathbf{u}_t^b)$ is the episodic memory target recorded by Q^{SA} . $\lambda \in [0, 1]$ is the coefficient to balance the trade-off between the two targets.

We can see that SAEM suffers from two deficiencies when using for MARL: high space complexity and high time complexity. Because SAEM needs $|U|^n$ lookup tables, the space complexity is $\mathcal{O}(c|U|^n)$, where $c > 1$ is a constant. Hence, the space complexity is exponentially higher than that in single-agent RL. The time complexity is also high. In the worst case, there are $|M|$ different joint actions, and hence $|M|$ tables are needed to be updated when using the set M to update Q^{SA} . Hence, the time complexity of updating Q^{SA} is $\mathcal{O}(|M|)$ in the worse case. This motivates us to design new EM mechanisms like state-based EM in the next section.

5 State-based Episodic Memory for MARL

This section introduces our newly proposed EM, called state-based episodic memory (SEM), for MARL.

5.1 Lookup Table in SEM

In SEM, we establish only one lookup table $Q^S(s)$, which is indexed by global states. Given the global state s , each entry in $Q^S(s)$ records the highest return ever obtained. At the end of episode, we store the episode $(\mathbf{o}_1, \mathbf{u}_1, s_1, r_1, \dots, \mathbf{o}_T, \mathbf{u}_T, s_T, r_T)$ into the replay buffer H and store $(s_t, R(s_t))$ into a set M , where $R(s_t) = \sum_{i=t}^T \gamma^{i-t} r_i$ is the discounted return received after taking joint action \mathbf{u}_t in state s_t . Q^S is updated according to the set M :

$$Q^S(s_t) \leftarrow \begin{cases} R_{\text{Max}}(s_t), & \text{if } s_t \notin Q^S; \\ \max\{Q^S(s_t), R_{\text{Max}}(s_t)\}, & \text{otherwise,} \end{cases} \quad (4)$$

where $R_{\text{Max}}(s_t) = \max_k \{R_k(s_t)\}$ and $(s_t, R_k(s_t)) \in M$, $k \in \{1, 2, \dots, K\}$. For any state, we update their values with the highest return. After updating Q^S , the set M is made empty. The value of each entry in the table Q^S is updated increasingly and the number of entries also increases during training. We limit the maximum size of the table Q^S and remove the least frequently assessed entry when Q^S is filled.

5.2 Training Procedure

In the centralized training procedure, the loss function of SEM is defined as:

$$L(\theta) = \sum_{b=1}^B \sum_{t=1}^T (1 - \lambda)(Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b; \theta) - y_t^b)^2 + \lambda(Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b; \theta) - E_t^b)^2, \quad (5)$$

where $y_t^b = r_t^b + \gamma \max_{\mathbf{u}} Q_{tot}(\boldsymbol{\tau}_{t+1}^b, \mathbf{u}; \theta^-)$ and $E_t^b = r_t^b + \gamma Q^S(s_{t+1}^b)$. θ^- is the parameter of the target network which is copied from θ . y_t^b is a target inferred by the target network and E_t^b is an episodic memory target. $\lambda \in [0, 1]$ is the coefficient to balance the trade-off between two targets. When λ is set to 0, our method degenerates to the original value-decomposed MARL algorithm.

During training, the vanilla target y_t^b , approximated by using the target network, might be over-estimated, which could mislead the training. The episodic memory target E_t^b (or $E_{s_t^b, \mathbf{u}_t^b}^{b, \mathbf{u}}$ in SAEM) is more stable than the vanilla target, because the episodic memory target is retrieved from the lookup table, rather than approximated by the function. Furthermore, the maximum operator in updating $Q^S(s)$ not only records the highest return but also plays the role of the maximum operator $\max_{\mathbf{u}} Q_{tot}(\boldsymbol{\tau}_{t+1}^b, \mathbf{u}; \theta^-)$, because we ignore the action \mathbf{u} given the state s in the $Q^S(s)$. If we cannot find the $Q^S(s_{t+1}^b)$ in the lookup table, we use $\max_{\mathbf{u}} Q_{tot}(\boldsymbol{\tau}_{t+1}^b, \mathbf{u}; \theta^-)$ to replace $Q^S(s_{t+1}^b)$. Algorithm 1 briefly presents the training procedure of SEM. The architecture of SEM is shown in Figure 1. Please note that each agent takes actions based on Q_a and our method can be combined with all existing value-decomposed MARL algorithms.

Algorithm 1 State-based Episodic Memory (SEM) for MARL

Initialize a replay buffer H , an empty set M , an episodic memory table Q^S .
 Initialize network parameter θ and $\theta^- = \theta$.
for each episode **do**
 for $t = 1, 2, 3, \dots, T$ **do**
 Receive observation $[o_t^a]_{a=1}^n$, global state s_t .
 Select a random action u_t^a with probability ϵ , otherwise $u_t^a = \arg \max_{u^a} Q_a(\tau^a, u^a)$ for each agent a .
 Take action $[u_t^a]_{a=1}^n$.
 end for
 Store the episode $(\mathbf{o}_1, \mathbf{u}_1, s_1, r_1, \dots, \mathbf{o}_T, \mathbf{u}_T, s_T, r_T)$ in H .
for $t = T, T-1, \dots, 1$ **do**
 Let $s_t = \phi(s_t)$.
 Compute R_t and store (s_t, R_t) into M .
end for
 Sample B episodes from the replay buffer H .
 Compute y_t^b and E_t^b .
 Update θ by minimizing the loss in (5).
 Update target network parameter $\theta^- = \theta$ with period I_n .
 Update $Q^S(s_t)$ using M according to (4) when M is filled and then empty M .
end for

