# Dual Self-Awareness Value Decomposition Framework without Individual Global Max for Cooperative MARL

**Zhiwei Xu, Bin Zhang, Dapeng Li, Guangchong Zhou, Zeren Zhang, Guoliang Fan**
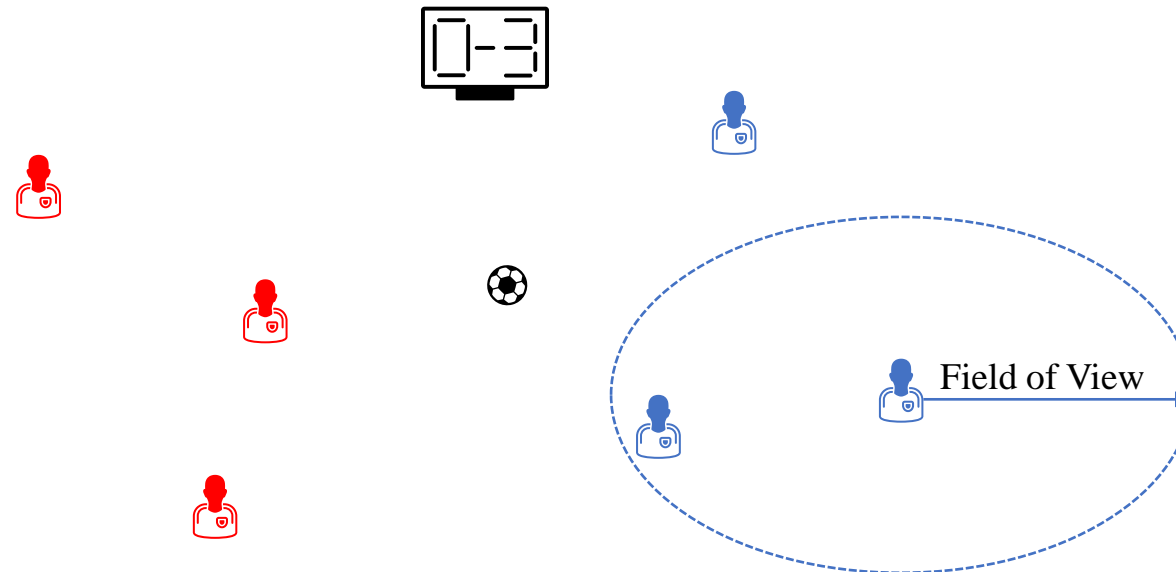
Institute of Automation, Chinese Academy of Sciences

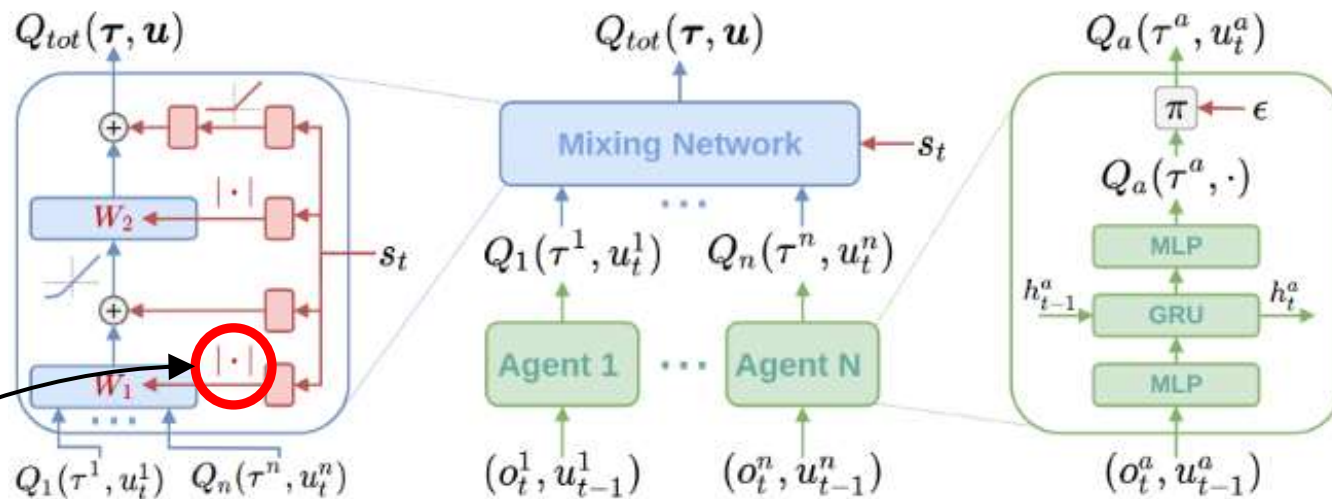School of Artificial Intelligence, University of Chinese Academy of Sciences

