



Multi-Agent RL

Zhi Wang

Nanjing University, China

2025-04-27



Contents



- Background on RL and MARL
 - Contrastive Role Representations for MARL
 - Interpretable MARL via Mixing Recurrent Soft Decision Trees
-

Reinforcement Learning



- Supervised Learning
 - (input, label)
- Unsupervised Learning
 - (input)
- **Reinforcement Learning**
 - sequential decision-making
- Computer Vision
 - Input: image pixels
- Natural Language Processing
 - Input: sentences
- **Reinforcement Learning**
 - Input: states

The Era of RL



- Video games: Human-level control through DRL, Nature 2015 (视频游戏)
- AlphaGo, Nature 2016; AlphaGo Zero, Nature 2017 (围棋)
- AlphaStar in StarCraft II, Nature 2019 (星际争霸II)
- DRL for legged robots, Science Robotics 2019 (机器人学习)
- Superhuman AI for multiplayer poker, Science 2019 (德州扑克, 多人非完全信息博弈)
- Discovering faster matrix multiplication algorithms, Nature 2022 (矩阵相乘算法发现, 基础数学)
- Magnetic control of tokamak plasmas, Nature 2022 (可控核聚变控制)
- Outracing champion Gran Turismo drivers, Nature 2022 (赛车模拟控制)
- Safety validation of autonomous vehicles, Nature 2023 (无人驾驶安全验证)
- Faster sorting algorithms discovering, Nature 2023 (排序算法发现, 基础信息科学)
- Champion-level drone racing, Nature 2023 (无人机竞速)
- Mastering diverse control tasks through world models, Nature 2025
-



RL = Artificial General Intelligence (AGI)?

Yet?

The Dilemma of RL



Transformers

Attention is all you need

[A Vaswani, N Shazeer, N Parmar...](#) - Advances in neural ... , 2017 - proceedings.neurips.cc

... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent

... **We** implement this inside of scaled dot-product **attention** by masking out (setting to $-\infty$) ...

☆ Save ⏪ Cite [Cited by 176805](#) Related articles All 73 versions ☰

Vision Transformers

[PDF] An **image** is worth 16x16 words: Transformers for **image** recognition at scale

[A Dosovitskiy, L Beyer, A Kolesnikov...](#) - arXiv preprint arXiv ..., 2020 - arxiv.org

... directly to **images**, with the fewest possible modifications. To do so, we split an **image** into patches ... only to small-resolution **images**, while we handle medium-resolution **images** as well. ...

☆ Save ⏪ Cite [Cited by 60296](#) Related articles All 21 versions ☰

Decision Transformers

Decision transformer: Reinforcement learning via sequence modeling

[L Chen, K Lu, A Rajeswaran, K Lee...](#) - Advances in neural ... , 2021 - proceedings.neurips.cc

... of the **Transformer** architecture, and associated advances in language modeling such as GPT-x and BERT. In particular, we present **Decision Transformer**, ... , **Decision Transformer** simply ...

☆ Save ⏪ Cite [Cited by 1942](#) Related articles All 13 versions ☰



ChatGPT (Generative Pre-Training)

Next-token prediction

Enter text:

The dog eats the apples.

**Transformer
Architecture**

The dog eats the apples .
464 3290 25365 262 22514 13

464 3290 25365 262 22514 13

**Self-supervised learning
Algorithm**

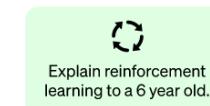
RL in ChatGPT



Step 1

Collect demonstration data and train a supervised policy.

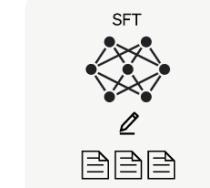
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



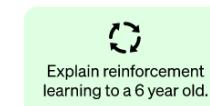
This data is used to fine-tune GPT-3.5 with supervised learning.



Step 2

Collect comparison data and train a reward model.

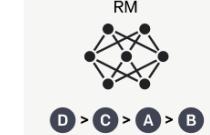
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



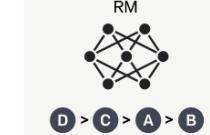
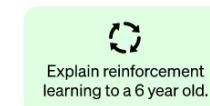
This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.



r_k

The PPO model is initialized from the supervised policy.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

RL: Fine-tuning in Step 3, playing an **auxiliary** role

The Dilemma of RL



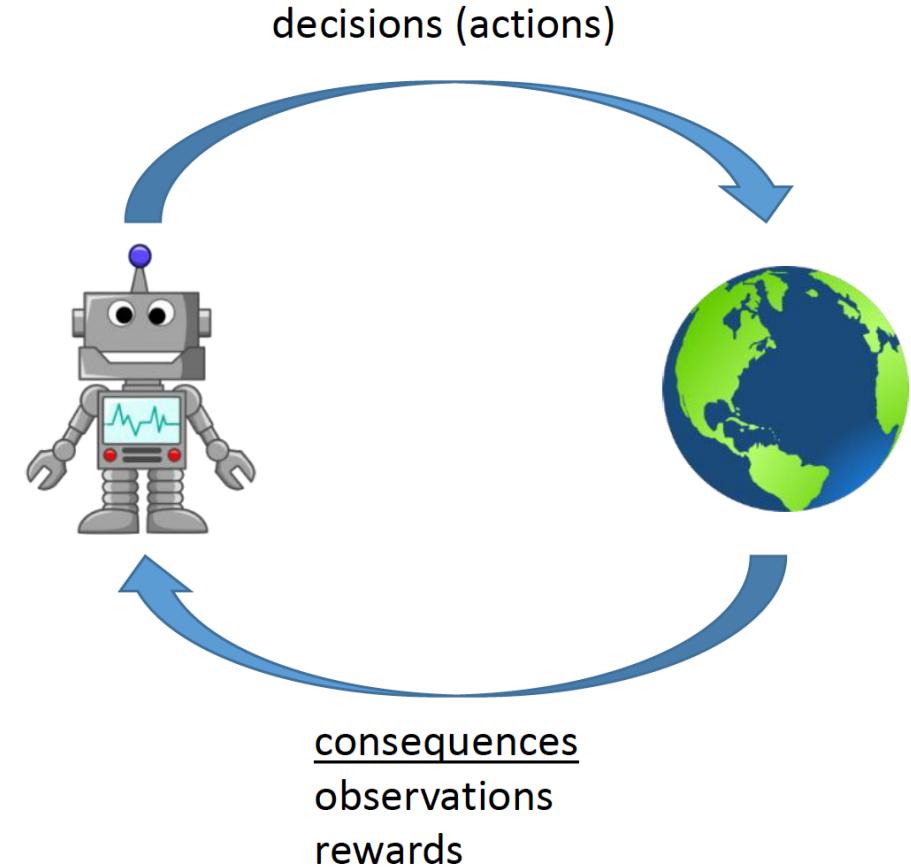
- Computer Vision
 - Input: **image pixels**
- Natural Language Processing
 - Input: **sentences**
- Reinforcement Learning
 - Input: **states, (states, actions)**

**Semantics
not aligned**

The Dilemma of RL

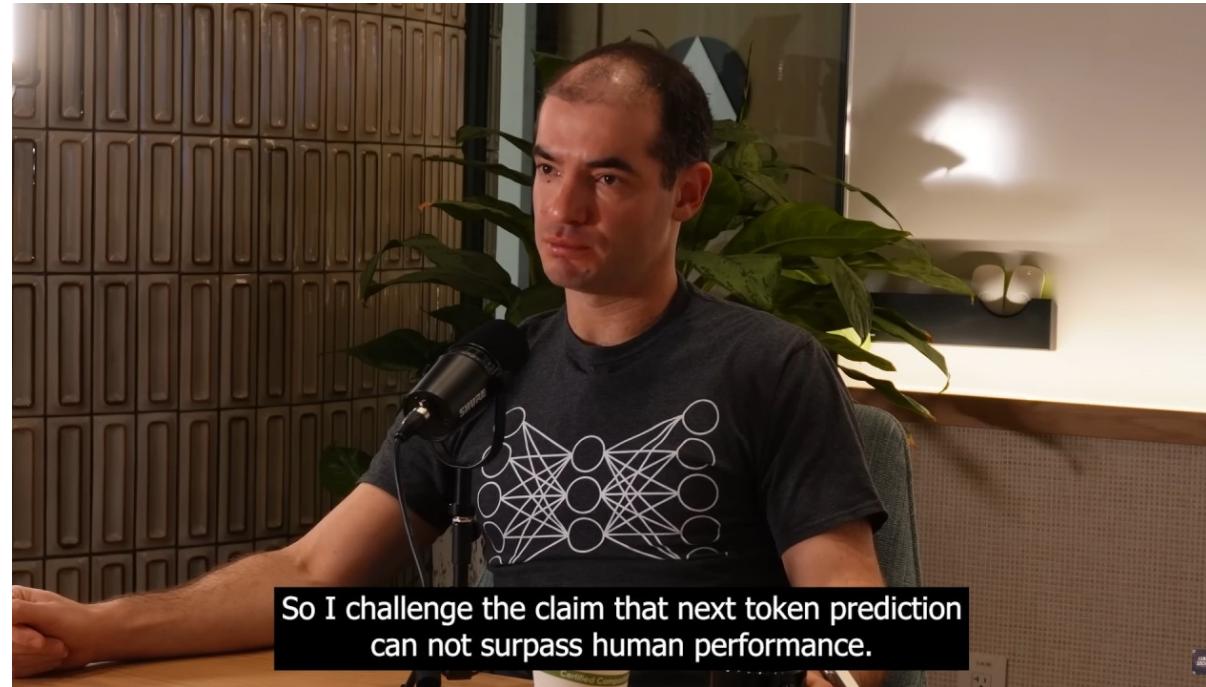


Data
From online interactions



Supervised learning

Maybe imitating the intelligence within existing data?





Supervised learning

Maybe imitating the intelligence within existing data?

Reinforcement learning

Can surpass the intelligence within existing data definitely

LLM: From Pre-Training to Post-Training

Pre-training will end

-- by Ilya Sutskever

@NeurIPS 2025

Pre-training as we know it will end

Compute is growing:

- Better hardware
- Better algorithms
- Larger clusters

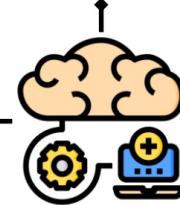
Data is not growing:

- We have but one internet
- **The fossil fuel of AI**

LLM: From Pre-Training to Post-Training

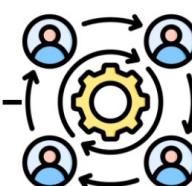
Reasoning, inference

- Supervised Fine-Tuning
- Reinforcement Fine-Tuning



Fine-Tuning

- Self-Refine for Reasoning
- Reinforcement Learning for Reasoning



Alignment

- Reinforcement Learning with Human Feedback
- Direct Preference Optimization
- Group Relative Policy Optimization



Reasoning

What comes next?