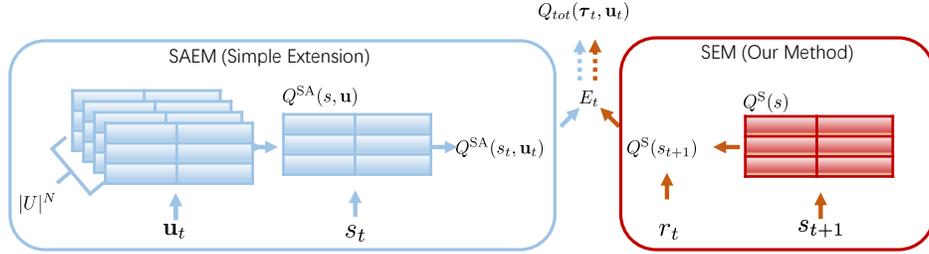


Figure 1: Comparison between SAEM architecture for MARL and SEM architecture. Best viewed in color.

5.3 Complexity

SEM only need one table to store and update, rather than $|U|^n$ tables in SAEM. The space complexity of SEM is $O(c_s)$, where c_s is a constant. When $Q^S(s)$ is updated by the set M , it needs to update the $Q^S(s)$ only once. Hence, the time complexity of SEM is $O(1)$. Hence, we can theoretically prove that SEM has lower space complexity and time complexity than SAEM, when using for MARL.

5.4 State Representation

Although SEM has lower space complexity than SAEM, the storage cost of all global states in the table $Q^S(s)$ might still be large if the dimension of state space is high. Similar to [2, 13, 8], we utilize random projection ϕ to project a state from the original global state space S with dimension F into a lower-dimensional space with dimension D . The random projection ϕ is denoted as $\phi(s) : s \rightarrow Vs$, where each entry in $V \in \mathbb{R}^{D \times F}$ is randomly drawn from a standard Gaussian. Based on Johnson-Lindenstrauss lemma [4], when the random matrix V is drawn from a standard Gaussian, this transformation approximately keeps relative distances in the original space. Projecting global states to lower-dimension vectors can accelerate the speed of table lookup. Please note that we still use $Q^S(s)$ to denote the table in SEM, omitting ϕ although random projection is adopted in this paper. This state representation can also be adopted in SAEM.

6 Experiment

6.1 StarCraft Multi-Agent Challenge

StarCraft multi-agent challenge (SMAC) ¹, based on the real-time strategy game StarCraft II and the SC2LE environment [21], is a popular benchmark for cooperative multi-agent RL. SMAC focuses on micromanagement challenges. In SMAC, it involves two armies, one controlled by the build-in AI and the other controlled by the user. Each unit can be controlled by an independent agent. At each time-step, agents receive their local observations, which depend on their sight range. The agents are allowed to take actions, including move[direction], attack[enemy_id], stop and no-op. Agents can only move in four directions: north, south, east, or west. If the enemy is within the agent’s shooting range, the agent can perform the action attack[enemy_id]. The maximum number of actions an agent can take ranges between 7 and 70, depending on the scenario. The goal of agents is to maximize the win rate for each battle scenario. The default setting of SMAC is to use the shaped reward. For evaluation, we run 32 evaluation episodes without any exploratory behaviors every 10000 time-steps. More details about the implementation are included in Appendix of the supplementary materials.

6.2 Effectiveness of Episodic Memory

We evaluate SAEM on 2c_vs_64zg, a hard map on SMAC, to verify the effectiveness of EM for MARL. We choose two classic value-decomposed MARL algorithms, VDN and QMIX, as the baselines. We combine SAEM with VDN and QMIX, which are denoted as SAEM-VDN and SAEM-QMIX, respectively. The coefficient λ is set to 0.1. The size of the lookup table is set to 1 million and the size of M is set to 5K. Other hyper-parameters are described in Appendix. The results are shown in Figure 2. We can find that the methods with SAEM have better performance than baselines without episodic memory, verifying that introducing episodic memory into the multi-agent setting is effective.

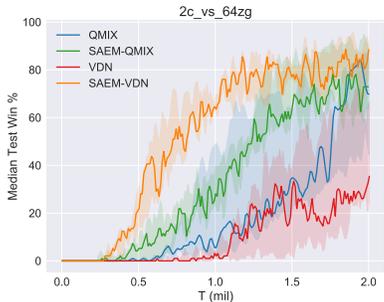


Figure 2: Results of SAEM-VDN, SAEM-QMIX and baselines (VDN and QMIX) on 2c_vs_64zg, including median test win rate as well as the 25-75% percentiles.

6.3 Results of SEM

We evaluate our method, SEM, on the eight maps of SMAC. These maps include 1c3s5z, 2s_vs_1sc, 2s3z, 3s5z, 27m_vs_30m, 2c_vs_64zg, MMM2 and bane_vs_bane. The snapshots and configurations of these maps are shown in the Appendix. We choose several value-decomposed MARL algorithms as baselines, including VDN [19], QMIX [15], QPLEX [22], WQMIX [14] ². The coefficient λ is set to 0.1. The size of the lookup table is set to 1 million and the size of set M is set to 5K, which means that the lookup table is updated every 5K time-steps. Other hyper-parameters are described in Appendix. We show our results, mean and median scores ³ over eight maps of SMAC, at 0.25M time-steps and 0.5M time-steps in Table 1. We can find that the mean scores and median scores of SEM combined with value-decomposed MARL algorithms all surpass those of the corresponding vanilla value-decomposed MARL algorithms. It verifies that SEM can improve sample efficiency compared with baselines that do not use episodic memory.