# Dec-POMDP

- A Dec-POMDP is a "Decentralized Partially Observable Markov Decision Process"
- It's a bunch of "agents" that are working together for a common reward
- Each agent only takes a local observation of the environment
- A fully cooperative multi-agent task can be described as a tuple $\langle S, U, p, r, Z, O, n, \gamma \rangle$



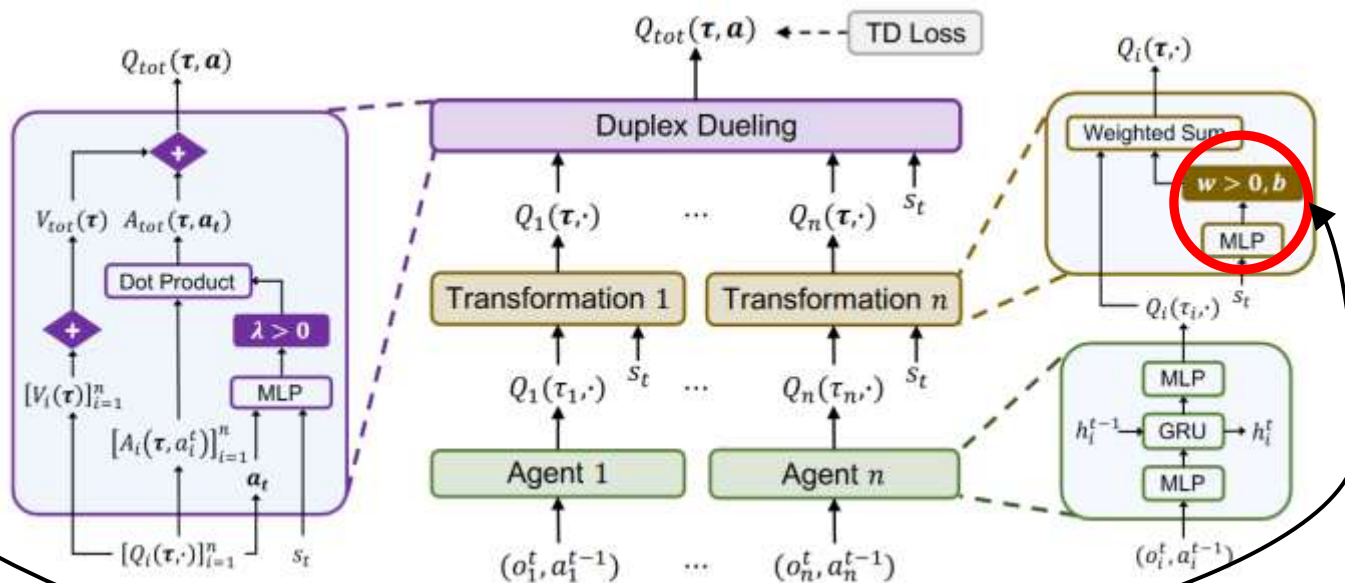Field of View

# Motivation

**Individual Global Max:**

$$\underset{\mathbf{u}}{\operatorname{argmax}} Q_{tot}(\boldsymbol{\tau}, \mathbf{u}) = \begin{pmatrix} \operatorname{argmax}_{u^1} Q_1(\tau^1, u^1) \\ \vdots \\ \operatorname{argmax}_{u^n} Q_n(\tau^n, u^n) \end{pmatrix}$$
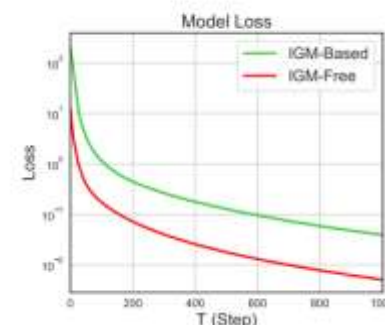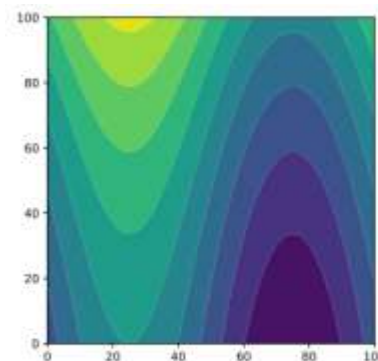
Absolute Operator