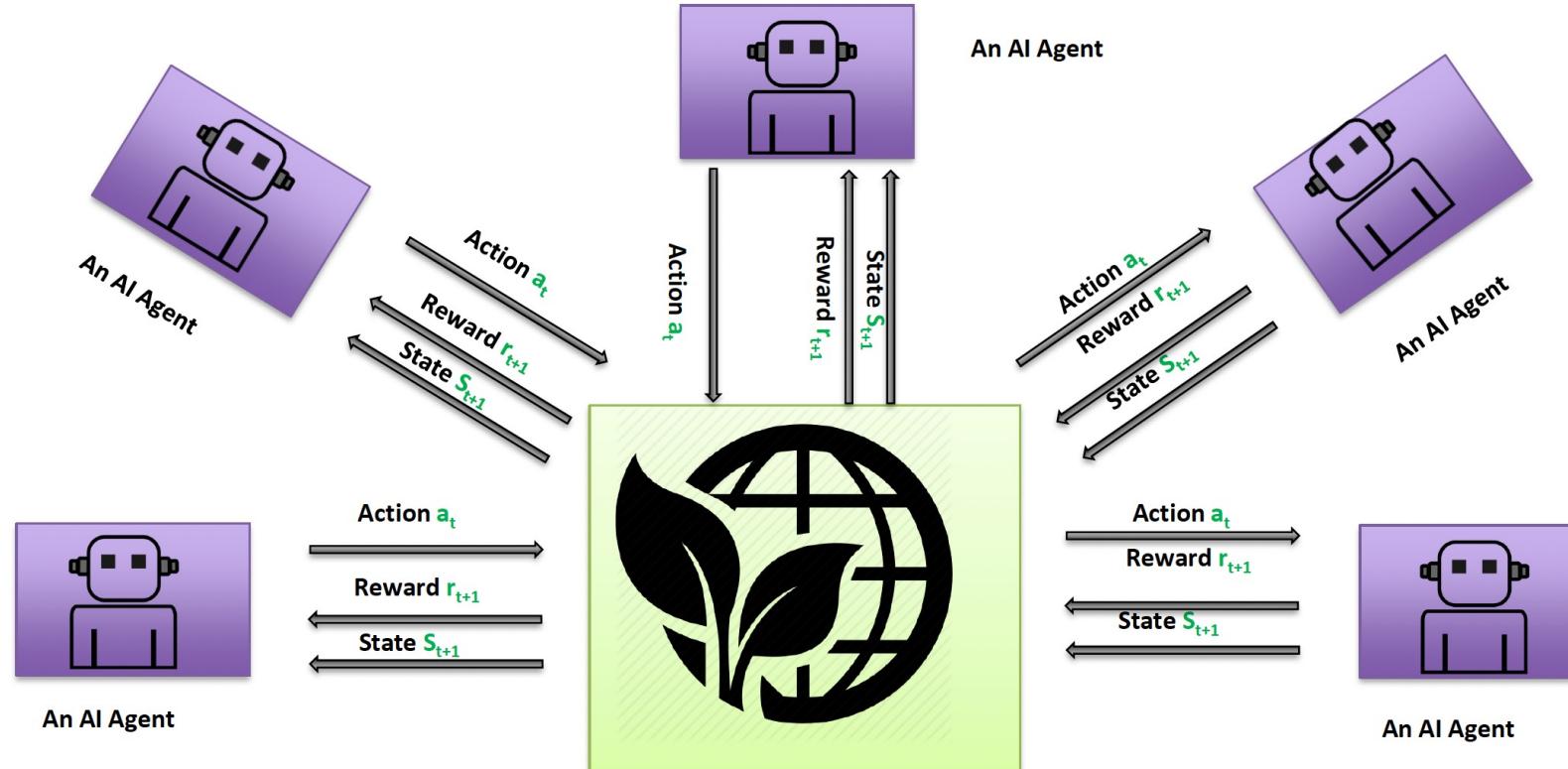
- “Agents”??
- “Synthetic data”
- Inference time compute ~ O1



Efficiency

- Model Compression
- Parameter-Efficient Fine-Tuning
- Knowledge Distillation

Multi-Agent RL



Environment

A set of autonomous agents that share a common environment

Difficulty in MARL



- **MARL is fundamentally difficult**
 - since agents not only interact with the environment but also with each other

- **If use single-agent Q-learning by considering other agents as a part of the environment**
 - Such a setting breaks the theoretical convergence guarantees and makes the learning unstable
 - i.e., the changes in strategy of one agent would affect the strategies of other agents and vice versa

Types of Multi-agent Systems



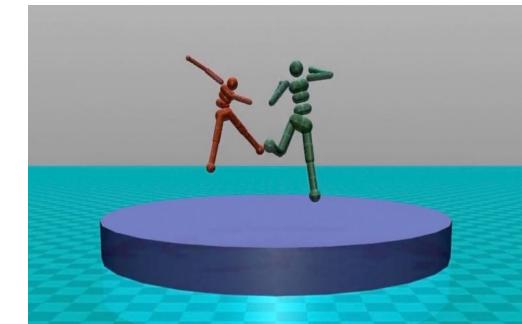
■ Cooperative

- Maximizing a shared team reward
- Coordination problems



■ Competitive

- Self-interested: maximizing an individual opposite reward
- Zero-sum games
- Minimax equilibria



■ Mixed

- Self-interested with different individual rewards (not opposite)
- General-sum games

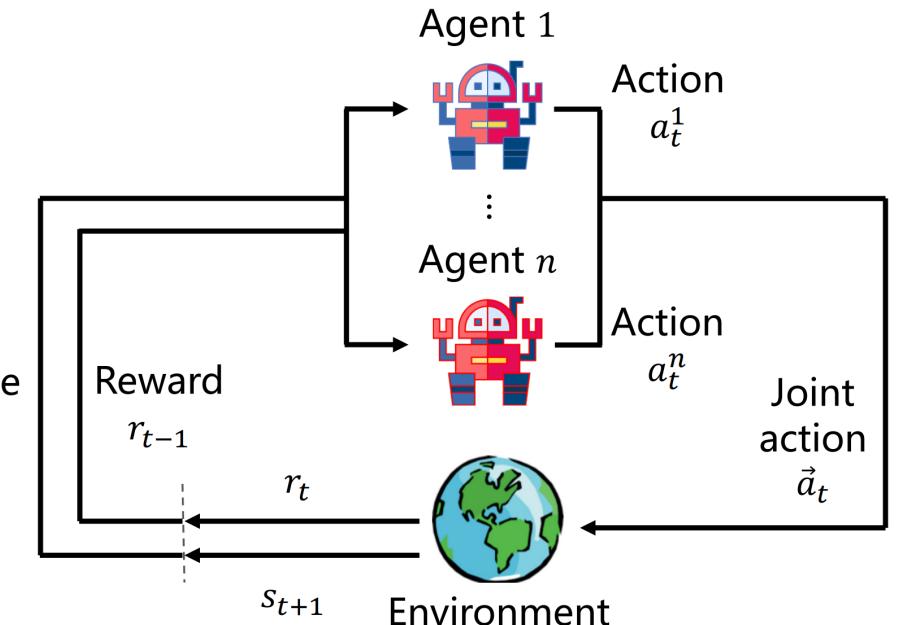


Cooperative Multi-Agent MDP



■ Assume agents can observe the global state

- Agent: $i \in I = \{1, 2, \dots, N\}$
- State: $s \in S$
- Action: $a_i \in A$, joint action $a = < a_1, \dots, a_n > \in A^N$
- Transition function: $P(s'|s, a)$
- Reward: $r(s, a)$
- Agent i 's policy $\pi_i(s): S \rightarrow A$
- Objective: finding a joint policy $\pi = < \pi_1, \dots, \pi_n >$ to
maximize expected return: $R = \sum_{t=0}^{\infty} \gamma^t r_t$
- Value function: $Q^\pi(s, a) = \mathbb{E}[R | s_0 = s, a_0 = a, \pi]$
- With optimal $Q^*(s, a)$, optimal $\pi^*(s) = \text{argmax}_a Q^*(s, a)$



Decentralized Partially Observable MDP



■ Agent can not observe the global states

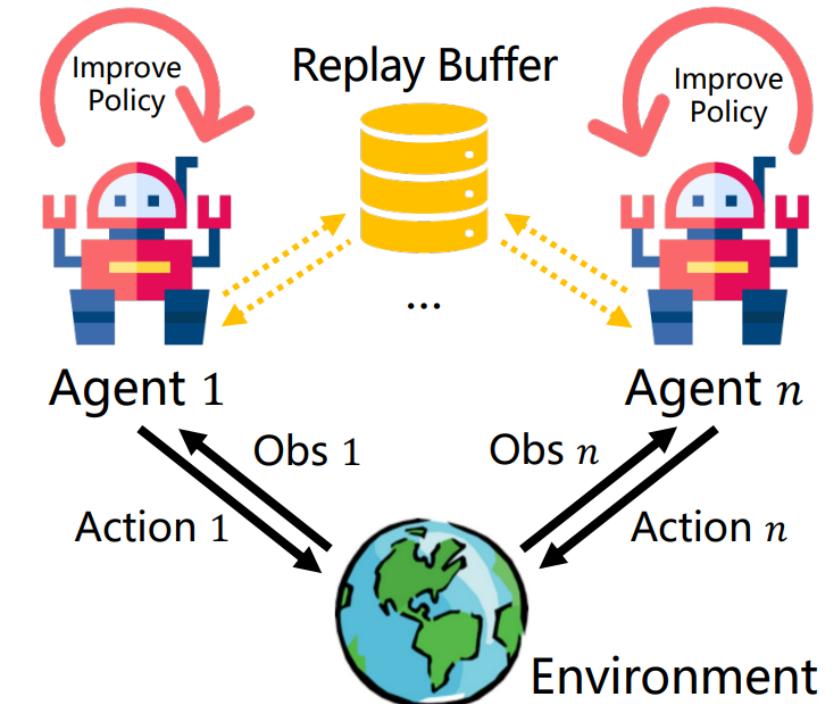
- Observation: $o_i \in \Omega$
- Observation function: $o_i \in \Omega \sim O(s, i)$

■ Decentralized policy for agent i :

- $\pi_i(\tau_i): T \rightarrow A$
- Action-observation history: $\tau_i \in T = (\Omega \times A)^*$