In Figure 3, we show the training curves of SEM-VDN, SEM-QMIX, VDN and QMIX on six maps ⁴. Our methods, SEM-QMIX and SEM-VDN, can converge faster than corresponding baselines and improve sample efficiency on all maps. On 2s3z, the performances of QMIX and VDN both drop sharply at about 0.2 million time-steps. Our method can alleviate this phenomenon obviously since episodic memory could help the agents to remember the best experience. On 27m_vs_30m, SEM-VDN can achieve 40% test battle won while VDN is bound to fail in all the battles. The

¹We use SC2.4.6.2.69232 (the same version as that in [16]), instead of the newer version SC2.4.10.

²The weighting function of WQMIX is the centrally-weighting function.

³The formula of mean scores and the formula of median scores are shown in the Appendix.

⁴Due to the limitation of pages, the training curves of SEM-VDN, SEM-QMIX, VDN, and QMIX on the other two maps are shown in the Appendix.

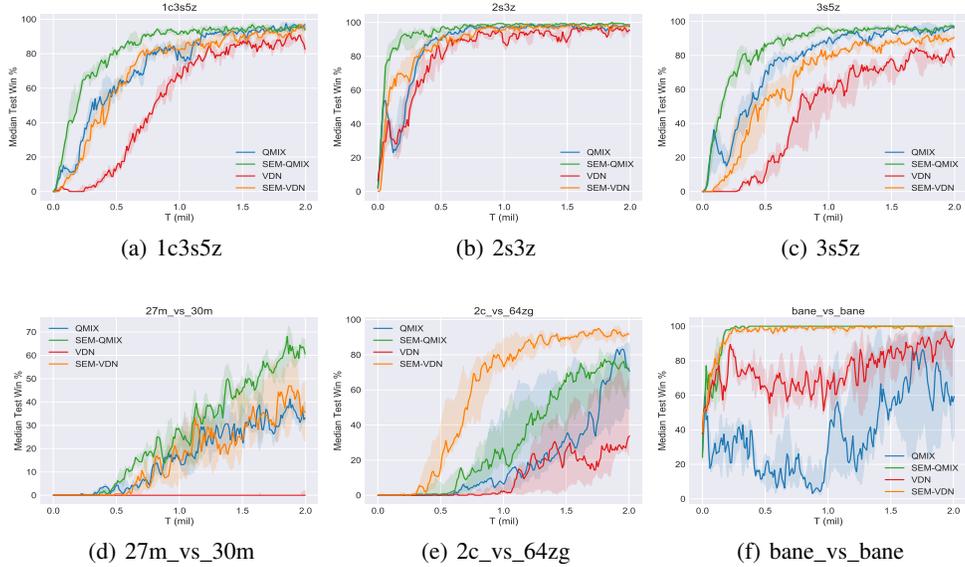


Figure 3: Results of our methods (SEM-VDN and SEM-QMIX) and baselines (VDN and QMIX), including the median performance as well as the 25-75% percentiles.

performance of SEM-QMIX has also increased by about 15% compared to that of QMIX. On 2c_vs_64zg, SEM-VDN can achieve 80% test battle won while VDN fails in all battles. SEM-QMIX can learn faster than QMIX. For bane_vs_bane, the results of QMIX and VDN exhibit a large variance. Our methods, SEM-QMIX and SEM-VDN, can both outperform the baselines (QMIX and VDN)

Table 1: Mean and median scores for eight maps on SMAC at 0.25M time steps and 0.5M time steps. "w/o" represents the vanilla baselines and "w/" represents the SEM combined with the baselines. Boldface numbers indicate best results. The results of each map is shown in the Appendix.

Baselines	0.25M				0.5M			
	Mean		Median		Mean		Median	
	w/o	w/	w/o	w/	w/o	w/	w/o	w/
VDN	26%	36%	1%	25%	33%	55%	14%	64%
QMIX	21%	44%	18%	47%	35%	55%	28%	76%
QPLEX	45%	50%	46%	69%	58%	63%	84%	90%
WQMIX	26%	42%	20%	40%	46%	52%	49%	54%

with a large margin and the median test win of our methods converges to 1 quickly.

6.4 Comparison among Different Targets

To gain the insight of introducing episodic memory into MARL, we try to conduct an in-depth analysis of the training process. We compare different targets in SAEM-VDN, SEM-VDN and VDN on 2c_vs_64zg. y_t^b is denoted as y , which is a target inferred by the target network. E_t^b is denoted as E_s , which is an episodic memory target replayed from Q^S . E_s^u denotes $E_{s_t}^{b,u_t^b}$, which is replayed from Q^{SA} . The results are shown in Figure 4. Both E_s and E_s^u are stable targets, because they are retrieved from the tables which record the highest return from historical episodes. In VDN, we can see that y is higher than the episodic memory targets $\{E_s, E_s^u\}$ with a large margin. The gap between y and $\{E_s, E_s^u\}$ in SEM-VDN and SAEM-VDN is smaller than that in VDN. It illustrates that the vanilla target y , approximated by using target network, is overestimated and the episodic memory target can provide centralized training with a stable target. Hence, our methods can improve the performance and sample efficiency. Furthermore, we can find that the curves of E_s^u and E_s are almost overlapped in Figure 4, which illustrates that using E_s to replace E_s^u is reasonable.

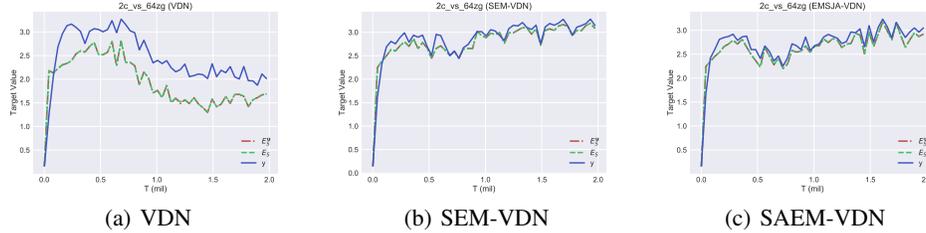


Figure 4: Comparison among different targets, E_s^u , E_s and y , in our methods and baselines. (a) VDN. (b) SEM-VDN. (c) SAEM-VDN. The results are summarized over 5 random runs. For clarity, we represent the median performance without the 25-75% percentiles.