# Motivation

**Example**

$$Q_i(u^i) = \frac{u^i}{100}, \quad \forall i \in \{1, 2\}$$

$$R(u^1, u^2) = \sin\left(2\pi Q_1(u^1)\right) + \exp\left(Q_2(u^2)\right)$$

**Hard to obtain Max Q —— (IGM-Free)**

**OR ?**

**Restrict the set of functions that can be represented**

**——(IGM-Based)**

IGM-Based

(a) 200 step  (b) 400 step  (c) 600 step  (d) 800 step  (e) 1000 step

IGM-Free

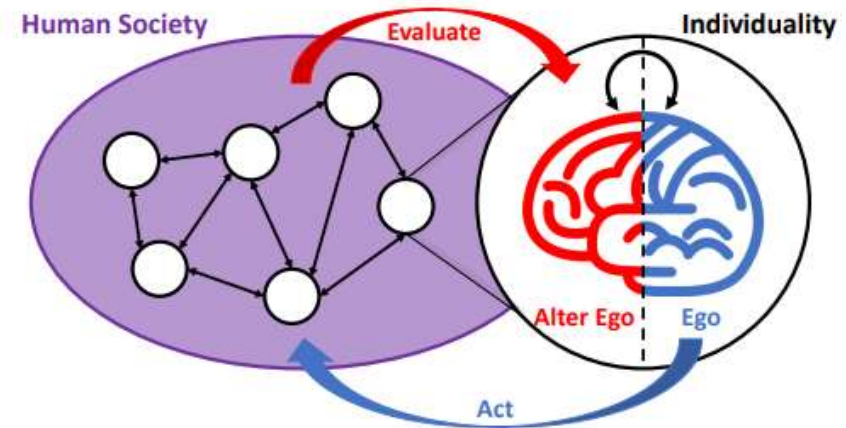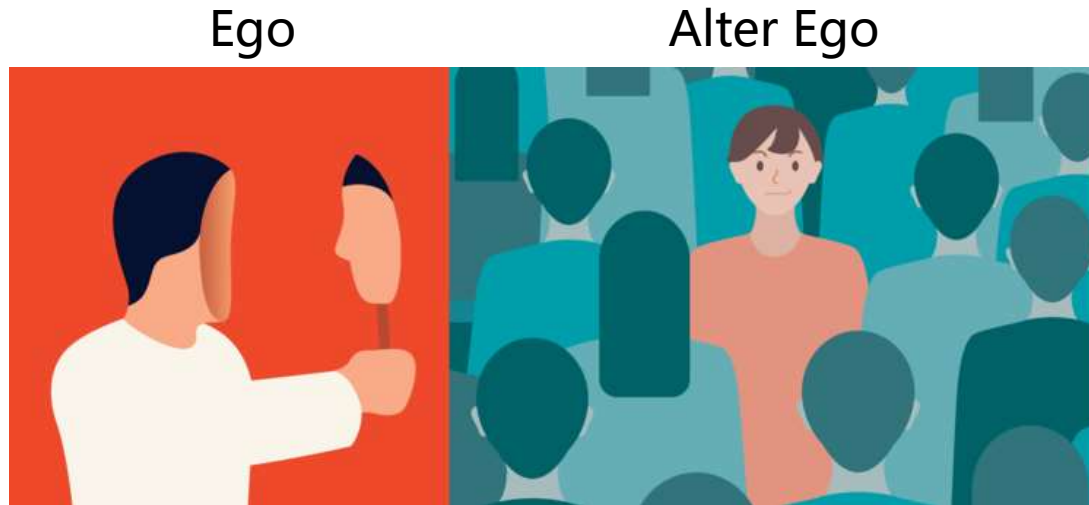(f) 200 step  (g) 400 step  (h) 600 step  (i) 800 step  (j) 1000 step

# Dual Self-Awareness Hypothesis

Ego        Alter Ego



There are concepts of **ego** and **alter ego** in psychology. The ego usually refers to the conscious part of the individual, and Freud considered the ego to be the executive of the personality. Some people believe that an alter ego pertains to a different version of oneself from the authentic self. Others define the alter ego in more detail as the evaluation of the self by others.

Enlightened by these psychological concepts, we propose a novel MARL algorithm, **D**ual self-**A**wareness **V**alue d**E**composition (DAVE).

# Related Work

# Dual Self-Awareness Framework

The objective of the IGM-free value decomposition method is as follows:

$$\arg\max_{\pi} Q_{\text{tot}}(s, u),$$

Without the IGM assumption, it is **NP-hard** because it cannot be solved and verified in polynomial time. Therefore, in our proposed dual self-awareness framework, each agent has an additional policy network to assist the value function network to find the action corresponding to the optimal joint policy.

**Ego Policy Model** & **Alter Ego Value Function Model**

# Dual Self-Awareness Framework

So we modify the objective of the cooperative multi-agent value decomposition problem from previous equation to

$$\arg\max_{\pi^{ego}} Q_{tot}^{alter}(s, u^{ego}), \qquad \text{s.t.} \quad u^{ego} \in U^{ego},$$

where $U^{ego} := \{u_i^{ego} \sim \pi^{ego}(s)\}_{i=1}^{M}$ and M is the number of samples.