■ Communication and sensory constraints

- Decentralized execution



Challenges of Cooperative MARL

■ Scalability

- Curse of dimensionality

■ Multi-Agent Credit Assignment

- Each agent's contribution to the team

■ Learning Efficiency

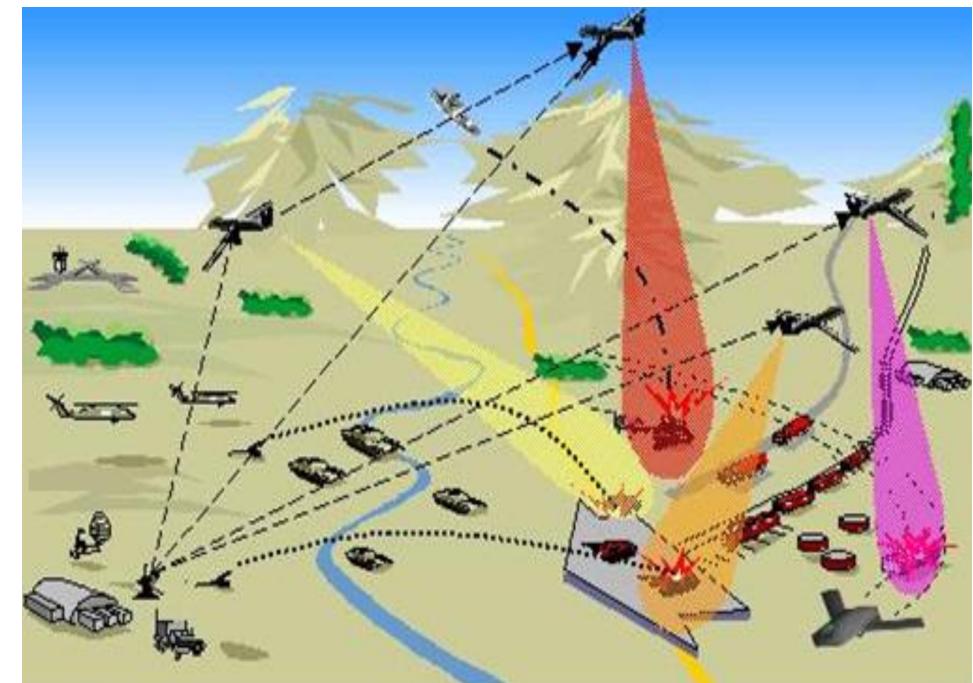
- Requiring extensive interactions

■ Limited Observability

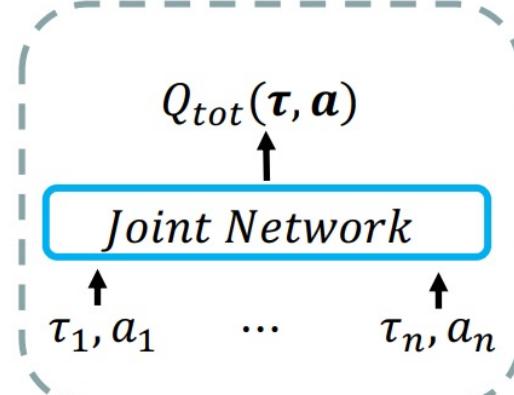
- Sensory constraints

■ Exploration

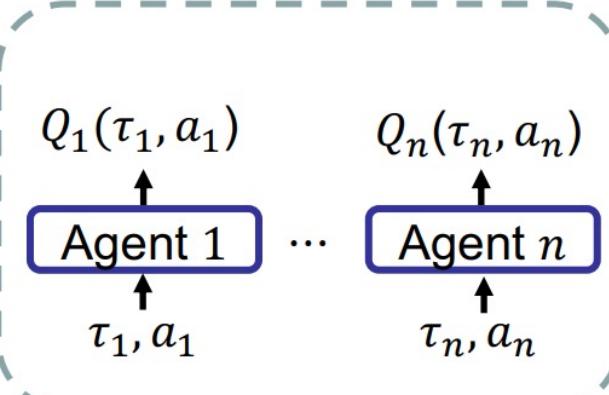
- An exponential joint policy space



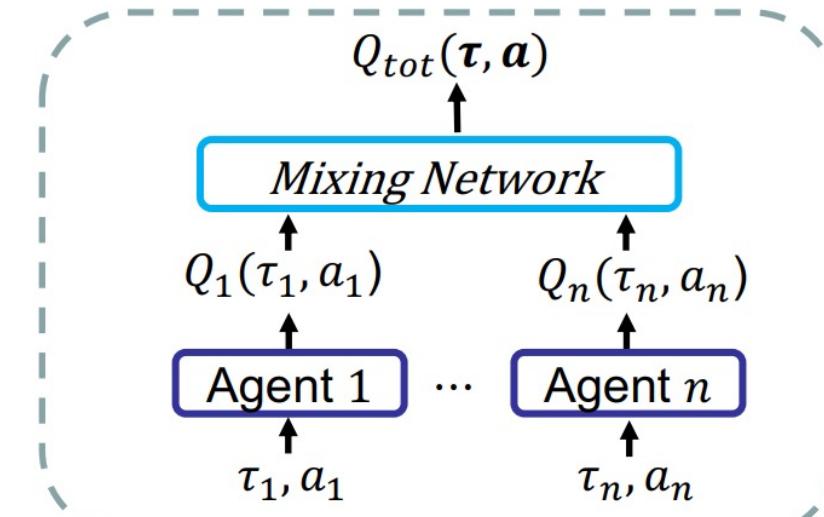
MARL Paradigms



Joint Learning



Independent Learning



Factored Learning

- Optimality

- Scalability

- Scalability

- Suboptimality
- Non-stationarity
- Credit assignment

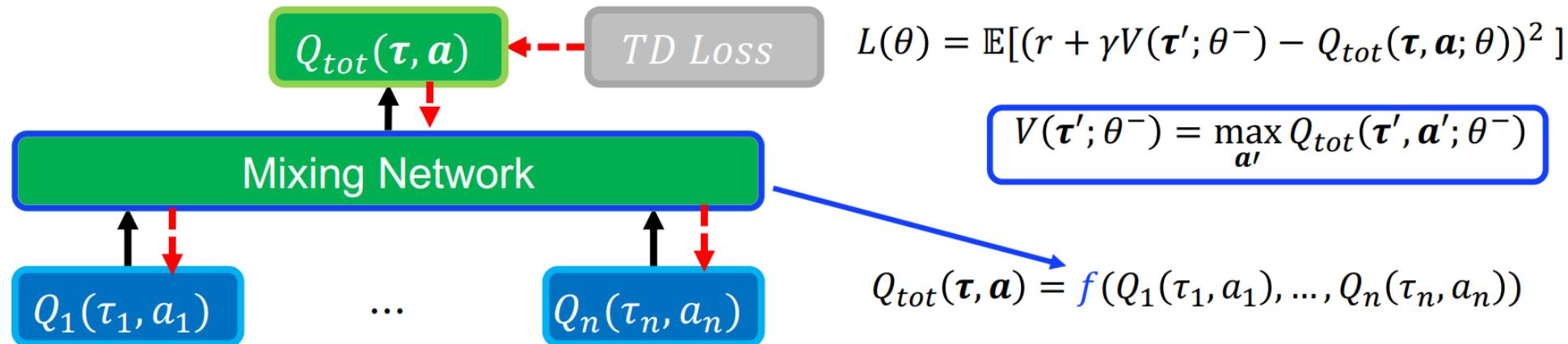
- Scalable centralized training
- Decentralized execution
- Implicit credit assignment
- Optimality & Convergence

- Centralized training

Factored Value Functions for MARL



- Scalable centralized training with decentralized execution



- Individual-Global Maximization (IGM) Constraint
 - $\underset{\mathbf{a}}{\operatorname{argmax}} Q_{tot}(\tau, \mathbf{a}) = (\underset{a_1}{\operatorname{argmax}} Q_1(\tau_1, a_1), \dots, \underset{a_n}{\operatorname{argmax}} Q_n(\tau_n, a_n))$
 - Consistent greedy action selection between joint and individuals

Factored Value Functions for MARL

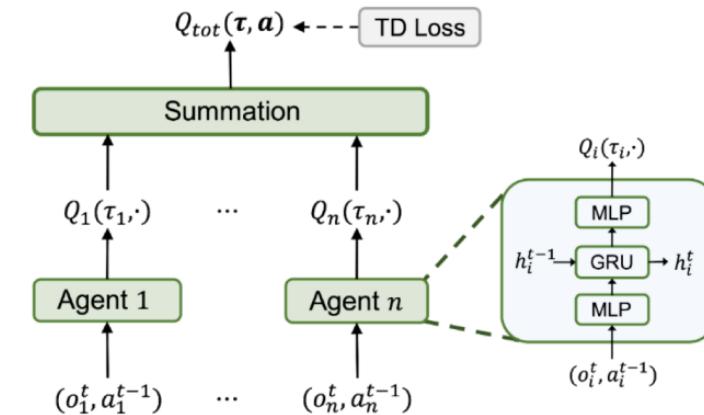


- Linear Mixing: $Q_{tot}(\tau, \mathbf{a}) = \sum_i Q_i(\tau_i, a_i)$ [Sunehag et. al., 2017]

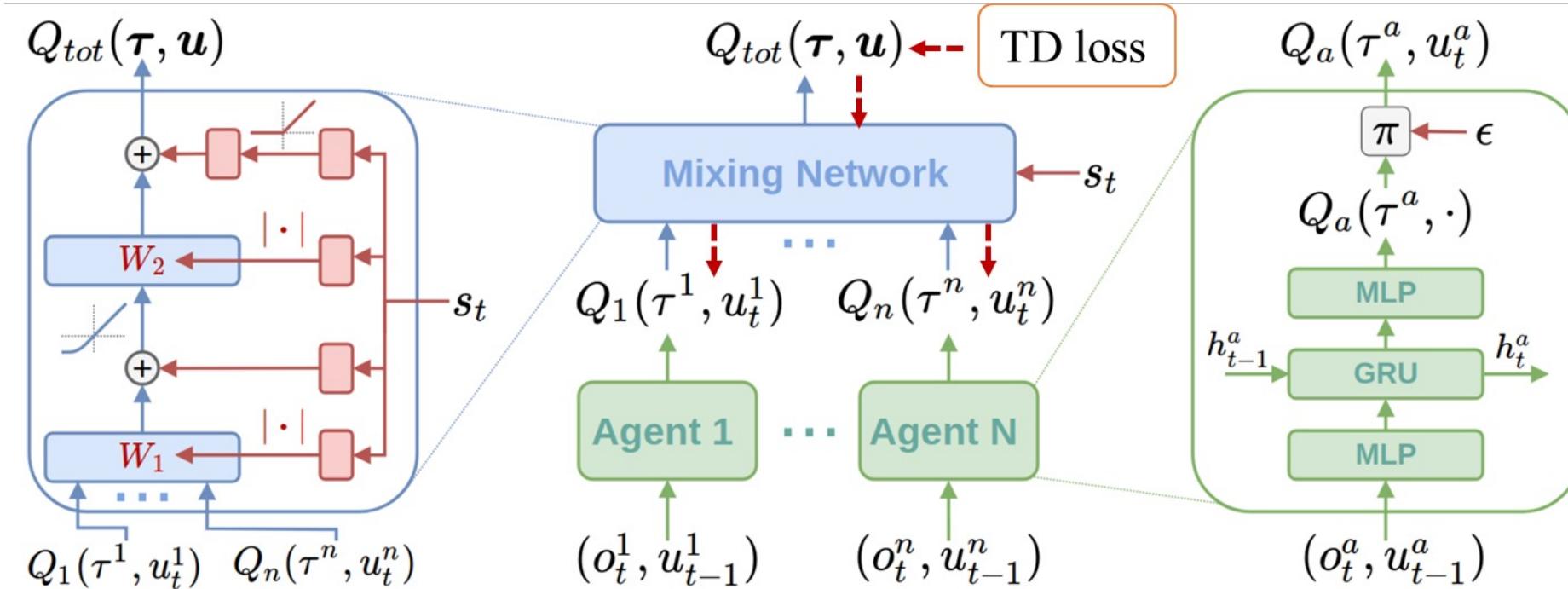
- Satisfying IGM Constraint

$$\underset{\mathbf{a}}{\operatorname{argmax}} Q_{tot}(\tau, \mathbf{a}) = \begin{pmatrix} \operatorname{argmax}_{a_1} Q_1(\tau_1, a_1) \\ \dots \\ \operatorname{argmax}_{a_n} Q_n(\tau_n, a_n) \end{pmatrix}$$

- No parameters in the mixing network
- No specific reward for each agent
- Implicit credit assignment through gradient backpropagation



Centralized Training and Decentralized Execution



■ Individual-Global Maximization (IGM) Constraint

- $\text{argmax}_a Q_{tot}(\tau, a) = (\text{argmax}_{a_1} Q_1(\tau_1, a_1), \dots, \text{argmax}_{a_n} Q_n(\tau_n, a_n))$
- Consistent greedy action selection between joint and individuals

Environment

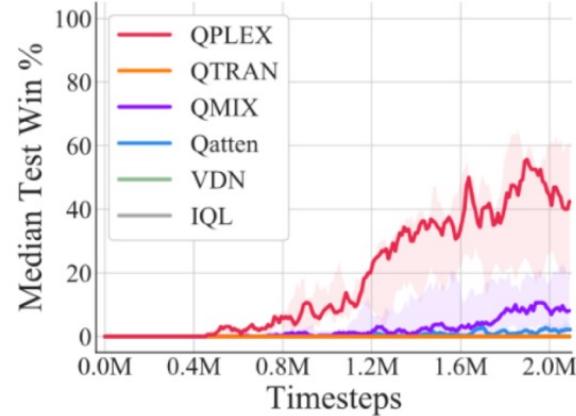


StarCraft Multi-Agent Challenge (SMAC)

