



西安交通大学
XI'AN JIAOTONG UNIVERSITY

Greedy-based Value Representation for Optimal Coordination in Multi-agent Reinforcement Learning

Lipeng Wan / Zeyang Liu / Xingyu Chen / Xuguang Lan* / Nannin Zheng
School of AI, Xi'an Jiaotong University
2022/6/21

1. Introduction # Fully-cooperative multi-agent reinforcement learning
2. Background # Challenge: relative overgeneralization (RO)
3. Analysis # Investigation of RO
4. Methodology # Solution: greedy-based value representation
5. Experiments # Matrix game / Predator-prey / Starcraft multi-agent challenge
6. Conclusions

Introduction



Fully-cooperative MARL

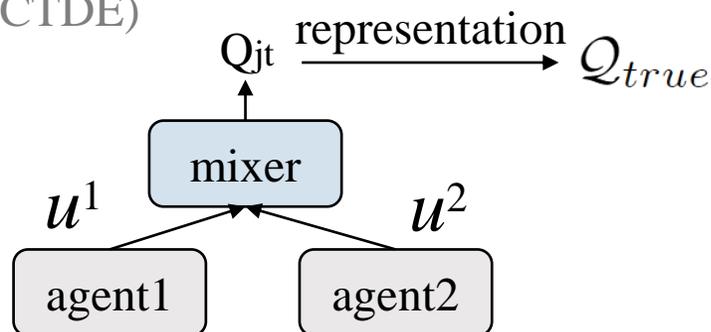
agent 1

		u_0	u_1	u_2
agent 2	u_0	r_1	r_2	r_3
	u_1	r_4	r_5	r_6
	u_2	r_7	r_8	r_9

shared rewards

Value decomposition

Credit assignment for Centralized Training with Decentralized Execution (CTDE)



Independent Global Max (IGM)

identity of joint greedy action and the set of individual greedy actions

Q_{jt}	u^1	u^2	
1	0	0	0.6
0	0	0	0
0	0	2	0.8
	0.4	0	1.2

true Q value function

Q_{true}

Q_{11}	Q_{12}	Q_{13}
Q_{21}	Q_{22}	Q_{23}
Q_{31}	Q_{32}	Q_{33}

Q_{jt}	u^1	u^2
1	0	0
0	0	0
0	0	2
	1	0
	0	2

VDN:

Linear Value Decomposition (LVD)

$$Q_{jt} = U^1 + U^2$$

QMIX:

Monotonic Value Decomposition (MVD)

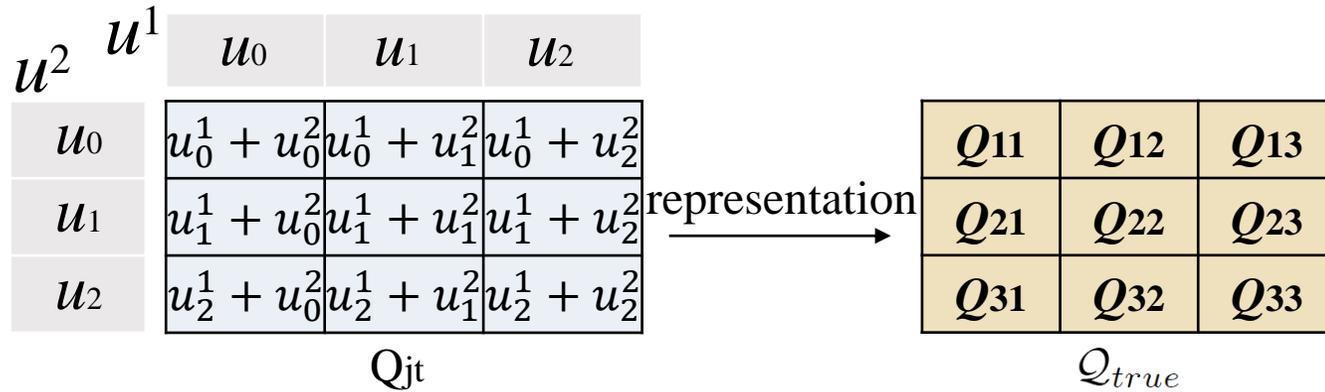
$$Q_{jt} = \omega_1 \cdot U^1 + \omega_2 \cdot U^2$$

* u^i : utility function
 Q_{jt} : joint Q value function

Background



Incomplete Representation Capability (IRC) of the joint Q value function



representation problem: solving an **overdetermined** equation system

$$\{u_i^1 + u_j^2 = Q_{ij}\}_{i,j \in [1,m]}$$

number of variables: $m \cdot n$ (6)

number of equations: m^n (9)