Define

$$u^{\star} = \arg\max_{u^{ego}} Q_{tot}^{alter}(s, u^{ego})$$

The loss function for the joint ego policy can be written as: $\mathcal{L}_{ego} = -\log \pi^{ego}(u^{\star} \mid s).$
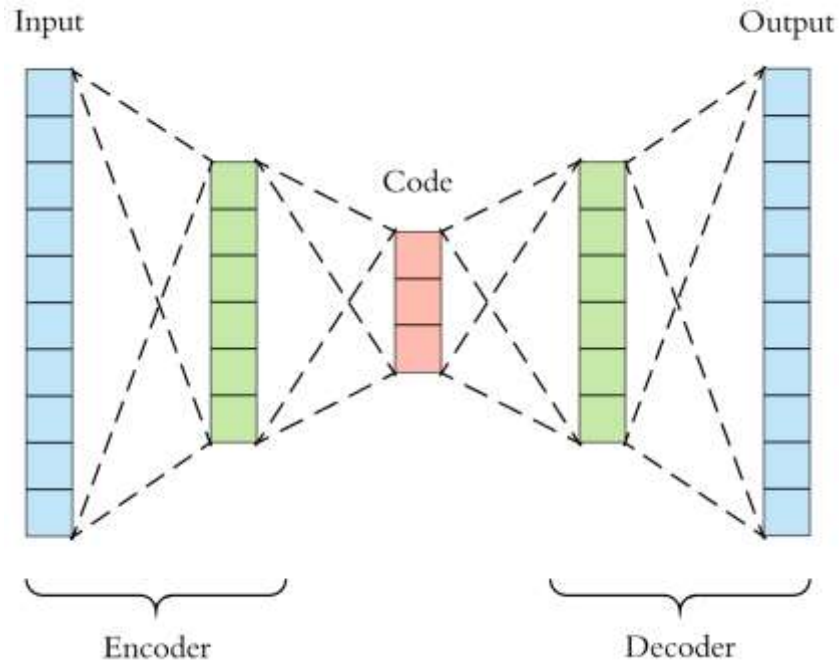
**Proposition A.1.** *As long as the ego policy assigns non-zero probabilities to all actions, this method will approach the objective described by Equation (1) as the number of samples $M$ increases.*

*Proof.* Let $(u^a)^*$ denote the individual actions corresponding to the global optimal joint state-action value function under state $s$, and the optimal joint action is $u^* = \{(u^1)^*, \ldots, (u^n)^*\}$. Then the probability that the sampling procedure draws $u^*$ is expressed as:

$$p(u^*) = 1 - (1 - \pi^{ego}(u^*|s))^M$$

$$= 1 - (1 - \prod_{i=1}^{n} \pi_a^{ego}((u^a)^*|\tau^a))^M,$$

where $\pi_a^{ego}(\cdot|\cdot) \in (0, 1)$ is true for any action. So the second term $(1 - \pi^{ego}(u^*|s))^M \in (0, 1)$ in the equation decreases as $M$ increases, indicating that $p(u^*)$ is positively correlated with the sample size $M$. $\square$

# Anti-Ego Exploration



Input ... Code ... Output
Encoder ... Decoder

Our method leverages the fact that the auto-encoder cannot effectively encode novel data and accurately reconstruct it.

The relationship between the ego and anti-ego policy of each agent is as follows:

$$\pi_a^{\text{ego}}(\tau^a) = \text{Softmax}(f(\tau^a)),$$
$$\hat{\pi}_a^{\text{ego}}(\tau^a) = \text{Softmin}(f(\tau^a)),$$

**Use Autoencoder to filter out the full-explored state-action pairs.**

$$\mathcal{L}_{\text{recon}}(s, \boldsymbol{u}) = \text{MSE}(s, s') + \sum_{a=1}^{n} \text{CE}(u^a, (u^a)'),$$

$$\hat{\boldsymbol{u}}^{\star} = \arg\max_{\hat{\boldsymbol{u}}^{\text{ego}}} \mathcal{L}_{\text{recon}}(s, \hat{\boldsymbol{u}}^{\text{ego}}),$$

where $\hat{\boldsymbol{u}}^{\text{ego}} \in \hat{\boldsymbol{U}}^{\text{ego}} := \{\hat{\boldsymbol{u}}_i^{\text{ego}} \sim \hat{\boldsymbol{\pi}}^{\text{ego}}(s)\}_{i=1}^{M}$