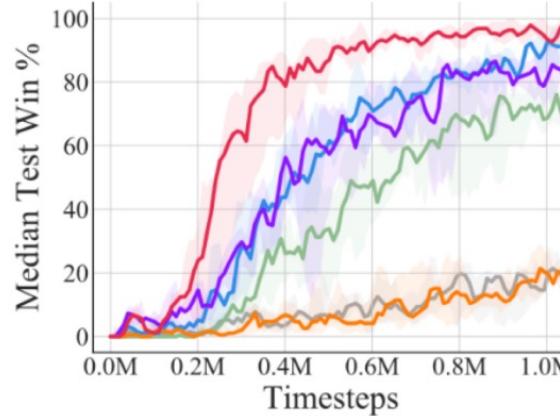


- GRF
- MA-Mujoco
- MOBA

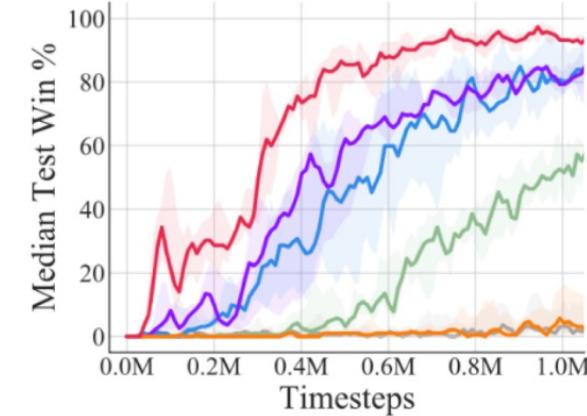
Online Learning Performance



(a) 5s10z



(b) 1c3s5z



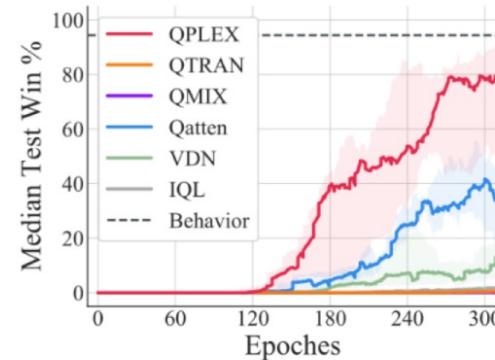
(c) 3s5z



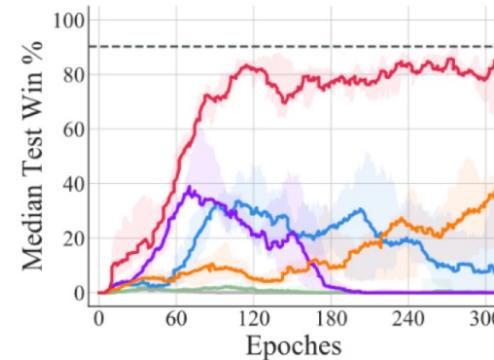
Data-Driven Offline MARL Learning



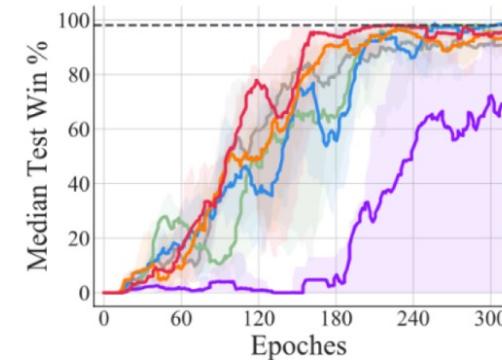
Data collected by a behavior policy learned by QMIX



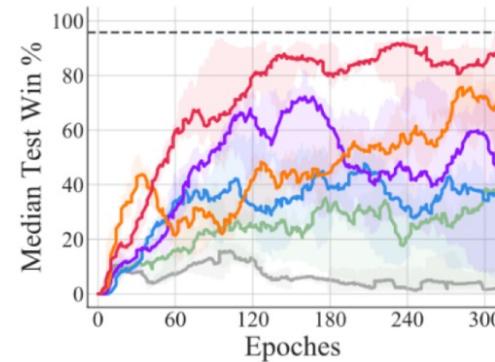
(a) 3s_vs_5z



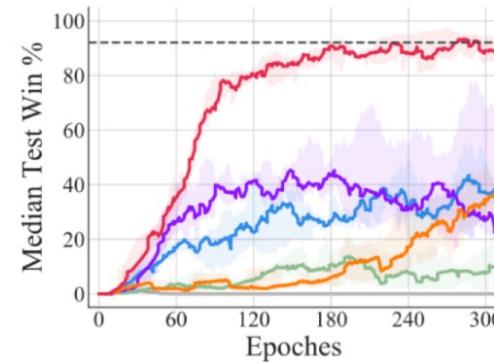
(b) 1c3s5z



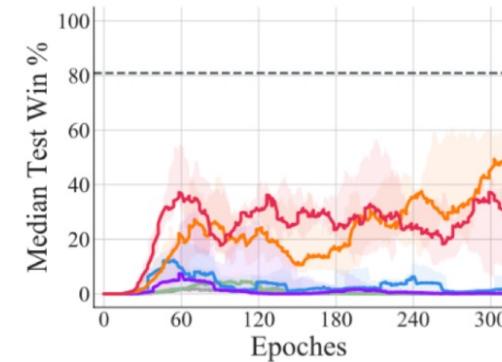
(c) 2s_vs_1sc



(d) 2s3z



(e) 3s5z



(f) 2c_vs_64zg

MARL with Factored Policies



- COMA: Counterfactual Multi-Agent Policy Gradients [Foerster et. al, 2017]
- MADDPG: Multi-agent actor-critic for mixed cooperative-competitive environments [Lowe et. al, 2017]
- MAPPO: Multi-Agent PPO [Yu et. al, 2021]
- HATRPO/HAPPO: Trust Region Policy Optimisation in Multi-Agent Reinforcement Learning [Kuba et. al, 2021]
- T-PPO: Towards Global Optimality in Cooperative MARL with Sequential Transformation [Ye et. al, 2021]

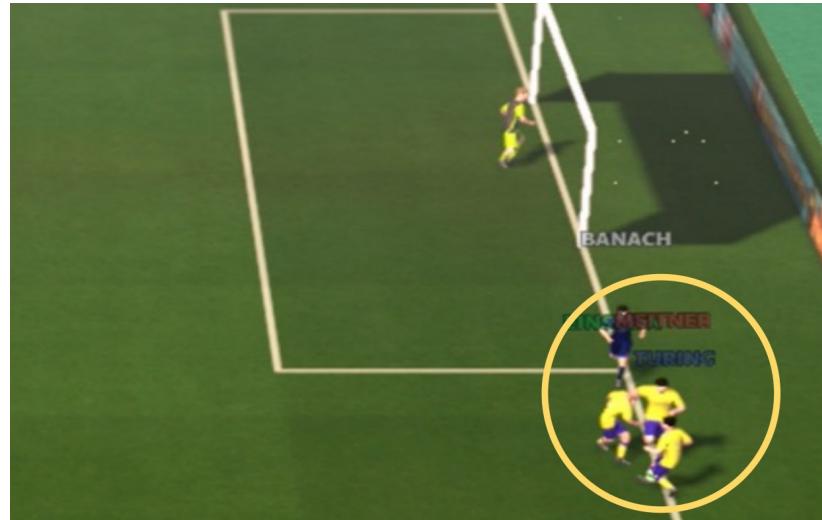
Contents



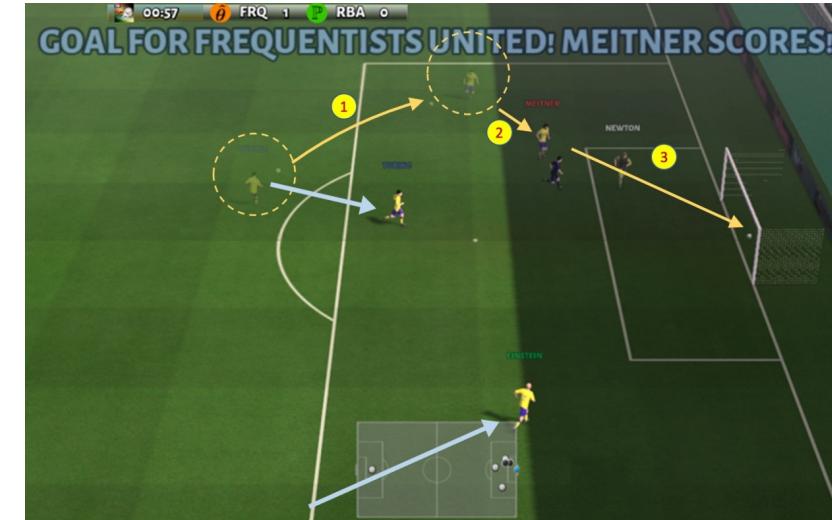
- Background on RL and MARL
 - Contrastive Role Representations for MARL
 - Interpretable MARL via Mixing Recurrent Soft Decision Trees
-

Limitation: Shared learning among agents

- Parameter sharing is critical for deep MARL methods
- However, agents tend to acquire **homogeneous** behaviors
- Dynamic sharing with diversity is essential for practical tasks



Similar behaviors (competing for ball)



Each agent has its responsibility to score

MARL with Representation Learning

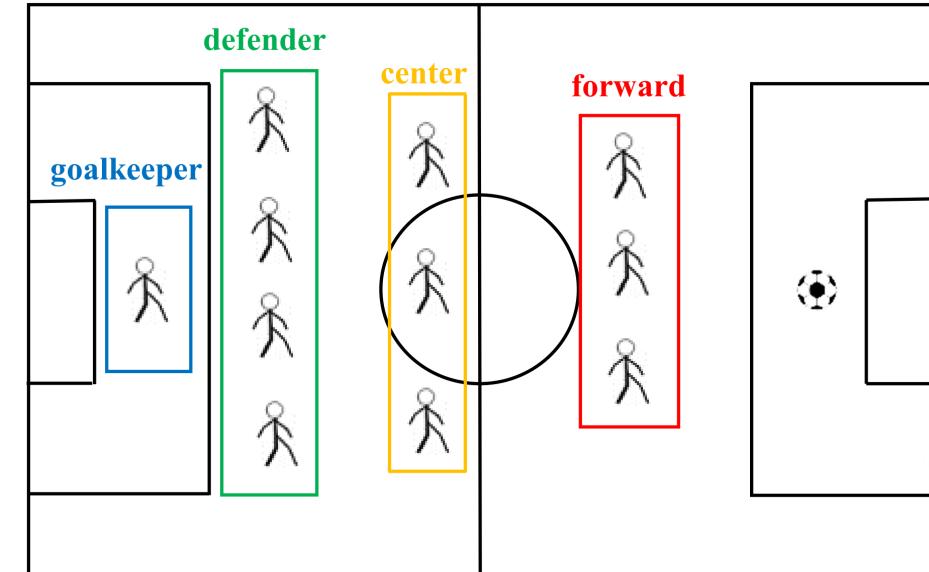


■ Agents with similar role have similar policies and share their learning

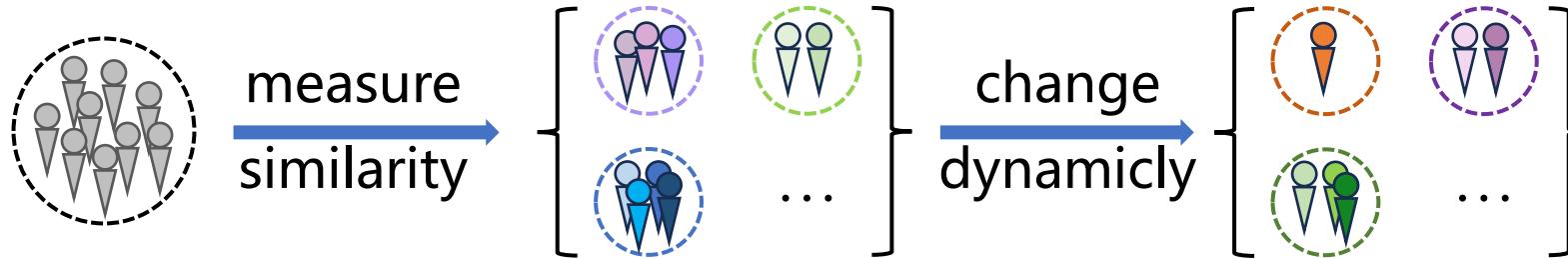
- Roles \leftrightarrow Subtasks \leftrightarrow Skills
- An example of subtask assignment in football: forward, center, defender, goalkeeper

■ Benefits of task decomposition and Role (Subtask or Skill) assignment

- Agent can change its roles in different situations \rightarrow Tackles agent homogenization
- Agent learn policy conditioned on their roles \rightarrow Facilitates efficient knowledge transfer



Challenge of Role Representation Learning



- ☒ **How to measure the similarity of the agents?**
- ☒ **How to define role representation?**
- ☒ **How to achieve the knowledge transfer?**
- ☒ **How to change the role dynamically?**