equation num > variable num

no exact solution

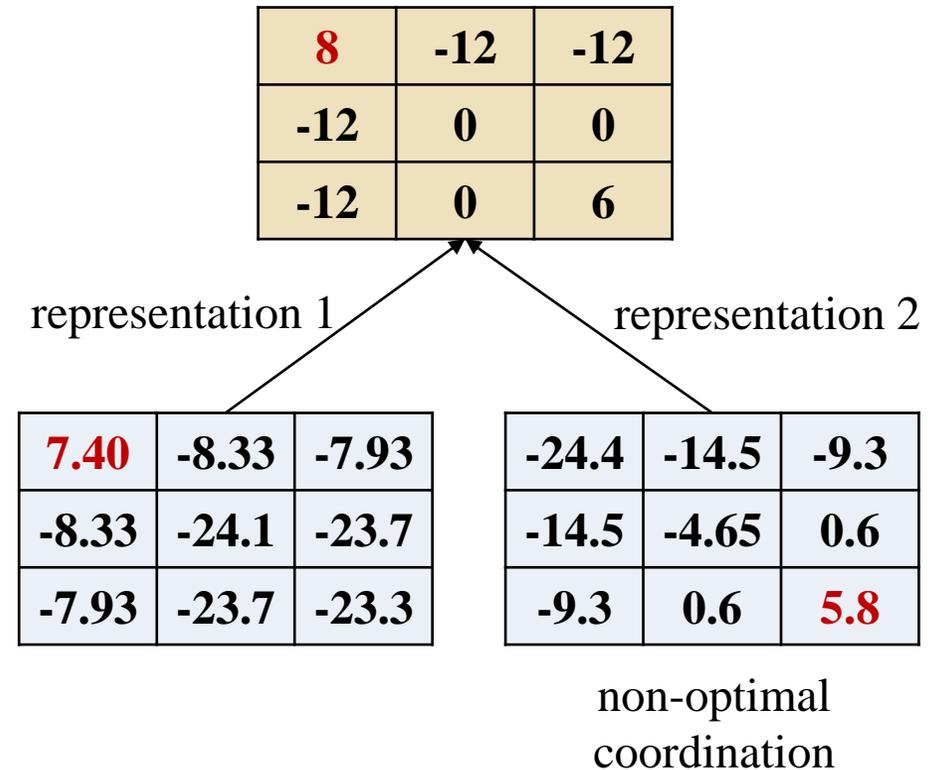


Q_{jt} is unable to converge to Q_{true}



Q_{jt} suffers from IRC

Relative Overgeneralization (RO)
multiple potential representations



* m: number of options of individual action space
n: number of agents

Existing solutions to RO

8	-12	-12
-12	0	0
-12	0	6

1. **biased representation**

(Weighted QMIX)

place more weight on the representation of good samples

7.9	-8.33	-7.93
-8.33	-24.1	-23.7
-7.93	-23.7	-23.3

concerns:

- rely on heuristic parameters
- can not ensure to solve the RO

2. **complete representation**

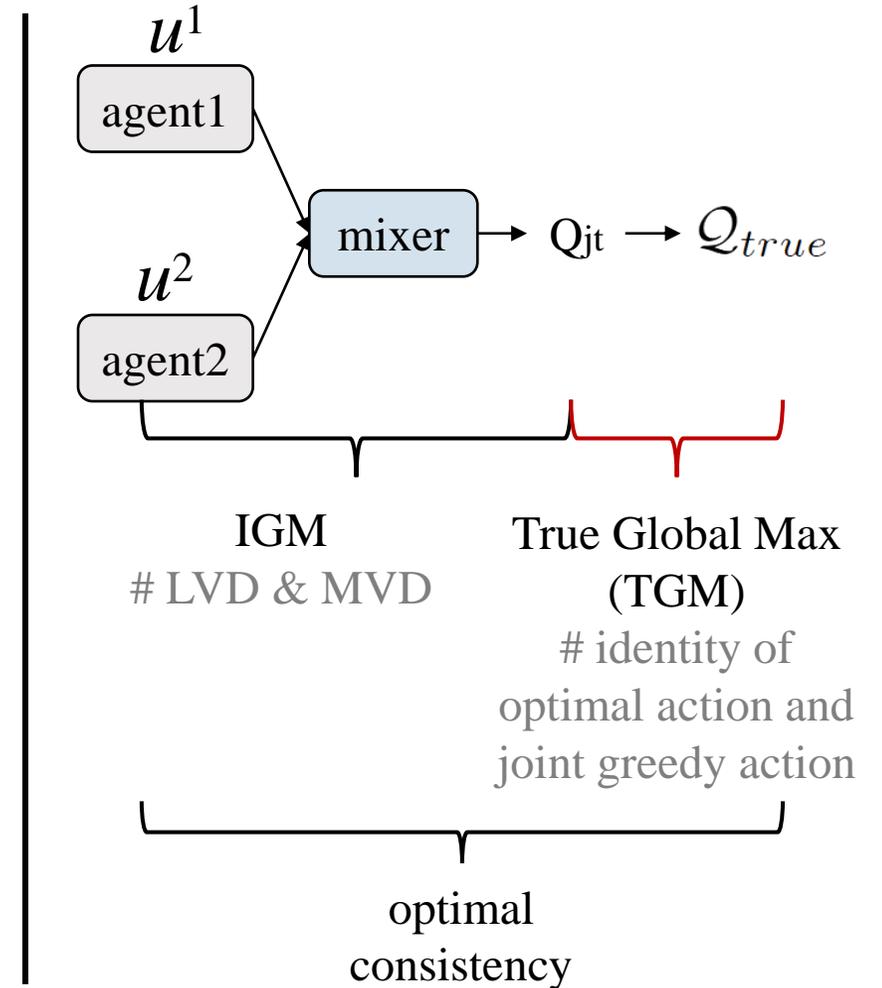
(QPLEX, Qtran)

design joint Q value functions with complete representation capability

8	-12	-12
-12	0	0
-12	0	6

concerns:

- CRC is complex and unnecessary



Analysis

Conditions of the TGM principle for LVD and MVD

According to the equation system

$$\{U_i^1 + U_j^2 = Q_{ij}\}_{i,j \in [1,m]}$$

Under the e-greedy visitation

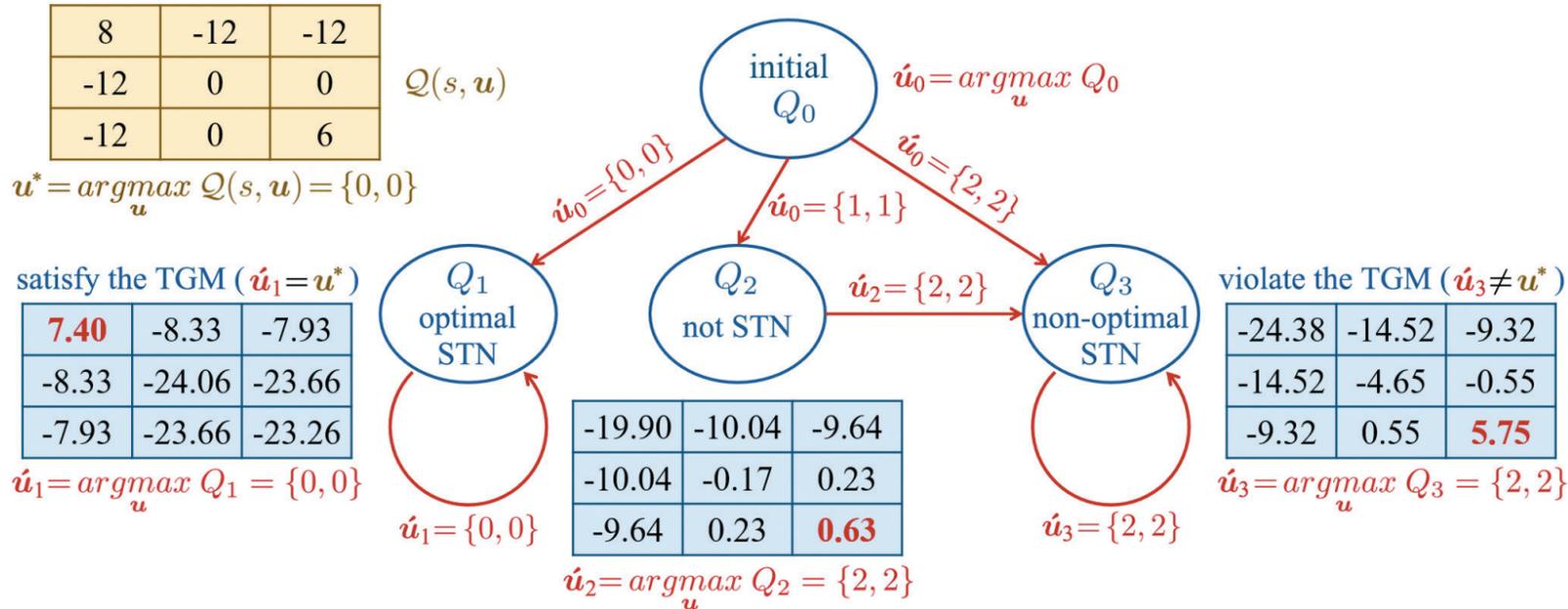
(\dot{i}, \dot{j} denotes the greedy action of agent 1, 2)

$$U_i^1 = \frac{\epsilon}{m} \sum_{k=1}^m (Q_{ik} - U_k^2) + (1 - \epsilon)(Q_{i\dot{j}} - U_j^2)$$

$$U_j^2 = \frac{\epsilon}{m} \sum_{k=1}^m (Q_{kj} - U_k^1) + (1 - \epsilon)(Q_{i\dot{j}} - U_i^1)$$

The joint Q value function can be acquired

$$Q(u_i^1, u_j^2, \tau) = \frac{\epsilon}{m} \sum_{k=1}^m (Q_{ik} + Q_{kj}) + (1 - \epsilon)(Q_{i\dot{j}} + Q_{i\dot{j}}) - \frac{\epsilon^2}{m^2} \sum_{i=1}^m \sum_{j=1}^m Q_{ij} - \frac{\epsilon(1 - \epsilon)}{m} \sum_{k=1}^m (Q_{i\dot{k}} + Q_{k\dot{j}}) - (1 - \epsilon)^2 Q_{i\dot{j}}$$