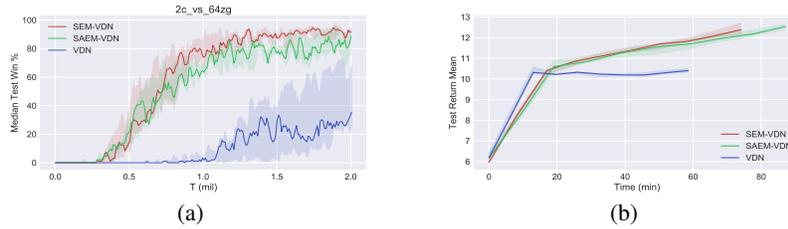


Figure 5: (a) Results of SAEM-VDN, SEM-VDN and VDN on 2c_vs_64zg.(b) Test return mean with respect to wall-clock time of SAEM-VDN, SEM-VDN and VDN during the first 100K time-steps.

6.5 Comparison between SEM and SAEM

Compared with SAEM, it has been theoretically proved in Section 5.3 that SEM has lower space complexity and time complexity. We choose 2c_vs_64zg and 27m_vs_30m as the test environment to compare the space complexity and time complexity of SEM and SAEM.

The space complexity of SEM and SAEM is affected by the number of lookup tables. Here, the size of all lookup tables is set to 1 million, and the dimension D is set to 4. For SEM, the storage cost is fixed because it has only one lookup table, which needs 0.029GB storage space to store. For SAEM, the number of lookup tables is equal to the number of joint actions. In 2c_vs_64zg, it contains two allied agents, and each agent has 70 available actions. SAEM needs 142GB storage space to store 70^2 lookup tables. In 27m_vs_30m, there are 27 agents, and each agent has 36 actions. SAEM needs $36^{27} \approx 10^{42}$ lookup tables, which needs 3×10^{40} GB storage space to store. We can see that the high storage cost makes the implementation of SAEM difficult and even impossible. But SEM has a much lower storage cost than SAEM.

For time cost, we choose 2c_vs_64zg to experimentally illustrate that SEM is more time-saving than SAEM. Here, the size of the set M is set to 5K. Other hyper-parameters are described in Appendix. In Figure 5(a), we compare the results of SAEM-VDN, SEM-VDN and VDN. In Figure 5(b), we compare the test return means with respect to the wall-clock time of SAEM-VDN, SEM-VDN and VDN during the first 100K time-steps using Nvidia Geforce GTX 2080 Ti graphics cards. We can see that SEM can use relatively less time than SAEM to achieve a similar return and SEM can save 14% of time cost compared with SAEM.

6.6 Learned Policies

In this section, we investigate the learned behaviors of the different policies in our methods to understand the differences between the strategies better. Here, we choose 27m_vs_30m, 2c_vs_64zg, MMM2 and bane_vs_bane to investigate. For other scenarios, both our methods and baselines can learn a good policy, although the baselines have worse sample efficiency than our methods.

In the 27m_vs_30m scenario, the allied army contains 27 Marines and the enemy army contains 30 Marines. For QMIX, the allied units can learn to stand in a line. For VDN, it bounds to fail the battle. In our methods, the allied army can stand in a line more evenly and more rapidly before the battle than baselines. Hence, these allied units can join the battle faster than those in baselines, as shown in Figure 6(a) and Figure 6(b).

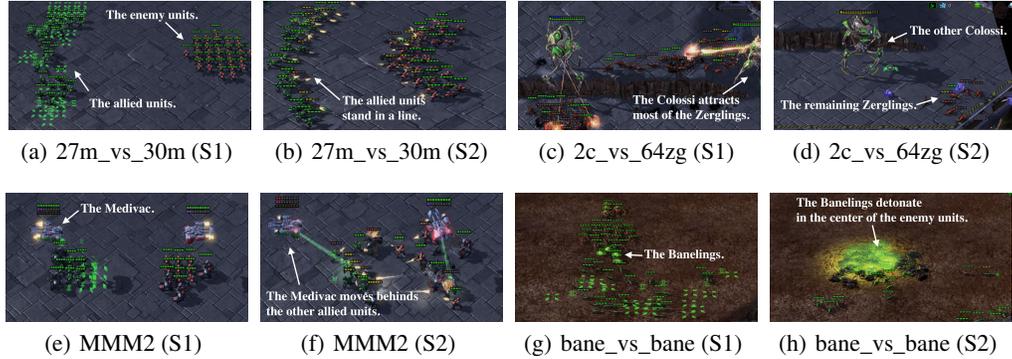


Figure 6: Illustration of the learned policy in several scenarios.

The `2c_vs_64zg` scenario contains two allied units (Colossi) but 64 enemy units (Zerglings), which leads to a much larger action space than the other scenarios. In this scenario, the enemy units are divided into two groups and the allied units are surrounded by enemy units from two opposite directions. For VDN, it fails most of the battles. For QMIX, although it can learn a good policy, it has worse sample efficiency than our methods. For our methods, a Colossi attracts most of the Zerglings to move far away from the other Colossi. The Colossi, with most of the Zerglings, kills some enemy units and then it is killed. The other Colossi kills the enemy units around it and then it searches for the remaining enemy units and kills them.