Our method: ACORM



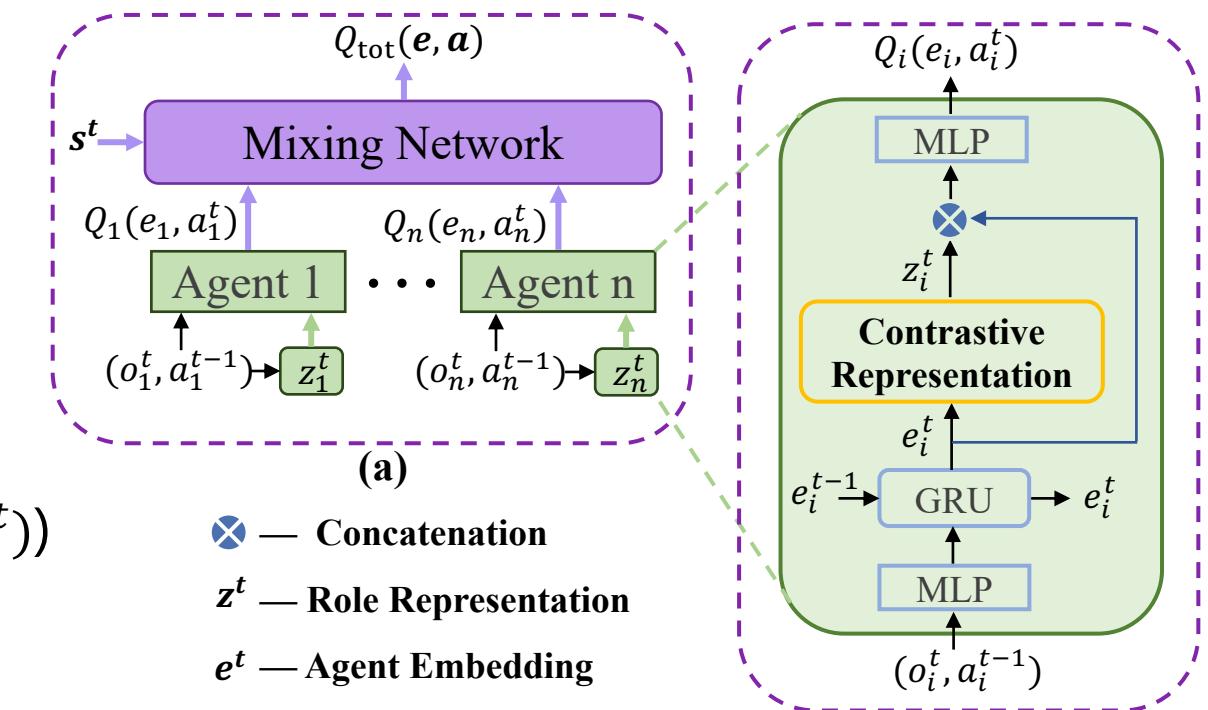
(i) Contrastive Role Representation

■ Learn agent embedding

- Extract complex agent behaviors from trajectory as $e_i^t = f_\Phi(o_i^t, a_i^{t-1}, e_i^{t-1})$

■ Learn role representation

- Reinforce role representation ($z^t \sim f_\theta(z^t | e^t)$) through contrastive learning



i) Contrastive Role Representation

■ Negative pairs generation

- Cluster the agent embedding
- the same cluster set as positive keys
- The different clusters set as negative

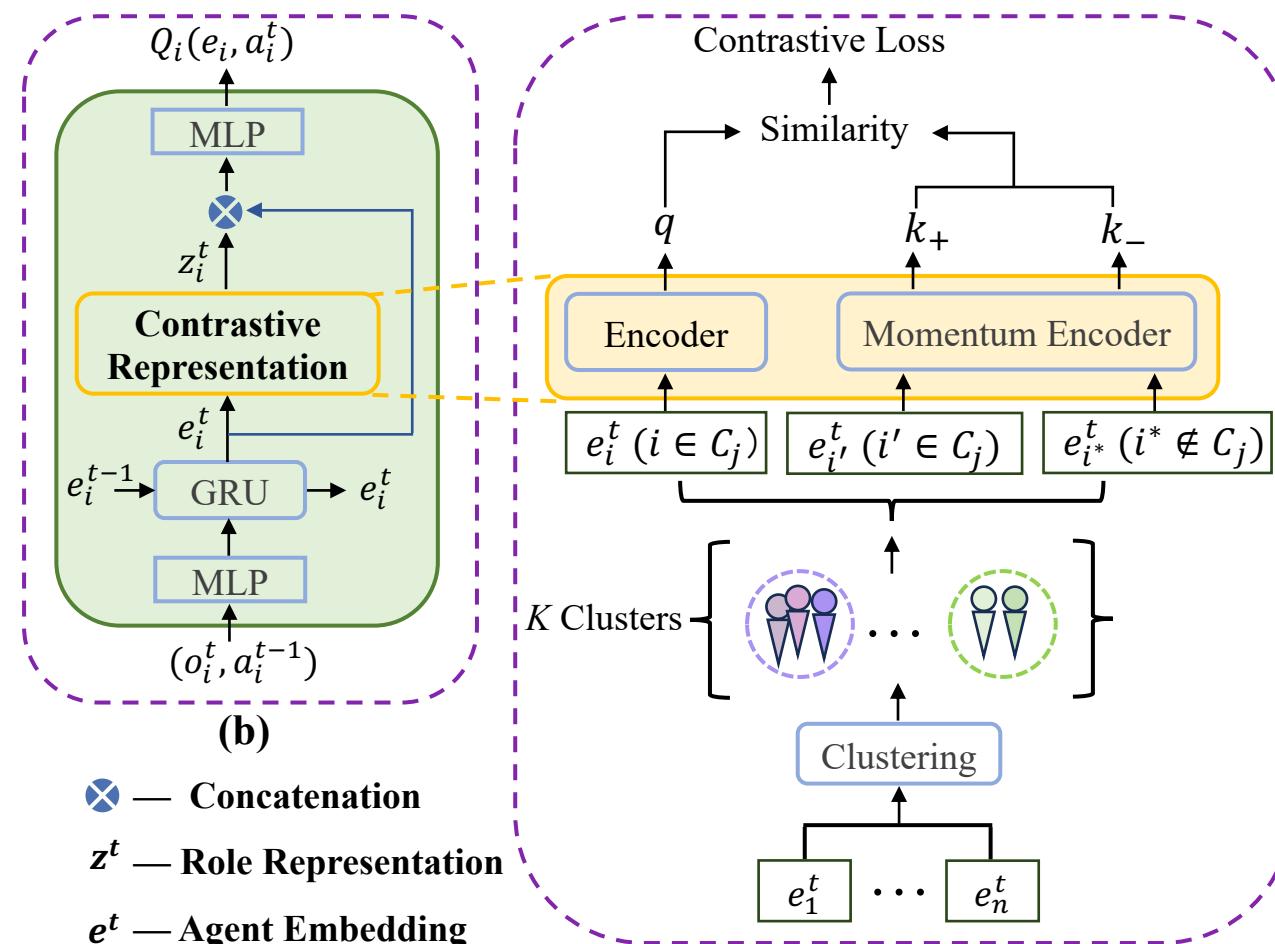
■ Calculate contrastive loss

- InfoNCE loss is rearranged as

$$\mathcal{L}_K = -\log \frac{\exp(q^\top W k_+)}{\exp(q^\top W k_+) + \exp(q^\top W k_-)}$$

■ Update momentum encoder

$$\theta_k \leftarrow \beta \theta_k + (1 - \beta) \theta_q$$



ii) Attention-Guided Role Coordination

■ Attention mechanism to enhance coordination

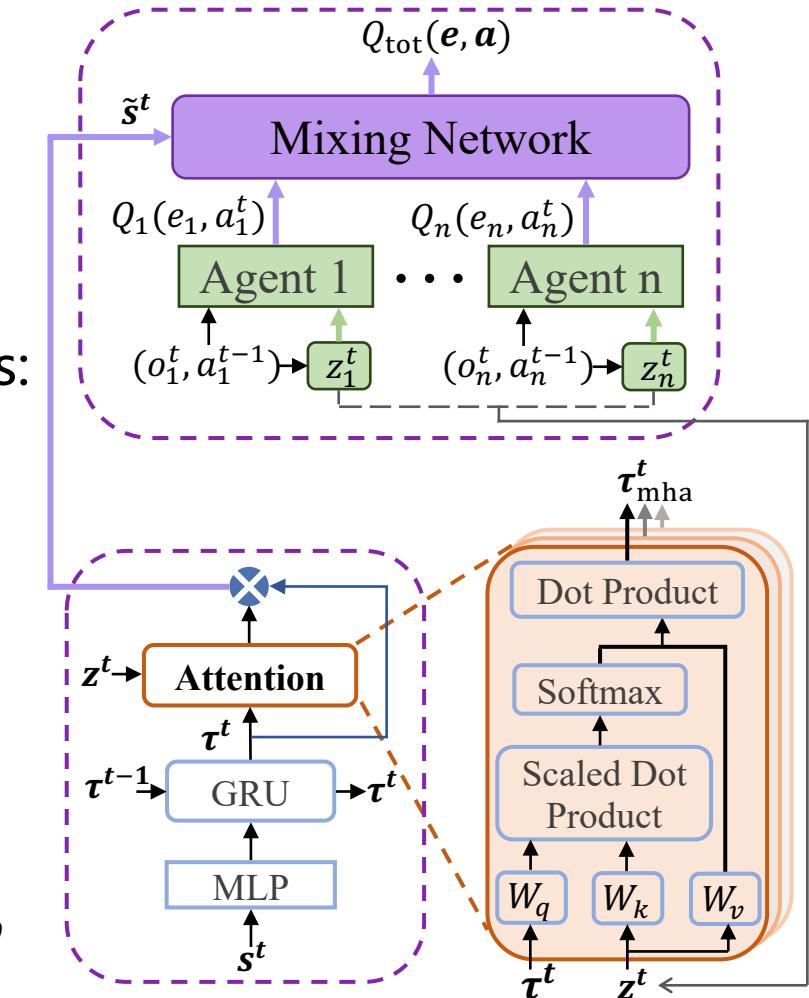
- Set state embedding as the query, role representation as the key and value
- calculate a weighted combination of role representations as:

$$\tau_{\text{atten}} = \sum_{i=1}^n \alpha_i v_i = \sum_{i=1}^n \alpha_i \cdot z_i W^V$$

- The attention weight α_i computes as:

$$\alpha_i = \frac{\exp\left(\frac{1}{\sqrt{d_k}} \cdot \tau W^Q \cdot (z_i W^K)^T\right)}{\sum_{j=1}^n \exp\left(\frac{1}{\sqrt{d_k}} \cdot \tau W^Q \cdot (z_j W^K)^T\right)}$$

- Obtain the aggregated output as: $\tau_{\text{mha}} = (\tau_{\text{atten}}^1, \dots, \tau_{\text{atten}}^H) W^O$



Attention-guided COntrastive Role Representations for MARL (ACORM)