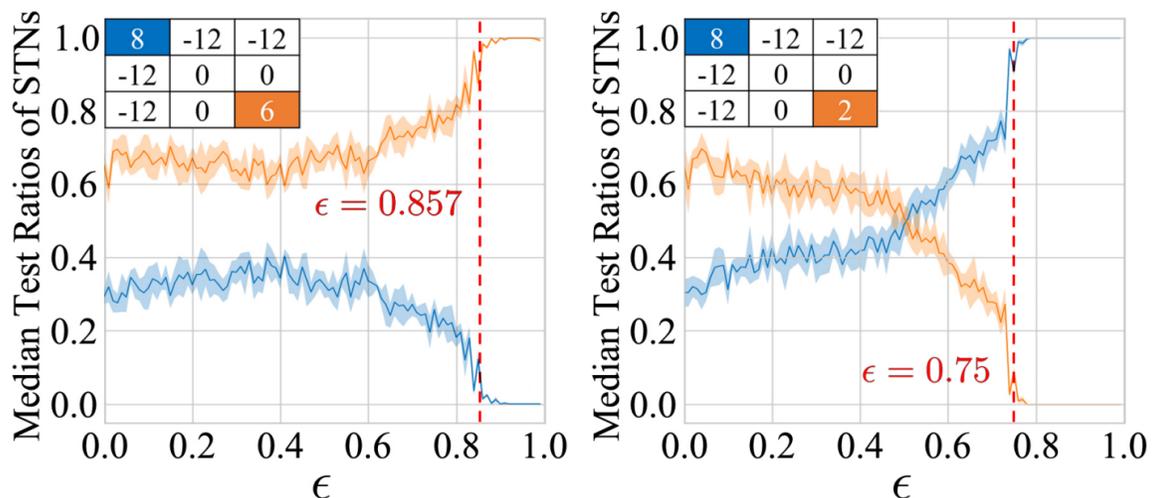


- The joint Q value function changes with the joint greedy action
- The joint Q value function is the same for LVD and MVD

* STN: Self-Transition Node

Analysis

- Improving exploration helps to eliminate extra STNs



Ensuring TGM = the optimal node is the only STN

sufficient condition:

Condition 1
 $Q(\mathbf{u}_s, \tau) < Q(\hat{\mathbf{u}}, \tau)$ s.t. $Q(s, \mathbf{u}_s) \leq Q(s, \hat{\mathbf{u}})$

Condition 2
 $Q(\mathbf{u}_s, \tau) > Q(\hat{\mathbf{u}}, \tau)$ s.t. $Q(s, \mathbf{u}_s) > Q(s, \hat{\mathbf{u}})$

$$\mathbf{u}^* = \{0, 0\}$$

8	-12	-12
-12	0	0
-12	0	6

$$Q(s, \mathbf{u})$$

optimal node ($\hat{\mathbf{u}} = \mathbf{u}^*$)

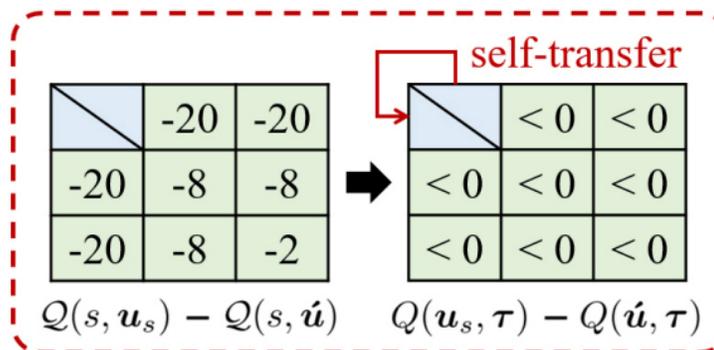
Condition 1

$$Q(\mathbf{u}_s, \tau) < Q(\hat{\mathbf{u}}, \tau)$$

$$s.t. Q(s, \mathbf{u}_s) \leq Q(s, \hat{\mathbf{u}})$$

STN

- greedy action
- actions satisfies condition 1
- actions satisfies condition 2



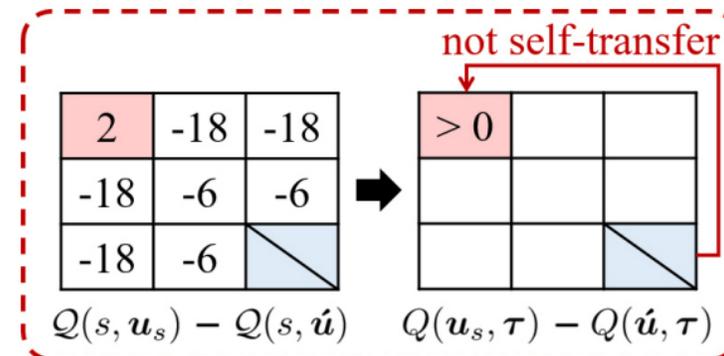
non-optimal node ($\hat{\mathbf{u}} \neq \mathbf{u}^*$)

Condition 2

$$Q(\mathbf{u}_s, \tau) > Q(\hat{\mathbf{u}}, \tau)$$

$$s.t. Q(s, \mathbf{u}_s) > Q(s, \hat{\mathbf{u}})$$

not STN



Inferior Target Shaping (ITS)

Instead of representing the specific Q value of an inferior action, reshaping it with a Q value “no better than current greedy”

$$Q_{its}(s, \mathbf{u}) = Q(\hat{\mathbf{u}}, \tau) - \alpha|Q(\hat{\mathbf{u}}, \tau)|$$

$$s.t. Q(s, \mathbf{u}) \leq Q(s, \hat{\mathbf{u}}) * (1 + e_{Q0}) \text{ and } \mathbf{u} \neq \hat{\mathbf{u}}$$

8	-12	-12
-12	0	0
-12	0	6

$Q(s, \mathbf{u})$

$Q_{its}(s, \mathbf{u})$	(a)	(b)	(c)
$\hat{\mathbf{u}}$	{0,0}	{1,1}	{2,2}
$Q(\hat{\mathbf{u}}, \tau)$	7.8	-0.2	5.8
$Q_{its}(s, \mathbf{u})$ for inferior actions: $Q(\hat{\mathbf{u}}, \tau) - \alpha Q(\hat{\mathbf{u}}, \tau) $	7.02	-0.22	5.22