For MMM2, it contains 1 Medivac, 2 Marauders and 7 Marines. The Medivac can heal damage for Marauders and Marines. Therefore, the key to winning the battle is the Medivac. In our method, the Medivac can move behind the allied units and avoid sacrifice, which can heal other allied units continuously. In baselines, the Medivac can also move behind the allied units, but it is too late in the battle to sacrifice finally and cannot heal other agents continuously.

On the `bane_vs_bane` scenario, it contains a large number of allied and enemy units. The allied army and the enemy army contain 20 Zerglings and 4 Banelings, respectively. Both VDN and QMIX struggle and exhibit large variance, as shown in Figure 3(f). Our methods can learn faster and finally converge. The essentially learned policy of our method is that the four allied Banelings walk into the enemy army’s center (Figure 6(g)) and then detonate where it is standing, damaging almost all of the enemy units (Figure 6(h)). This learned policy is concise and practical, which alleviates the instability to a large extent.

6.7 Sensitivity to Hyper-Parameters

In order to better understand our method, we investigate the effect of the balance coefficient λ and the size of the lookup table $|Q^S|$ on `bane_vs_bane`, shown in Figure 7. The coefficient λ is chosen from $\{0, 0.01, 0.05, 0.1, 0.2, 0.5, 1.0\}$ and $|Q^S|$ is chosen from $\{10^4, 10^5, 10^6, 2 \times 10^6\}$. When λ is set to 0, our method degenerates to the baseline. In most cases, our method performs better than the baseline. When $\lambda = 1.0$, our method only uses episodic memory to supervise the training procedure. When $\lambda \in \{0.05, 0.1, 0.2, 0.5, 1\}$, our method performs better and learns a good policy faster than the baseline. For $|Q^S|$, we can see that when $|Q^S|$ is small, it would lose some information and then deteriorate the performance. When $|Q^S| \in \{10^5, 10^6, 2 \times 10^6\}$, our method performs well. We also investigate the sensitivity to the other hyper-parameters in the appendix, including the update frequency of Q^S and the dimension of random projection for state representation.

7 Conclusion

In this paper, we have proposed a novel and effective method, SEM, to improve sample efficiency in MARL. To our best knowledge, this is the first work that introduces episodic memory into the multi-agent setting. Compared with the existing EM mechanisms, SEM has lower space complexity and time complexity. Experimental results on SMAC have verified the effectiveness and efficiency of SEM.

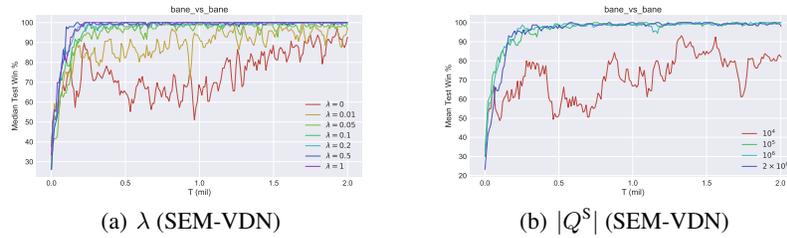


Figure 7: Median performance of SEM-VDN with different value of the coefficient λ and $|Q^S|$. The results are summarized over 5 random runs. For clarity, we represent the median performance without the 25-75% percentiles.

References

- [1] C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Dennison, D. Farhi, Q. Fischer, S. Hashme, C. Hesse, R. Józefowicz, S. Gray, C. Olsson, J. Pachocki, M. Petrov, H. P. de Oliveira Pinto, J. Raiman, T. Salimans, J. Schlatter, J. Schneider, S. Sidor, I. Sutskever, J. Tang, F. Wolski, and S. Zhang. Dota 2 with large scale deep reinforcement learning. *CoRR*, abs/1912.06680, 2019.
- [2] C. Blundell, B. Uria, A. Pritzel, Y. Li, A. Ruderman, J. Z. Leibo, J. W. Rae, D. Wierstra, and D. Hassabis. Model-free episodic control. *CoRR*, abs/1606.04460, 2016.
- [3] J. N. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson. Counterfactual multi-agent policy gradients. In *AAAI*, 2018.
- [4] W. B. Johnson and J. Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics*, 26(189-206):1, 1984.
- [5] J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- [6] M. Lengyel and P. Dayan. Hippocampal contributions to control: The third way. In *NeurIPS*, 2007.
- [7] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. In *ICLR*, 2016.
- [8] Z. Lin, T. Zhao, G. Yang, and L. Zhang. Episodic memory deep q-networks. In *IJCAI*, 2018.
- [9] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *NeurIPS*, 2017.
- [10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [11] F. A. Oliehoek and C. Amato. *A Concise Introduction to Decentralized POMDPs*. Springer Briefs in Intelligent Systems. Springer, 2016.
- [12] F. A. Oliehoek, M. T. J. Spaan, and N. A. Vlassis. Optimal and approximate q-value functions for decentralized pomdps. *Journal of Artificial Intelligence Research*, 32:289–353, 2008.
- [13] A. Pritzel, B. Uria, S. Srinivasan, A. P. Badia, O. Vinyals, D. Hassabis, D. Wierstra, and C. Blundell. Neural episodic control. In *ICML*, 2017.
- [14] T. Rashid, G. Farquhar, B. Peng, and S. Whiteson. Weighted QMIX: expanding monotonic value function factorisation for deep multi-agent reinforcement learning. In *NeurIPS*, 2020.
- [15] T. Rashid, M. Samvelyan, C. S. de Witt, G. Farquhar, J. N. Foerster, and S. Whiteson. QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning. In *ICML*, 2018.
- [16] M. Samvelyan, T. Rashid, C. S. de Witt, G. Farquhar, N. Nardelli, T. G. J. Rudner, C. Hung, P. H. S. Torr, J. N. Foerster, and S. Whiteson. The starcraft multi-agent challenge. In *AAMAS*, 2019.
- [17] S. Shalev-Shwartz, S. Shammah, and A. Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *CoRR*, abs/1610.03295, 2016.