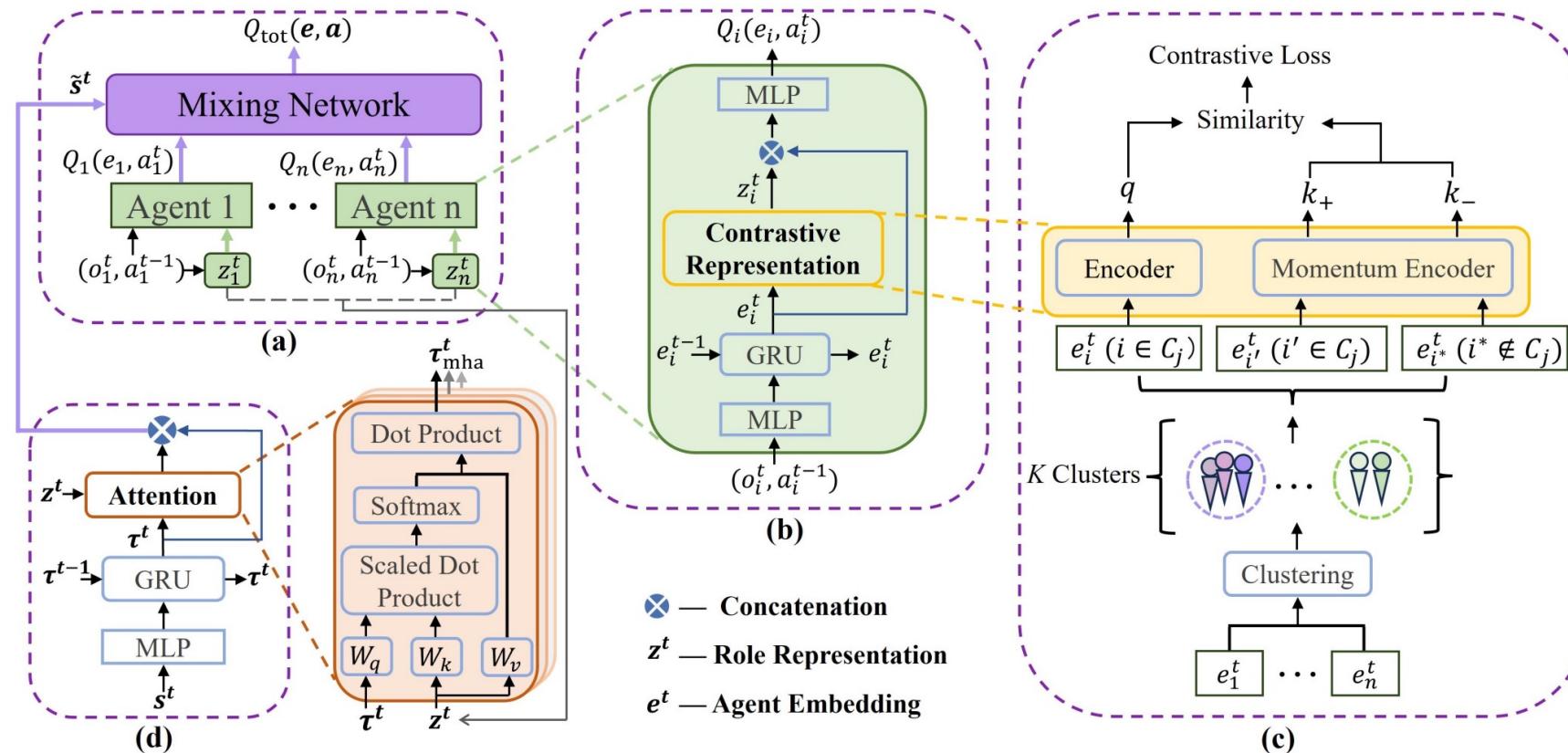
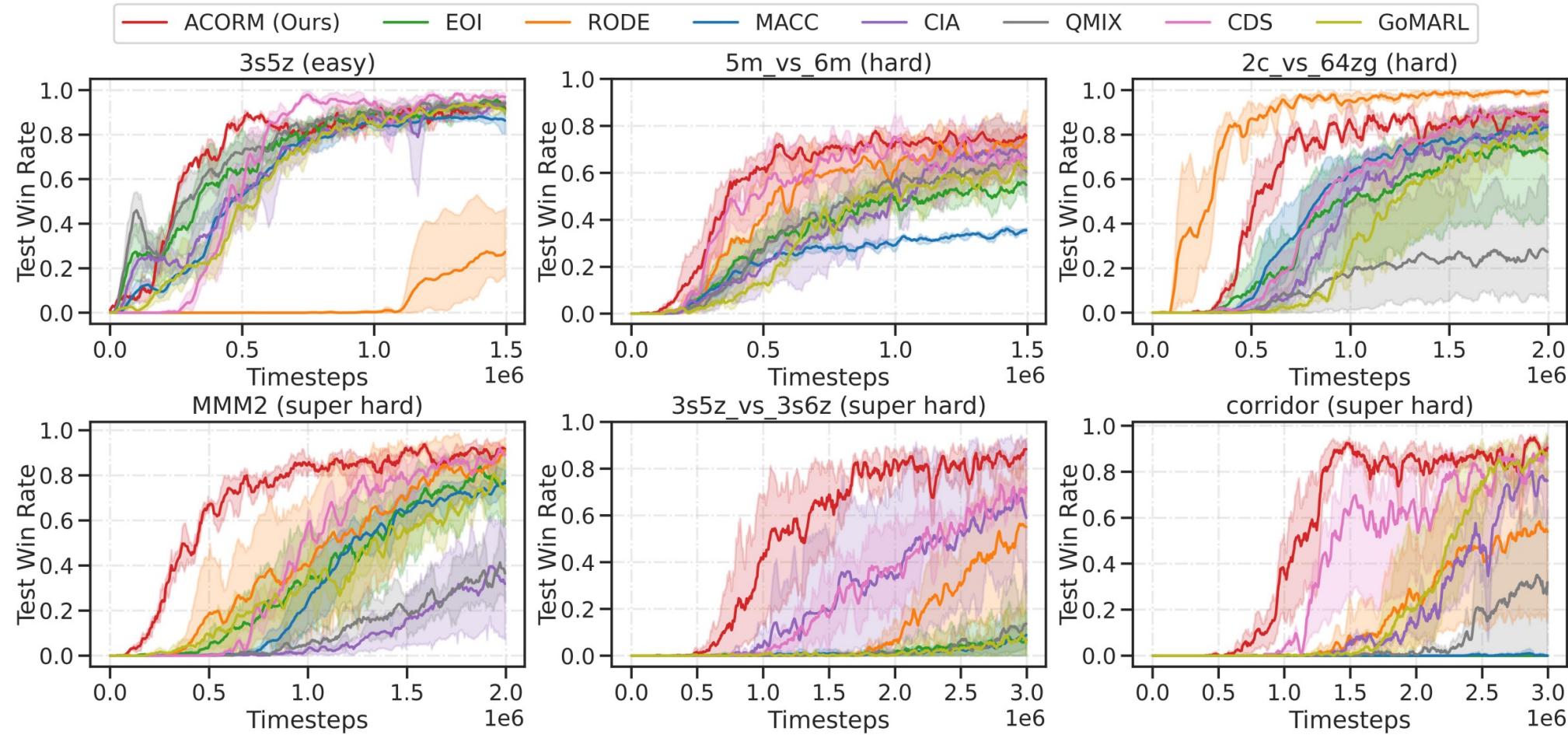
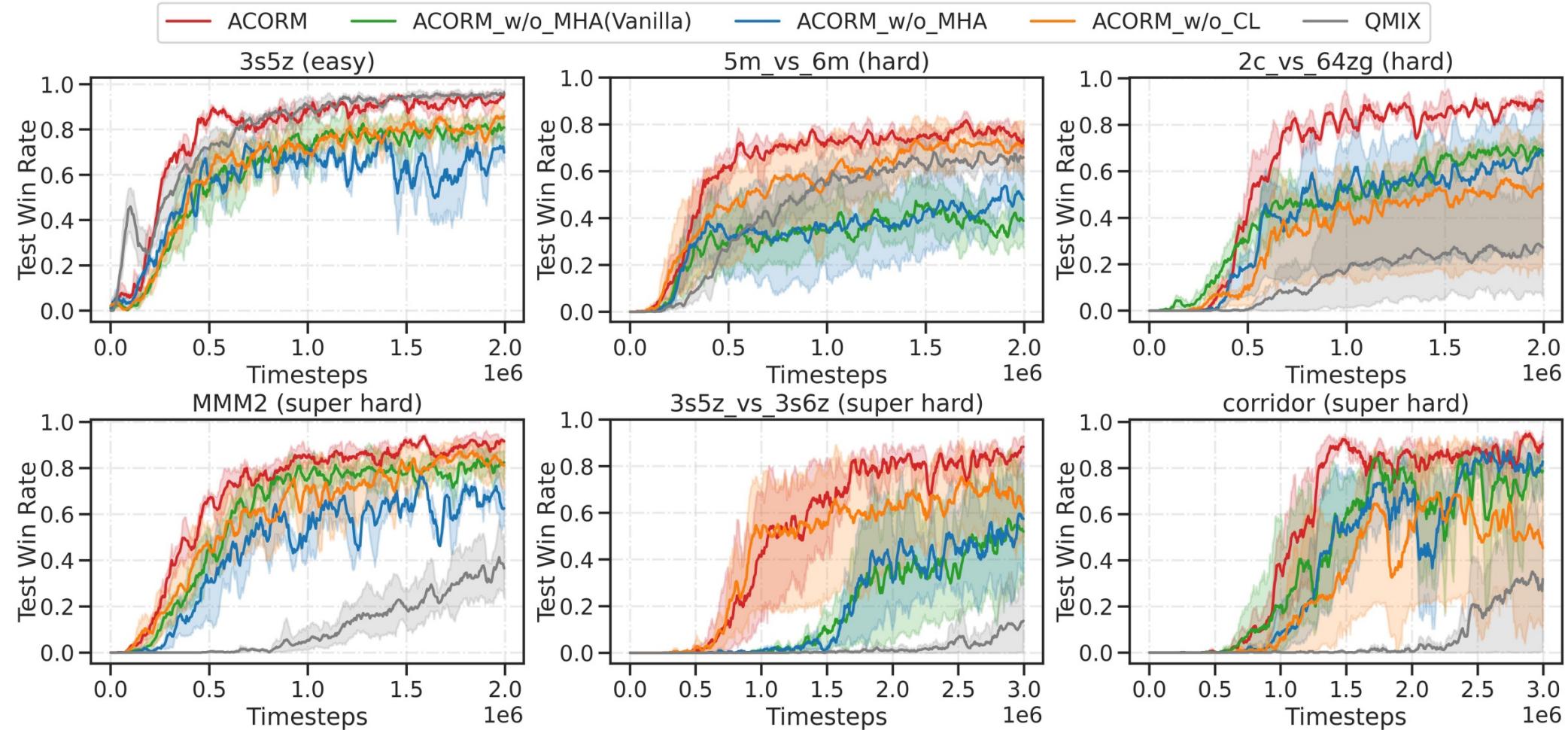


Figure 1: The ACORM framework based on QMIX. (a) The overall architecture. (b) The structure of shared individual Q-network. (c) The detail of contrastive role representation learning, where z_i is the query q , and $z_{i'}/z_{i^*}$ are positive/negative keys k_+/k_- . (d) The attention module that incorporates learned role representations into the mixing network's input for better value decomposition.

Performance on SMAC



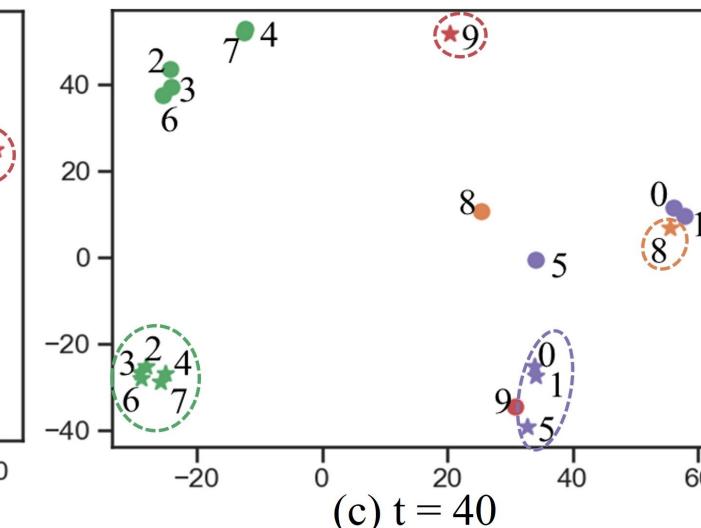
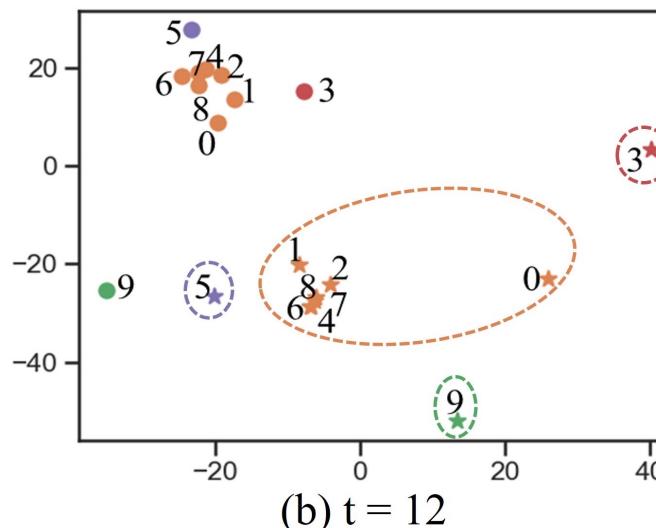
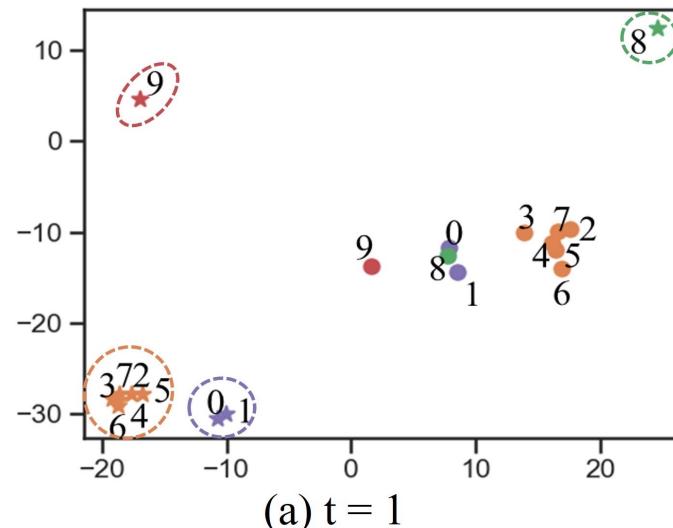
Ablation study