- greedy action
- inferior action
- superior action

8	7.02	7.02
7.02	7.02	7.02
7.02	7.02	7.02

(a)

8	-0.22	-0.22
-0.22	0	0
-0.22	0	6

(b)

8	5.22	5.22
5.22	5.22	5.22
5.22	5.22	6

(c)

$$Q(\mathbf{u}_s, \tau) - Q(\hat{\mathbf{u}}, \tau)$$

$$= n(\eta_1 - \eta_2) [Q(s, \hat{\mathbf{u}}) - (1 - \alpha)Q(\hat{\mathbf{u}}, \tau)] + n\eta_1 e_Q Q(s, \hat{\mathbf{u}}) \quad (1)$$

- **Optimal node** (fig (a))

$$Q(\mathbf{u}_s, \tau) - Q(\hat{\mathbf{u}}, \tau)$$

$$= n(\eta_1 - \eta_2) \frac{\alpha}{w} Q(\hat{\mathbf{u}}, \tau) + n\eta_1 e_Q Q(s, \hat{\mathbf{u}}) < 0$$

The optimal node is always STN

- **Non-optimal node** (fig (b), (c)). Let Eq.1>0:

$$\text{We have } \frac{\eta_1}{\eta_2} > \frac{\alpha}{\alpha + e_{Q0}} \quad (2)$$

Improving the probability of superior actions helps removing non-optimal STNs

$$* e_Q = \frac{Q_{its}(s, \mathbf{u}_s) - Q(s, \hat{\mathbf{u}})}{Q(s, \hat{\mathbf{u}})} \quad \eta_1 = \left(\frac{\epsilon}{m}\right)^{n-1} \quad \eta_2 = \left(1 - \epsilon + \frac{\epsilon}{m}\right)^{n-1}$$

Raising the probabilities of superior actions

1. Improving exploration • **Too large exploration rate**

$$\frac{\eta_1}{\eta_2} > \frac{\alpha}{\alpha + e_{Q0}} \quad \text{equals} \quad \epsilon > \frac{m}{\left(\frac{e_{Q0}}{\alpha}\right)^{\frac{1}{n-1}} + 1 + m - 1}$$

2. Reweighting samples

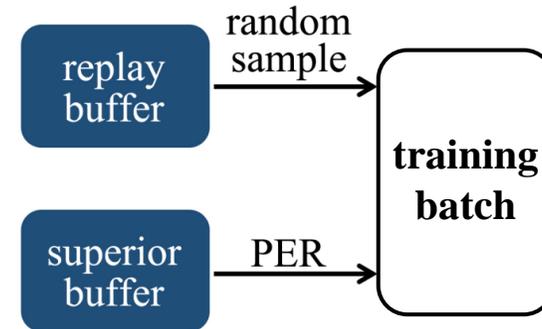
required sample weight $w > \frac{\alpha(\eta_2 - \eta_1)}{e_{Q0}\eta_1}$

m^n	3^2	3^3	3^4
w_0 (Eq.57)	3.60	50.32	659.50
$Q(\mathbf{u}_s, \boldsymbol{\tau}) - Q(\hat{\mathbf{u}}, \boldsymbol{\tau})$ (Eq.56)	0.01 ± 0.06 (0)	-0.02 ± 0.30 (0)	-0.48 ± 0.75 (0)
Test $Q(\hat{\mathbf{u}}, \boldsymbol{\tau})$	5.95 ± 0.02	5.90 ± 0.06	5.93 ± 0.03

- Required sample weight grows exponentially as the number of agents

Superior Experience Replay (SER)

remove the correlation of a superior action's proportion in training batch with its exploration probability



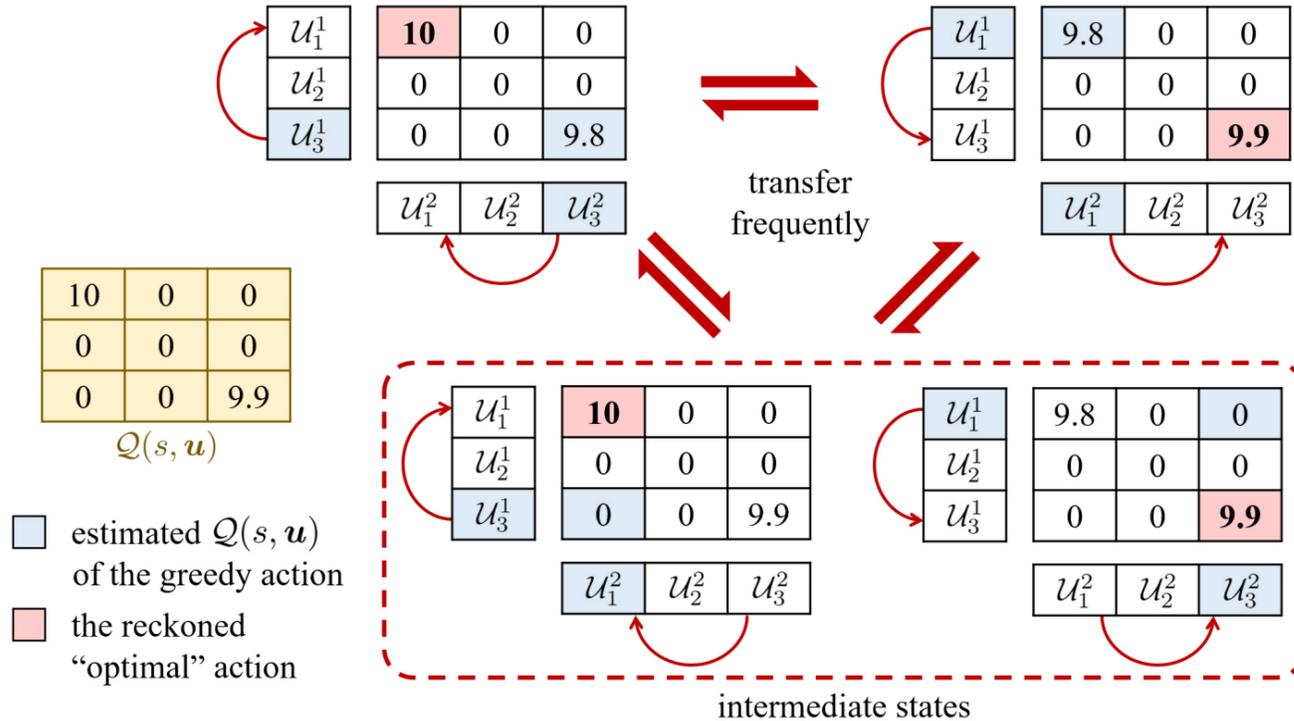
Relative weight for samples from superior buffer