- [18] K. Son, D. Kim, W. J. Kang, D. Hostallero, and Y. Yi. QTRAN: learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *ICML*, 2019.
- [19] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, and T. Graepel. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *AAMAS*, 2018.
- [20] M. Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *ICML*, 1993.
- [21] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. P. Agapiou, J. Schrittwieser, J. Quan, S. Gaffney, S. Petersen, K. Simonyan, T. Schaul, H. van Hasselt, D. Silver, T. P. Lillicrap, K. Calderone, P. Keet, A. Brunasso, D. Lawrence, A. Ekermo, J. Repp, and R. Tsing. Starcraft II: A new challenge for reinforcement learning. *CoRR*, abs/1708.04782, 2017.
- [22] J. Wang, Z. Ren, T. Liu, Y. Yu, and C. Zhang. QPLEX: duplex dueling multi-agent q-learning. *CoRR*, abs/2008.01062, 2020.

A StarCraft Multi-Agent Challenge

A.1 StarCraft Multi-Agent Challenge Setup

We used the open-source implementations of our baseline algorithms, including VDN [19], QMIX [15], QPLEX [22] and WQMIX [14] based on the PyMARL framework [16]. In Section 6, we choose eight maps of SMAC, including 1c3s5z, 2s_vs_1sc, 2s3z, 3s5z, 27m_vs_30m, 2c_vs_64zg, MMM2, bane_vs_bane, as the test environment. The snapshots of eight maps are shown in Figure 8, and the configurations of these eight maps are described in the Table 2. The hyper-parameters of SEM are illustrated in Table 3. For common hyper-parameters, we adopt the default implementation of PyMARL [16]. The replay buffer H stores episodes and its size is set to 5000. We sample $B = 32$ episodes uniformly from the replay buffer. The neural networks are all trained using RMSprop when the learning rate is set to 5×10^{-4} . After every 200 training episodes, the target networks are updated. During training, each agent a uses ϵ -greedy action selection for exploration. ϵ is annealed linearly from 1.0 to 0.05 over 50K time steps, and then is fixed as a constant. We set γ to 0.99 for all experiments. The enemy agents are controlled by built-in game AI. Depending on the exact scenario, each run of SEM-VDN takes between 11 to 26 hours using an Nvidia Geforce GTX 2080 Ti graphics card.

Table 2: SMAC challenges

Map Name	Ally Units	Enemy Units
1c3s5z	1 Colossus, 3 Stalkers & 5 Zealots	1 Colossus, 3 Stalkers & 5 Zealots
2s_vs_1sc	2 Stalkers	1 Spine Crawler
2s3z	2 Stalkers & 3 Zealots	2 Stalkers & 3 Zealots
3s5z	3 Stalkers & 5 Zealots	3 Stalkers & 5 Zealots
27m_vs_30m	27 Marines	30 Marines
2c_vs_64zg	2 Colossi	64 Zerglings
MMM2	1 Medivac, 2 Marauders & 7 Marines	1 Medivac, 2 Marauders & 8 Marines
bane_vs_bane	20 Zerglings & 4 Banelings	20 Zerglings & 4 Banelings

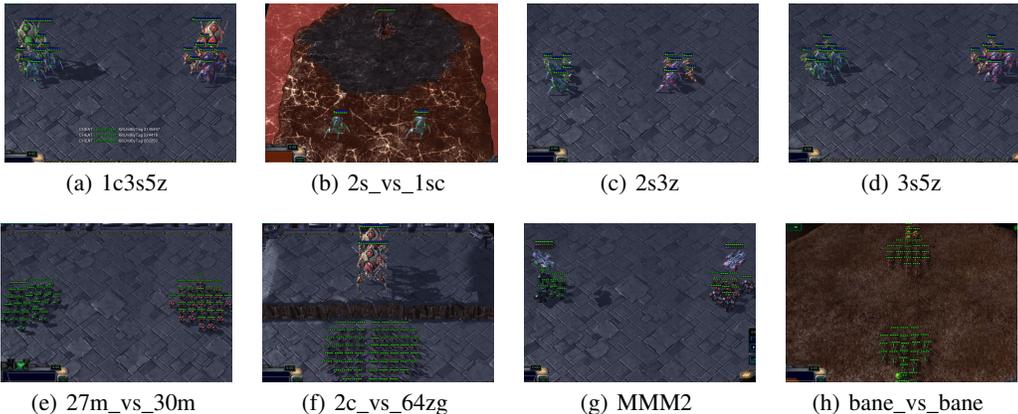


Figure 8: Snapshots of some StarCraft scenarios that we consider.

Table 3: The hyper-parameters in SEM on SMAC.

Hyper-parameter	Value
The balance coefficient λ	0.1
The dimension of random projection for state representation D	4
The size of lookup table $ Q^S $	10^6
The update frequency of lookup table $ M $	5000

A.2 SEM Combined with Baselines

We choose several value-decomposed MARL algorithms as baselines, including VDN [19], QMIX [15], QPLEX [22] and WQMIX [14]. We combine SEM with these value-decomposed MARL algorithms, respectively denoted as SEM-VDN, SEM-QMIX, SEM-QPLEX, and SEM-WQMIX.

A.2.1 SEM-VDN

The detailed architecture of VDN is illustrated in [19]. For SEM-VDN, the loss is as specified in Eq (5). The special hyper-parameters of SEM-VDN are illustrated in Table 3. The other hyper-parameters of SEM-VDN are the same as that in VDN [19].