Visualize role representations



● agent embedding ★ role representation ● cluster 1 ● cluster 2 ● cluster 3 ● cluster 4



Visualize attention mechanism



	0	1	2	3	4	5	6	7	8	9	
head 0	16	14	7	7.3	10	12	7.6	7.2	13	6.5	
head 1	11	17	16	7.7	9.8	8.8	8.4	10	7.6	3.8	
head 2	8.1	28	8.9	7	3.2	4.8	3.1	7.9	2.8	26	
head 3	9.2	7.8	13	7.1	9.6	8.8	11	9.5	6.5	18	
	0	1	2	3	4	5	6	7	8	9	role

(a) t = 4

	0	1	2	3	4	5	6	7	8	9	
head 0	12	17	7	9.2	13	5.2	15	5.5	12	3.8	
head 1	15	14	9.5	12	7.8	9.8	7.5	7	6	12	
head 2	11	9.6	6.6	6.6	2.8	19	2.7	3.7	1.5	36	
head 3	6.1	11	4.4	8.9	6.6	27	6.9	7.1	4.8	17	
	0	1	2	3	4	5	6	7	8	9	role

(b) t = 10

	0	1	2	3	4	5	6	7	8	9	
head 0	5.5	5.3	15	17	8.4	14	9.2	11	4.8	9.8	
head 1	4.6	4.1	16	16	9.4	15	11	9.1	4.4	9.3	
head 2	8.1	8.4	9.3	10	11	11	9.8	17	5.7	9.8	
head 3	5.3	1.9	15	17	4	12	5.5	9.5	18	12	
	0	1	2	3	4	5	6	7	8	9	role

(c) t = 36

Conclusions



- A general **role representation learning framework** (分工)
- Leverage role representations to realize more expressive **credit assignment** (协作)
- Tackle **agent homogenization** and facilitate efficient **knowledge transfer**

Paper: <https://openreview.net/forum?id=LWmuPfEYhH>

Code: <https://github.com/NJU-RL/ACORM>

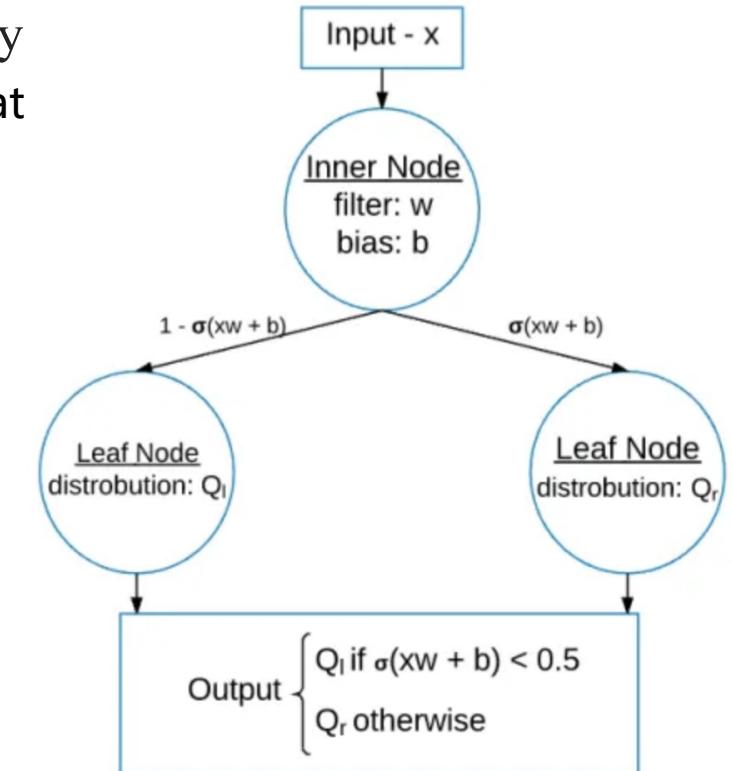
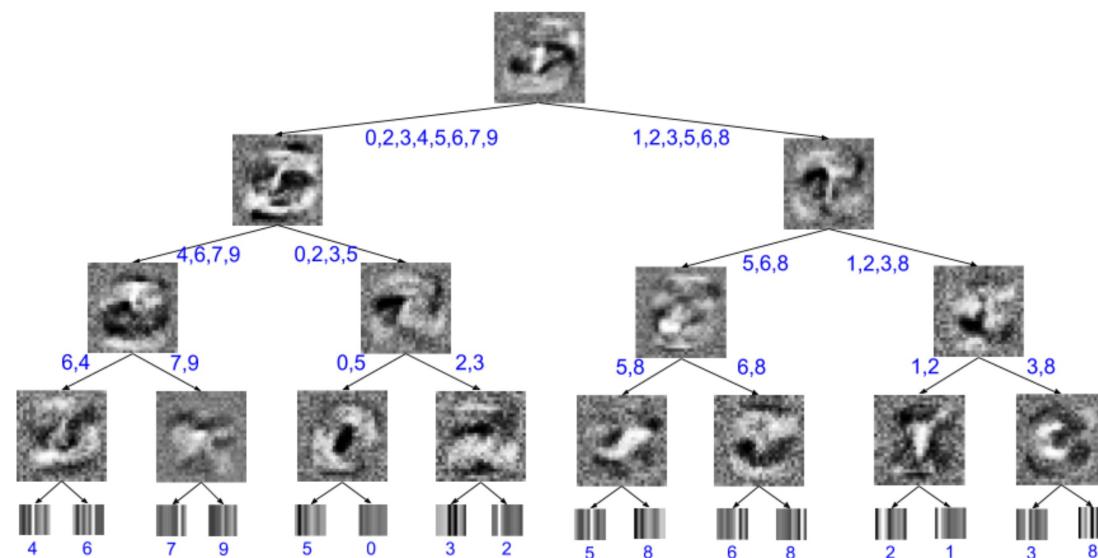
Contents



- Background on RL and MARL
 - Contrastive Role Representations for MARL
 - Interpretable MARL via Mixing Recurrent Soft Decision Trees
-

Soft Decision Trees

- Decision trees have long been valued for their simplicity and interpretability
 - mimic human decision-making processes by splitting data into branches at binary decision points, making them intuitive to understand and explain
- The term “soft” decision trees extends this concept further
 - incorporate elements of neural networks to enhance flexibility and performance



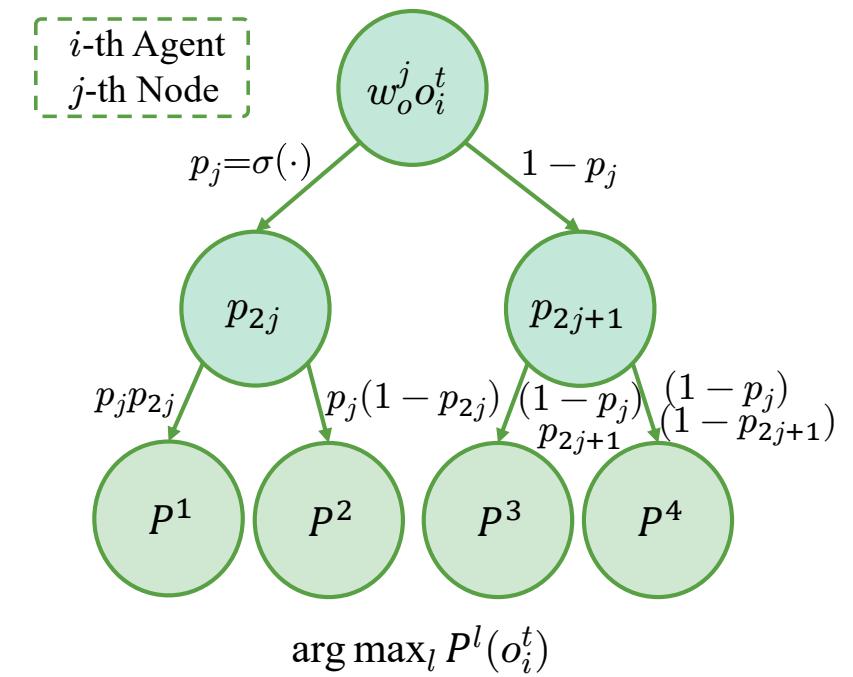
A soft binary decision tree with a single inner node and two leaf nodes

Soft Decision Trees for MARL



■ Motivation

- Soft decision trees for the Q-function
 - Differentiable structure, soft decision boundaries
 - Good representation ability
 - Good interpretability for decision problems