$$\mathcal{L}_{\text{ego}} = -\left(\log \boldsymbol{\pi}^{\text{ego}}(\boldsymbol{u}^{\star} \mid s) + \lambda \log \boldsymbol{\pi}^{\text{ego}}(\hat{\boldsymbol{u}}^{\star} \mid s)\right)$$
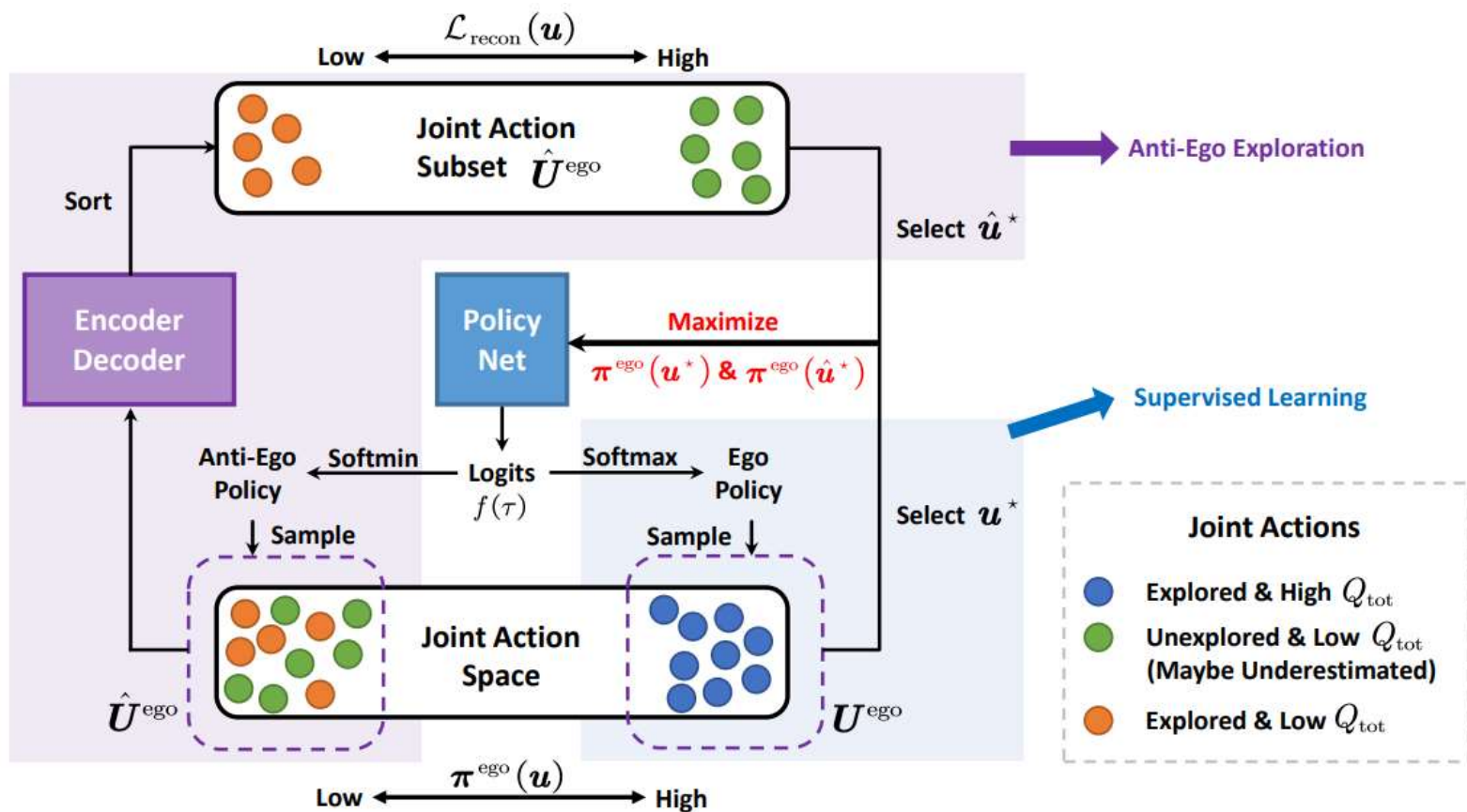
# Details of the Ego Policy Update



Figure 9: Diagram of the update of ego policy, where $U^{\text{ego}} := \{u_i^{\text{ego}} \sim \pi^{\text{ego}}(s)\}_{i=1}^M$ and $\hat{U}^{\text{ego}} := \{\hat{u}_i^{\text{ego}} \sim \hat{\pi}^{\text{ego}}(s)\}_{i=1}^M$. Dots indicate different joint actions. The selected actions are $u^\star = \arg\max_{u^{\text{ego}}} Q_{\text{tot}}^{\text{alter}}(s, u^{\text{ego}})$ and $\hat{u}^\star = \arg\max_{\hat{u}^{\text{ego}}} \mathcal{L}_{\text{recon}}(s, \hat{u}^{\text{ego}})$.