A.2.2 SEM-QMIX

The architecture of QMIX has been illustrated in [15]. The hyper-parameters of QMIX are the same as that in QMIX [15]. In SEM-QMIX, the loss is shown in Eq (5). The special hyper-parameters of SEM-QMIX are the same as that in Table 3.

A.2.3 SEM-QPLEX

The overall architecture of QPLEX consists of two key components: a duplex dueling component and an individual action-value function, described as [22]. In the centralized training, the parameters of the whole network are learned by minimizing TD loss as specified in Eq (1). The loss of SEM-QPLEX is as specified in Eq (5). The special hyper-parameters of SEM-QPLEX are illustrated in Table 3. The other hyper-parameters of SEM-QPLEX are the same as that in QPLEX [22].

A.2.4 SEM-WQMIX

The architecture of WQMIX is described as [14]. The loss of WQMIX is as follows:

$$L(\theta) = \sum_{b=1}^B \sum_{t=1}^T w(s_t^b, \mathbf{u}_t^b) (Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b, s_t^b) - y_t^b)^2 + \sum_{b=1}^B \sum_{t=1}^T \left(\hat{Q}^*(s_t^b, \boldsymbol{\tau}_t^b, \mathbf{u}_t^b) - y_t^b \right)^2, \quad (6)$$

where $y_t^b = r_t^b + \gamma \hat{Q}^*(s_{t+1}^b, \boldsymbol{\tau}_{t+1}^b, \arg_{\mathbf{u}} Q_{tot}(s_{t+1}^b, \boldsymbol{\tau}_{t+1}^b, \mathbf{u}, s_{t+1}^b))$. For $w(s_t^b, \mathbf{u}_t^b)$, we use centrally-weighting function (CW), described as follows:

$$w(s_t^b, \mathbf{u}_t^b) = \begin{cases} 1 & y_t^b > \hat{Q}^*(s_t^b, \boldsymbol{\tau}_t^b, \hat{\mathbf{u}}^*) \text{ or } \mathbf{u}_t^b = \hat{\mathbf{u}}^*, \\ \alpha & \text{otherwise.} \end{cases} \quad (7)$$

We can combine our method, SEM, with WQMIX, called SEM-WQMIX. The loss of SEM-WQMIX is as follows:

$$\begin{aligned} L(\theta) = & (1 - \lambda) \sum_{b=1}^B \sum_{t=1}^T w(s_t^b, \mathbf{u}_t^b) (Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b, s_t^b) - y_t^b)^2 \\ & + \lambda \sum_{b=1}^B \sum_{t=1}^T w_e(s_t^b, \mathbf{u}_t^b) (Q_{tot}(\boldsymbol{\tau}_t^b, \mathbf{u}_t^b, s_t^b) - E_t^b)^2 \\ & + (1 - \lambda) \sum_{b=1}^B \sum_{t=1}^T \left(\hat{Q}^*(s_t^b, \boldsymbol{\tau}_t^b, \mathbf{u}_t^b) - y_t^b \right)^2 + \lambda \sum_{b=1}^B \sum_{t=1}^T \left(\hat{Q}^*(s_t^b, \boldsymbol{\tau}_t^b, \mathbf{u}_t^b) - E_t^b \right)^2, \end{aligned} \quad (8)$$

where $y_t, w(s, \mathbf{u})$ are the same as the loss of WQMIX and $E_t^b = r_t^b + \gamma Q^S(s_{t+1}^b)$. Similar to $w(s, \mathbf{u})$, $w_e(s, \mathbf{u})$ is defined as follows:

$$w_e(s_t^b, \mathbf{u}_t^b) = \begin{cases} 1 & E_t^b > \hat{Q}^*(s_t^b, \boldsymbol{\tau}_t^b, \hat{\mathbf{u}}^*) \text{ or } \mathbf{u}_t^b = \hat{\mathbf{u}}^*, \\ \alpha & \text{otherwise.} \end{cases} \quad (9)$$

In WQMIX and SEM-WQMIX, α is set to 0.75. The other hyper-parameters are the same as that in [14]. For the hyper-parameters of SEM in SEM-WQMIX, they are illustrated as Table 3.

B Results of SEM on SMAC

In Table 1, we show the mean and median scores of SEM combined with several value-decomposed methods (VDN, QMIX, QPLEX, WQMIX) on eight maps. Here, we illustrate the formula of the mean score and the formula of the median score. We denote $\text{Map}=\{1c3s5z, 2s_vs_1sc, 2s3z, 3s5z, 27m_vs_30m, 2c_vs_64zg, \text{MMM2}, \text{bane_vs_bane}\}$. The formula of the mean score and the formula of the median score are shown as follows:

$$\text{Mean Score } (t) = \text{Mean}(\{P_i^t\}_{i \in \text{Map}}), \quad (10)$$

$$\text{Median Score } (t) = \text{Median}(\{P_i^t\}_{i \in \text{Map}}), \quad (11)$$

where P_i^t is the mean test battle won for the map i at t time steps. The mean test won at 0.25M time steps is shown in Table 4 and the mean battle won at 0.5M time steps is shown in Table 5.

Table 4: The mean test battle won $P^{0.25M}$ at 0.25M time steps on eight maps of SMAC.