$$w_{ser} > \frac{\alpha}{e_{Q0}} (\eta_2 - \eta_1) \eta_s - \eta_1 \eta_s$$

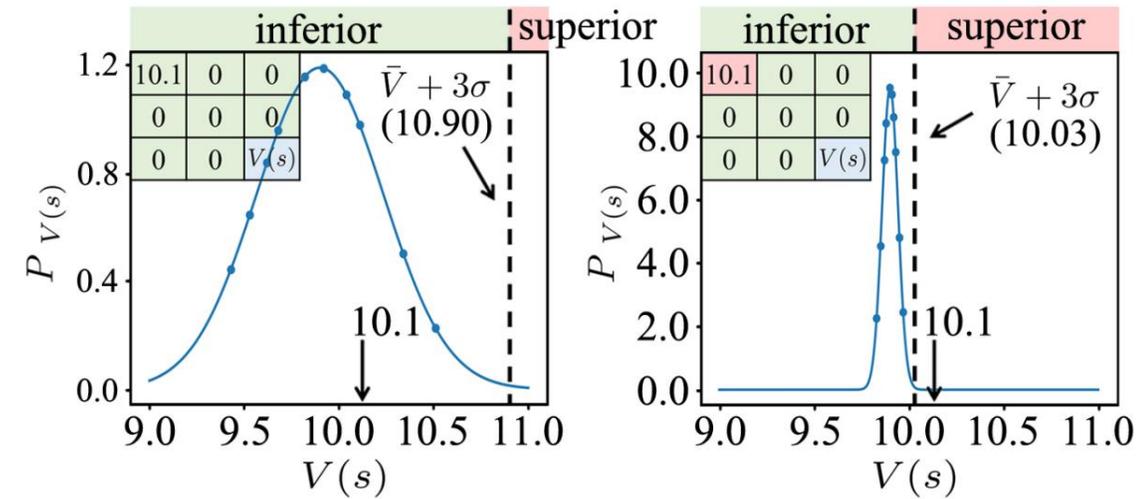
Methodology



optimality-stability paradox



adaptive trade-off between optimality and stability
filter out superior actions with large estimation uncertainty