# Algorithmic Description

---

**Algorithm 1** Training Procedure for DAVE

---

**Hyperparameters**: Sample size $M$, discount factor $\gamma$, exploration coefficients $\lambda_{\text{init}}$

Initialize the parameters of the neural networks shown in Figure 2

 1: **for** each episode **do**
 2:     Get the global state $s_1$ and the local observations $z_1 = \{z_1^1, z_1^2, \ldots, z_1^n\}$ of all agents
 3:     **for** $t \leftarrow 1$ to $T - 1$ **do**
 4:         **for** $a \leftarrow 1$ to $n$ **do**
 5:             Select action $u_t^a$ according to the ego policy $\pi_a^{\text{ego}}$
 6:         **end for**
 7:         Carry out the joint action $\boldsymbol{u}_t = \{u_t^1, \ldots, u_t^n\}$
 8:         Get the global reward $r_{t+1}$, the next local observations $z_{t+1}$, and the next state $s_{t+1}$
 9:     **end for**
10:     Store the episode in the replay buffer $\mathcal{D}$
11:     Sample a batch of episodes $\mathcal{B} \sim \text{Uniform}(\mathcal{D})$
12:     Sample and obtain the joint action set $\boldsymbol{U}^{\text{ego}} := \{\boldsymbol{u}_i^{\text{ego}} \sim \boldsymbol{\pi}^{\text{ego}}(s)\}_{i=1}^{M}$ for each trajectory in $\mathcal{B}$
13:     Update the parameters of the alter ego value function and the IGM-free mixing network according Equation (5)
14:     Obtain the anti-ego policies $\hat{\pi}_a^{\text{ego}}$ of each agent
15:     Sample and obtain $\hat{\boldsymbol{U}}^{\text{ego}} := \{\hat{\boldsymbol{u}}_i^{\text{ego}} \sim \hat{\boldsymbol{\pi}}^{\text{ego}}(s)\}_{i=1}^{M}$ by sampling $M$ times from the anti-ego policy for each state $s$ in $\mathcal{B}$
16:     Find the most novel joint action $\boldsymbol{u}^{\star}$ for each state $s$ in $\mathcal{B}$ according Equation (8)
17:     Update the parameters of the ego policy according Equation (9)
18:     Update the parameters of the auto-encoder according Equation (6)
19:     Update the parameters of the target network periodically
20: **end for**

---

# Experiment (Matrix Game)



Figure 3: Payoffs of the two matrix games.

| | k=0 | | | k=6 | | | | k=7.5 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | -12 | 0 | **8** | -12 | 0 | 6 | **8** | -12 | 0 | 7.5 | **8** |
| MADDPG | 35% | 45% | 20% | 27% | 34% | 31% | 8% | 27% | 35% | 34% | 4% |
| MAPPO | 0% | 100% | 0% | 0% | 0% | 100% | 0% | 0% | 0% | 100% | 0% |

Table 2: Proportion of different convergence results in Matrix Game I.

| | k=0 | | | k=25 | | | | k=100 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 2 | **10** | -25 | 0 | 2 | **10** | -100 | 0 | 2 | **10** |
| MADDPG | 51% | 0% | 49% | 3% | 45% | 45% | 7% | 1% | 20% | 76% | 3% |
| MAPPO | 0% | 0% | 100% | 0% | 0% | 100% | 0% | 0% | 0% | 100% | 0% |

Table 3: Proportion of different convergence results in Matrix Game II.
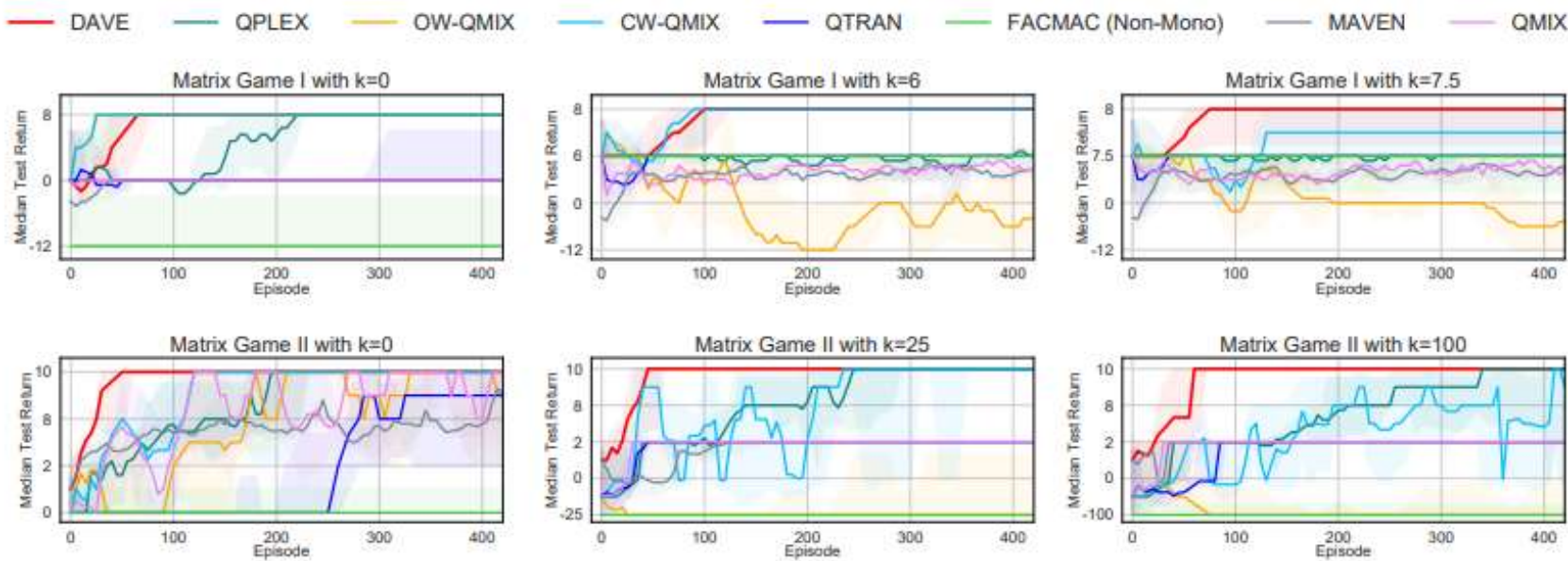