Method	1c3s5z	2s_vs_1sc	2s3z	3s5z
VDN	1%	66%	53%	1%
SEM-VDN	25%	66%	71%	24%
QMIX	34%	7%	61%	30%
SEM-QMIX	70%	28%	89%	66%
QPLEX	61%	99%	83%	31%
SEM-QPLEX	60%	77%	90%	77%
WQMIX	17%	63%	44%	23%
SEM-WQMIX	64%	17%	89%	63%

Method	27m_vs_30m	2c_vs_64zg	MMM2	bane_vs_bane
VDN	0%	0%	0%	85%
SEM-VDN	0%	0%	0%	99%
QMIX	0%	0%	0%	32%
SEM-QMIX	0%	0%	0%	99%
QPLEX	0%	0%	0%	87%
SEM-QPLEX	0%	0%	0%	99%
WQMIX	0%	0%	0%	61%
SEM-WQMIX	0%	0%	3%	100%

C Training Curves of SEM-VDN and SEM-QMIX

In the main paper, we show the training curves of SEM-VDN and SEM-QMIX on six maps. The training curves of SEM-VDN and SEM-QMIX on the other two maps, MMM2 and 2s_vs_1sc, are shown in Figure 9.

D Sensitivity to the Hyper-parameters

D.1 Sensitivity to the Hyper-parameter $|M|$

We choose bane_vs_bane to investigate the influence of episodic memory table update frequency. The lookup table Q^S updates when M is filled and then the set M is made empty. In other words, the lookup table updates after every $|M|$ time steps. $|M|$ is chosen from $\{1000, 2500, 5000, 10000\}$. The results are shown in Figure 10. We find that a larger $|M|$ might deteriorate the performance because the entries in the lookup table Q^S have not been updated timely with better value. When $|M| \in \{2500, 5000\}$, our methods perform better.

D.2 Sensitivity to the Hyper-parameter D

We investigate the influence of the dimension of random projection for state representation. D is chosen from $\{1, 2, 4, 10\}$. The results are shown in Figure 11. We can find that if D is too small,

Table 5: The mean test battle won $P^{0.5M}$ at 0.5M time steps on eight maps of SMAC.

Method	1c3s5z	2s_vs_1sc	2s3z	3s5z
VDN	14%	94%	81%	14%
SEM-VDN	67%	100%	89%	61%
QMIX	61%	31%	91%	62%
SEM-QMIX	83%	69%	91%	84%
QPLEX	87%	99%	98%	80%
SEM-QPLEX	87%	98%	99%	92%
WQMIX	48%	88%	78%	49%
SEM-WQMIX	72%	19%	95%	75%

Method	27m_vs_30m	2c_vs_64zg	MMM2	bane_vs_bane
VDN	0%	0%	0%	60%
SEM-VDN	1%	21%	0%	98%
QMIX	5%	0%	1%	24%
SEM-QMIX	6%	0%	6%	100%
QPLEX	4%	1%	0%	97%
SEM-QPLEX	14%	17%	0%	100%
WQMIX	0%	2%	0%	99%
SEM-WQMIX	5%	35%	15%	100%

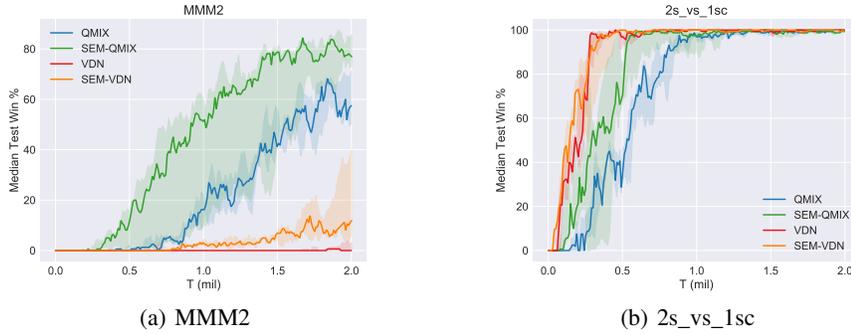


Figure 9: Results of our methods (SEM-VDN and SEM-QMIX) and baselines (VDN and QMIX), including the median performance as well as the 25-75% percentiles.

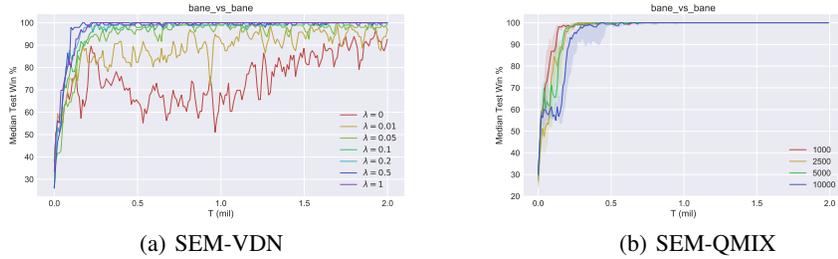
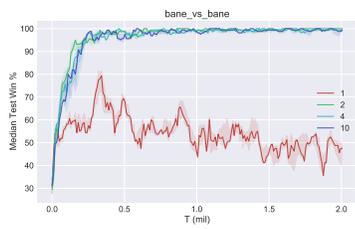
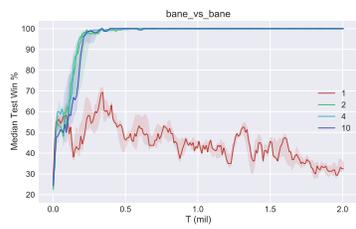


Figure 10: Median performance of SEM-QMIX and SEM-VDN when the episodic memory table is updated after every $|M|$ time steps. $|M|$ is chosen from $\{1000, 2500, 5000, 10000\}$. The results are summarized over 5 random runs.

performance degrades considerably and if D is too large, the storage overhead is significant. When $D \in \{2, 4\}$, it is a reasonable choice.



(a) SEM-VDN



(b) SEM-QMIX

Figure 11: Median performance of SEM-QMIX and SEM-VDN when the dimension of the random projection D is chosen from $\{1, 2, 4, 10\}$.