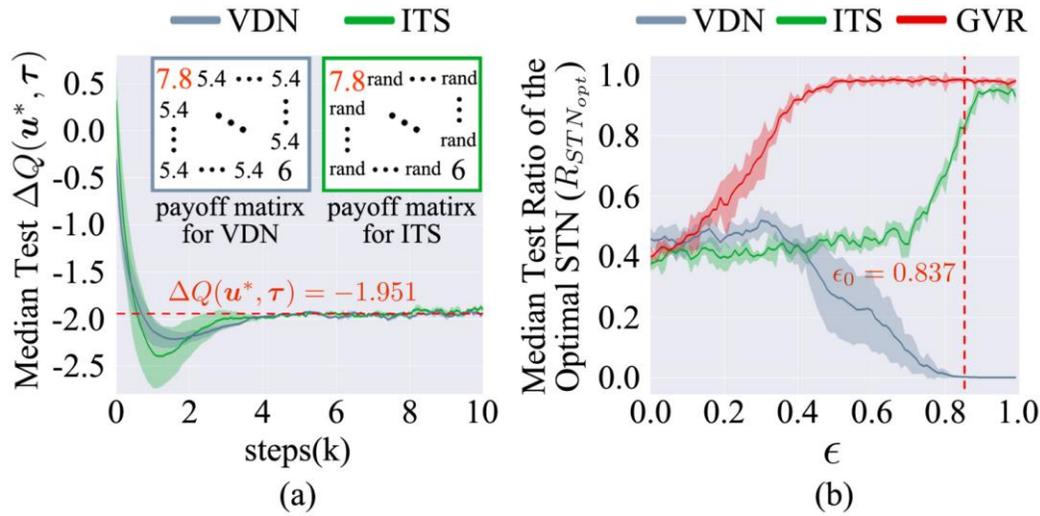


Experiments



Matrix Games

- Verification of optimal consistency under sufficient exploration

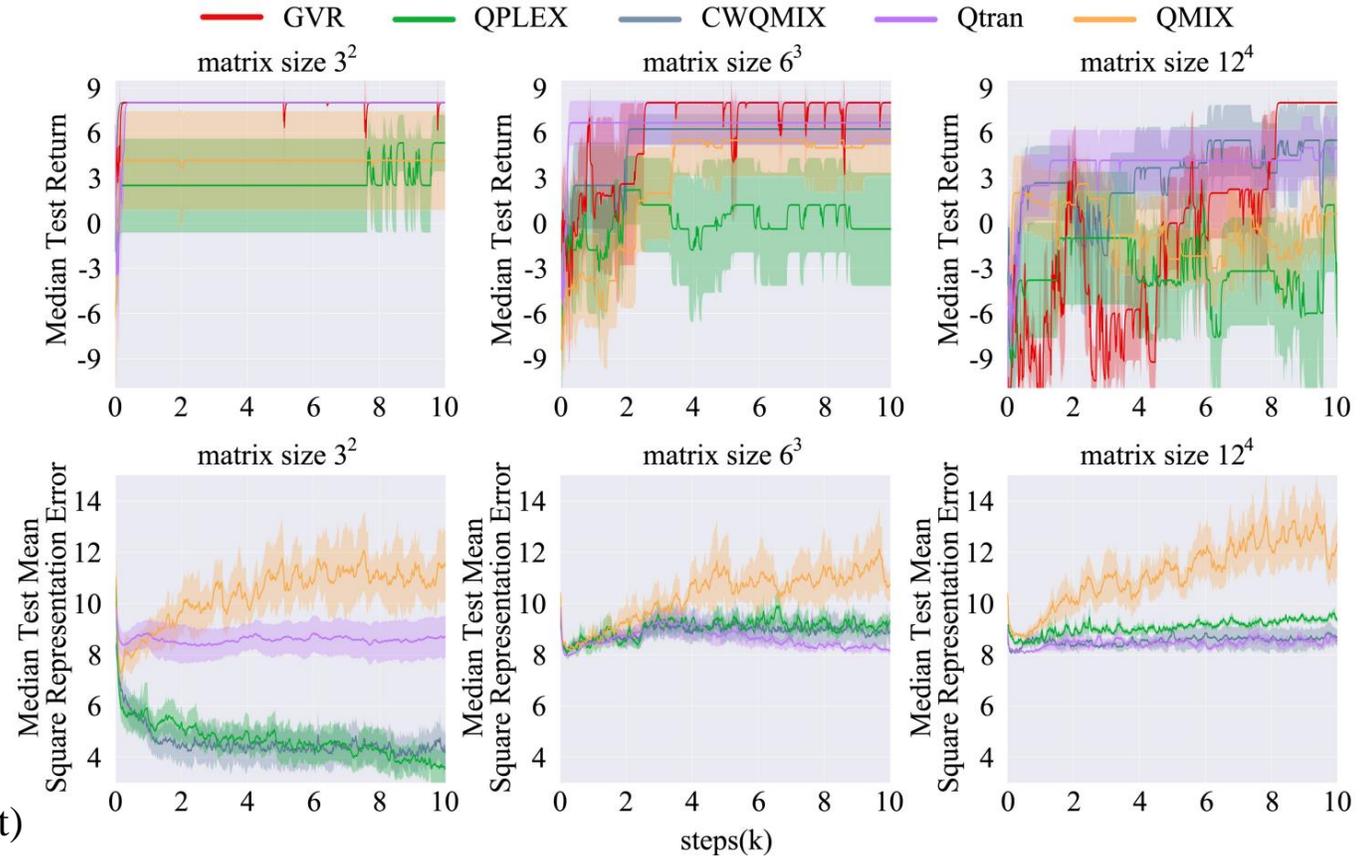


Experimental settings

$$Q(s, \mathbf{u}) = \begin{cases} 6(1 + e_Q) & \mathbf{u} = \{0, 0\} \\ 6 & \mathbf{u} = \{m, m\} \\ \text{random}(-20, 6) & \text{others} \end{cases}$$

$\alpha = 0.1, e_Q = 0.3$ (left) $\alpha = 0.1, \epsilon = 0.2, e_Q = 1/3$ (right)