Figure 4: The learning curves of DAVE and other baselines on the matrix games. Note that the ordinates are non-uniform.
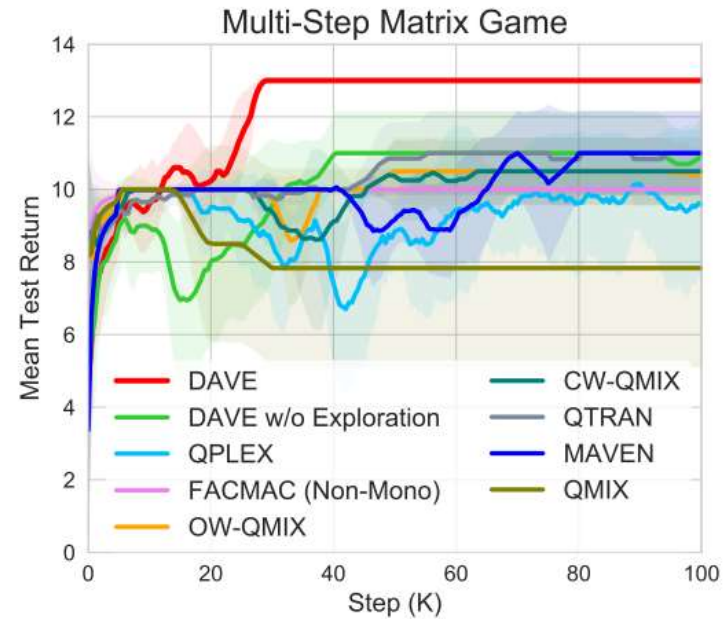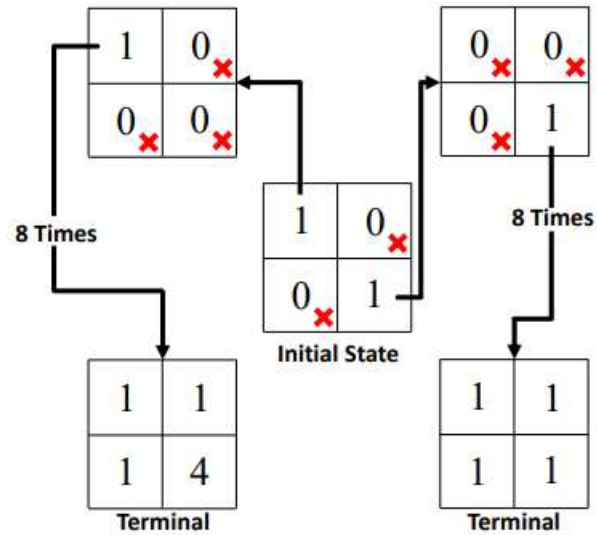
# Experiment (Matrix Game)



Figure 5: **Left**: Illustrations of the multi-step matrix game. **Right**: Performance over time on the multi-step matrix game.

# Experiment (SMAC)



Figure 6: Performance comparison with baselines in different scenarios.

# Experiment (SMACv2)



Figure 13: Comparisons of median win rate for variants of QMIX on SMACv2.



Figure 14: Comparisons of median win rate for variants of QPLEX on SMACv2.

# Experiment (MA-MuJoCo)



(a) Hopper (3)  (b) HalfCheetah (6)  (c) Humanoid (17)

Figure 16: Illustration of benchmark tasks in Multi-Agent MuJoCo.