- Comparison with baselines

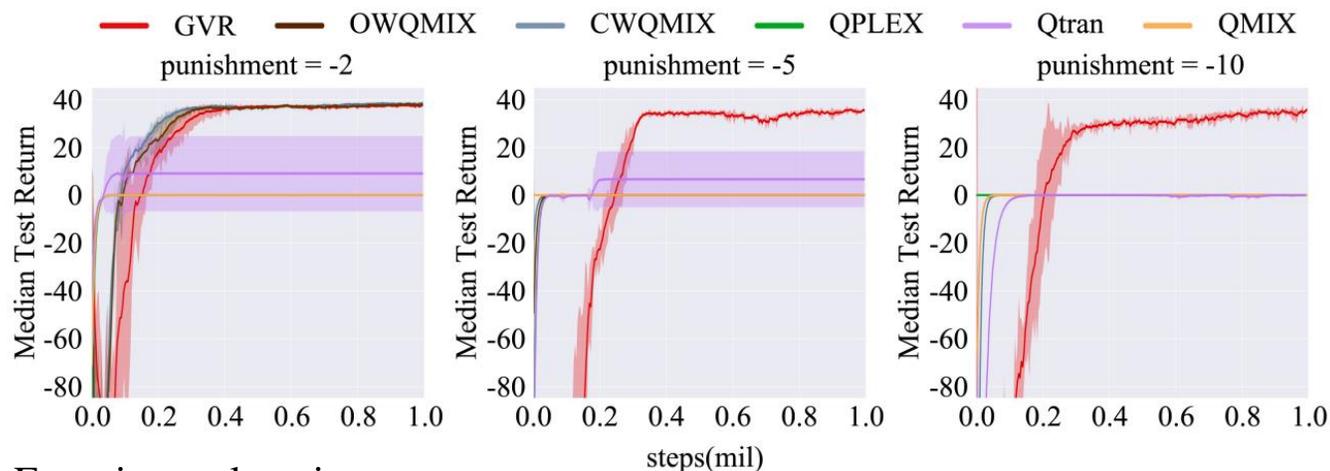


Experiments



Predator-Prey

- Preys are assigned with random policies
- suppose k agents capture the same prey at the same time-step
 - reward=0 if $k=0$
 - reward<0 (=punishment) if $k=1$
 - reward>0 if $k>1$



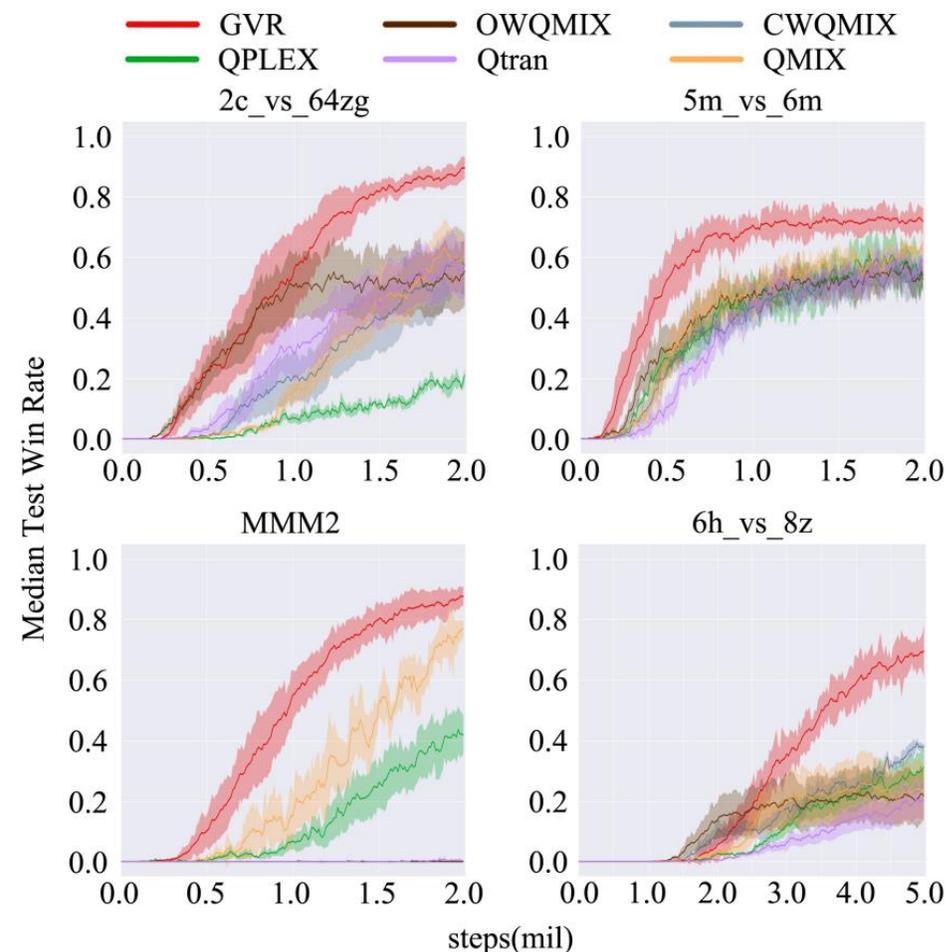
Experimental settings:

RNN shared by all agents as the utility functions

MVD for GVR, $\alpha=0.2$

exploration rate damps from 1 to 0.05 during 1m steps (for 6h_vs_8z) or 50k steps (for the others)

Starcraft Multi-agent Challenge (SMAC)



1. Relative over generalization arises from different representations under different action distributions
 - There is no difference between LVD and MVD from the perspective of value representation.
 - When ϵ grows to 1, the action distributions under different greedy actions converge to the uniform distribution, where the representations also converge.

2. Solving the RO is equivalent to ensure the optimal consistency
 - A sufficient condition to ensure the optimal consistency is: compared to the greedy action, the joint Q value of an inferior action is smaller; the joint Q value of a superior action is larger.
 - Inferior target shaping ensures the optimal node is an STN, but removing the non-optimal STN rely on further raising the probability of superior actions.

3. The paradox in optimality and stability
 - A trade-off between the optimality and stability can be achieve by moderate tolerance on sub-optimality.



Thanks !