Contains

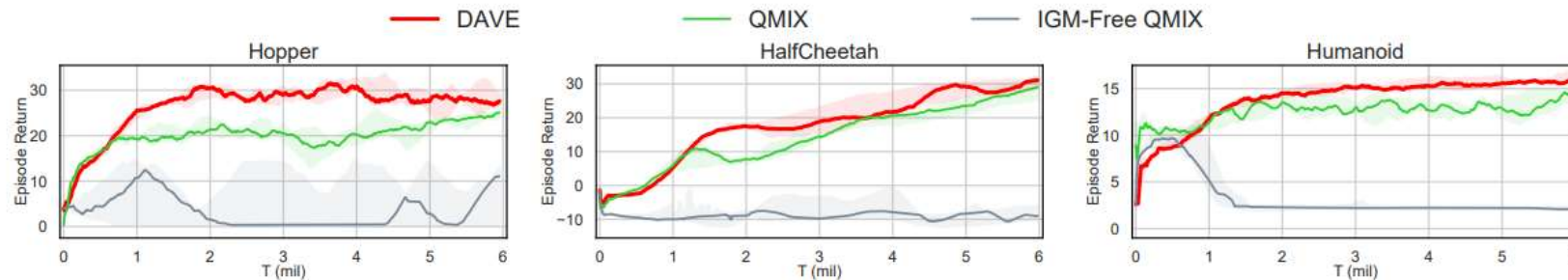$$31^{17} \approx 2.26 \times 10^{25}$$

joint actions



Figure 17: Median episode return on different MA-MuJoCo tasks. Note that only QMIX follows the IGM assumption.

# Experiment



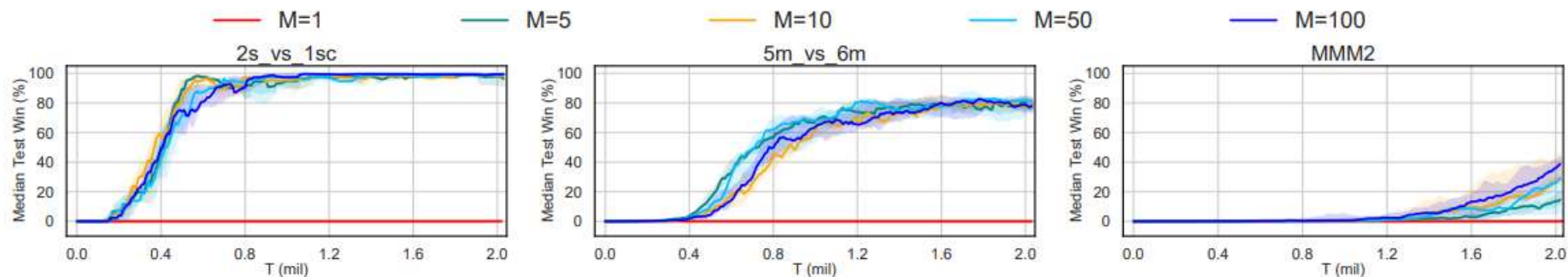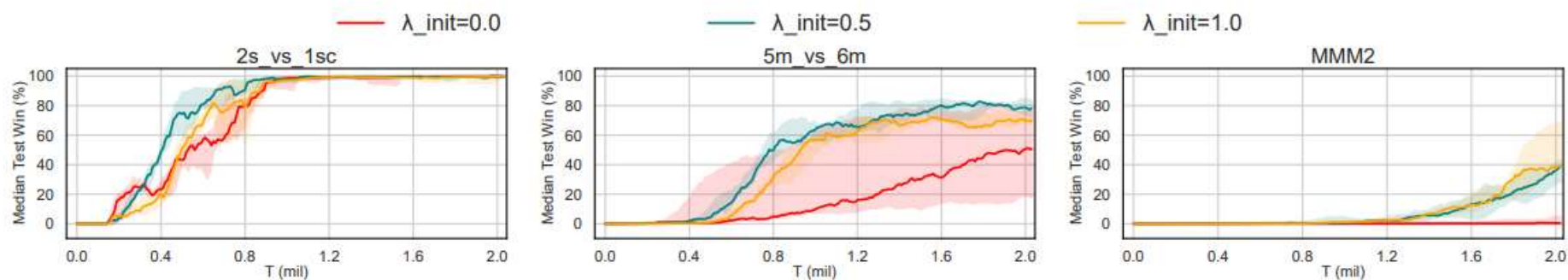Figure 7: Influence of the sample size $M$ for DAVE.



Figure 8: Results of DAVE with different $\lambda_{\text{init}}$. The action space for the three scenarios gradually increases from left to right.

# Summary

**Pros**

The <u>first</u> multi-agent value decomposition method that <u>completely abandons IGM</u>

Can be applied to most IGM-based value decomposition methods and <u>turn them into IGM-free ones</u>

Can achieve desirable performance in various cooperative tasks, including non-monotonic and complex tasks

Use Autoencoder to avoid the algorithm becoming stuck in a local optimum

**Cons**

Choice of $\lambda$

May be hard to solve the tasks with large action spaces