Recurrent Soft Decision Trees

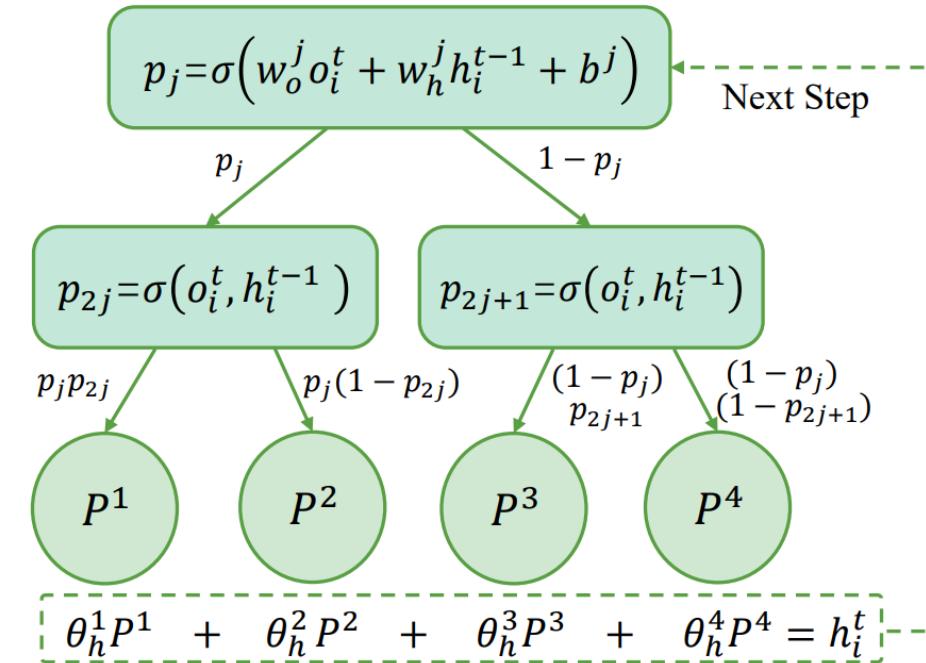


■ Key insight

- Incorporate history information akin to RNNs

$$h_i^t = \sum_{l \in Leaf\ Nodes} P^l(o_i^t, h_i^{t-1}) \theta_h^l$$

$$Q_i(\tau_i, \cdot) = h_i^t w_q$$



Ensemble Recurrent Soft Decision Trees

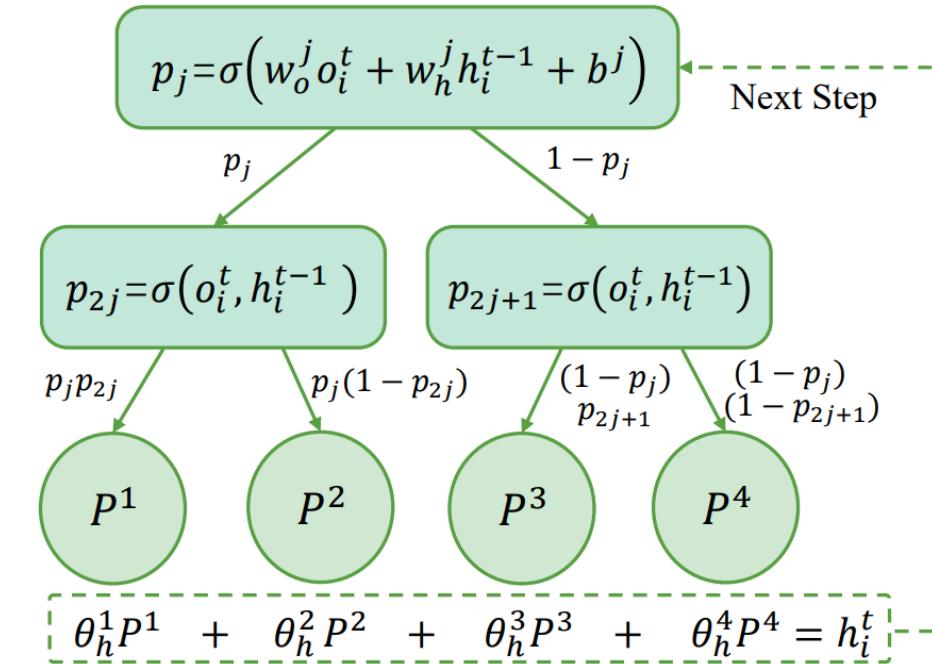


■ Key insight

- Increase representation power
- Ensure interpretability
- Linear ensembling

$$h_i^t = \sum_{l \in LeafNodes} P^l(o_i^t, h_i^{t-1}) \theta_h^l \quad Q_i(\tau_i, \cdot) = h_i^t w_q$$

$$h_i^t = [h_{i,(1)}^t, h_{i,(2)}^t, \dots, h_{i,(H)}^t]$$
$$Q_i(\tau_i, \cdot) = h_{i,(1)}^t w_{q,(1)} + h_{i,(2)}^t w_{q,(2)} + \dots + h_{i,(H)}^t w_{q,(H)}$$



Mixing Tree Architecture



■ Key insight

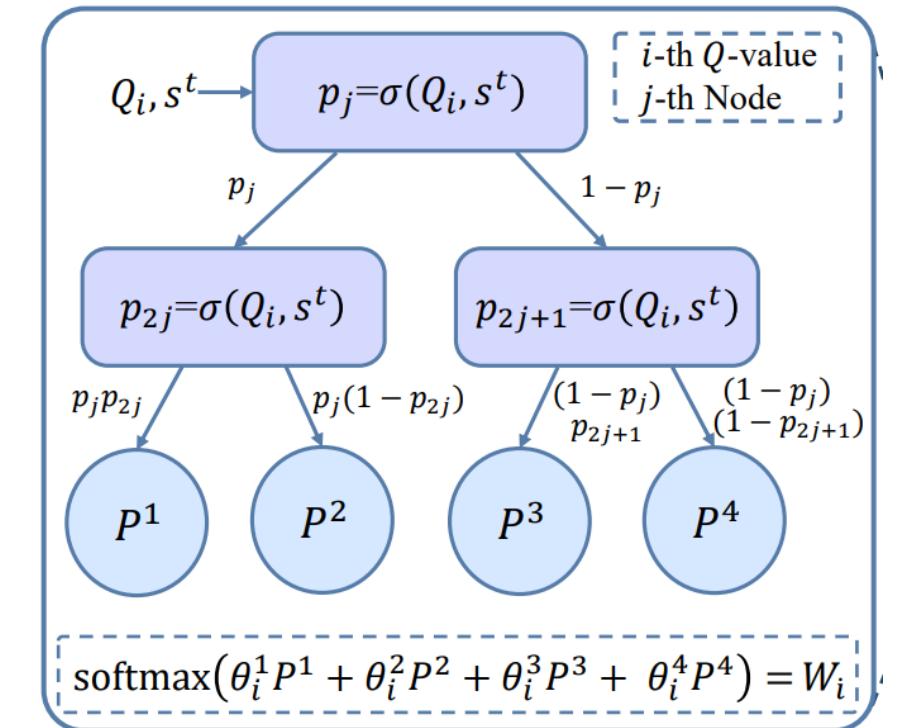
- **Value decomposition** using soft decision tree

$$p_j(Q_i, s^t) = \sigma(w_q^j Q_i + w_s^j s^t + b^j)$$

$$\phi_i = \sum_{l \in Leaf\ Nodes} P^l(Q_i, s^t) \theta_i^l,$$

$$W_i = \frac{\exp(\sum_{k=1}^H \phi_{i,(k)} w_{\phi,(k)})}{\sum_{i=1}^n \exp(\sum_{k=1}^H \phi_{i,(k)} w_{\phi,(k)})},$$

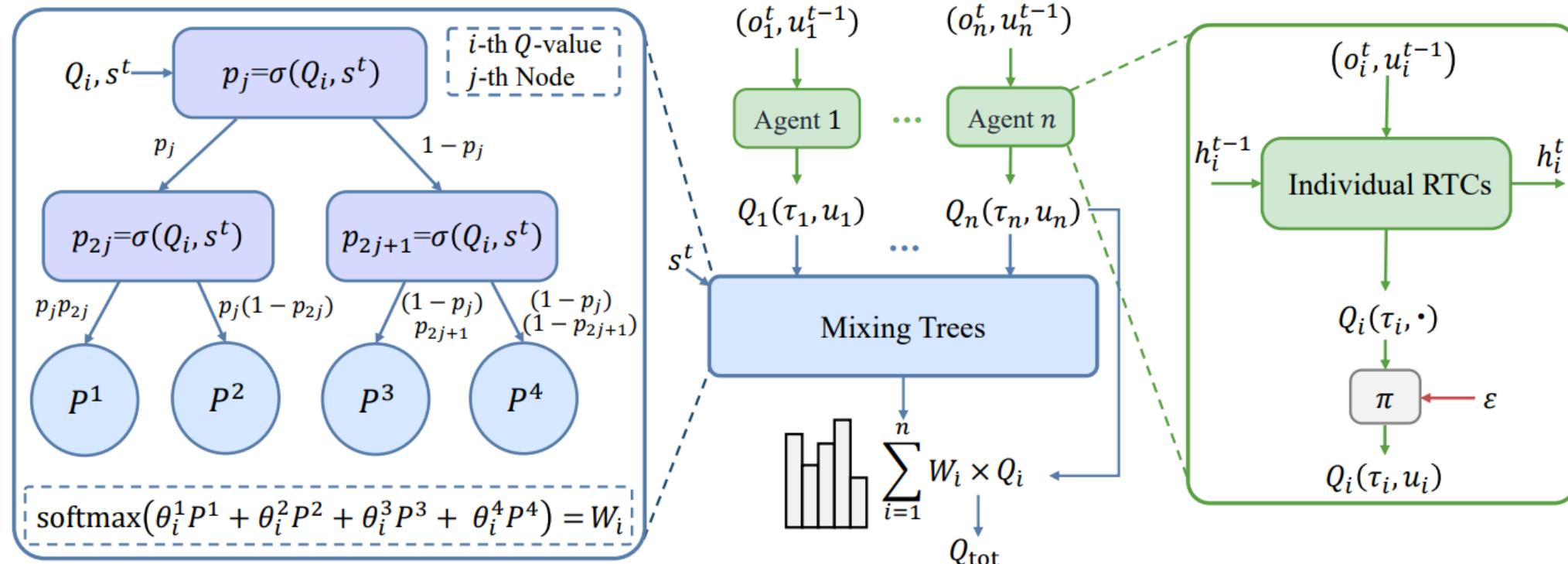
$$Q_{tot}(\tau, u) \approx \sum_{i=1}^n W_i Q_i(\tau_i, u_i)$$



MIXing Recurrent soft decision Trees (MIXRTs)



■ Put it together

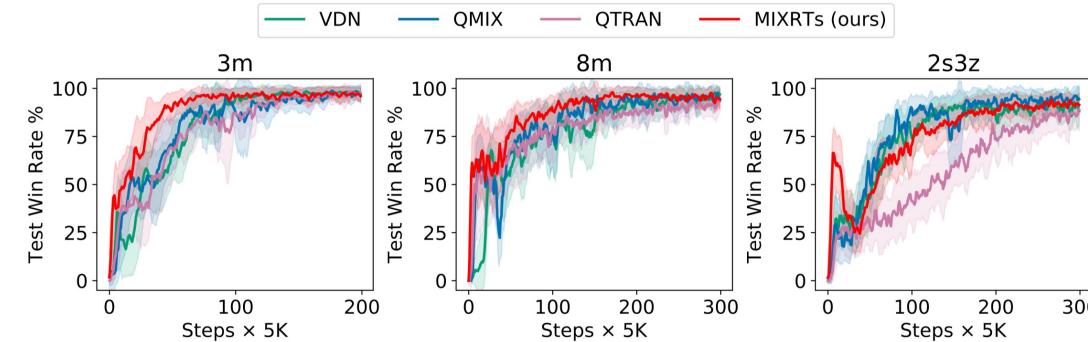


Primary Results: Performance



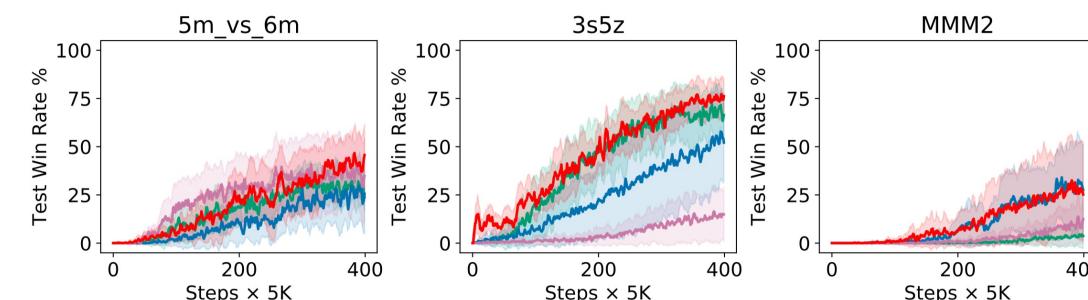
➤ 在简单场景上的性能

- 快速掌握简单任务
- 模型取得较高的胜率



➤ 在困难/超困难场景上的性能

- 实现具有竞争力的性能
- 在学习过程中更加稳定

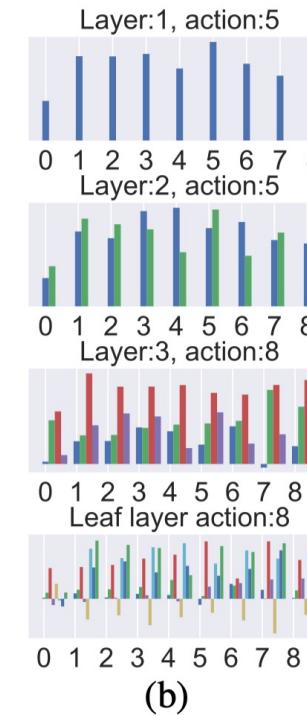
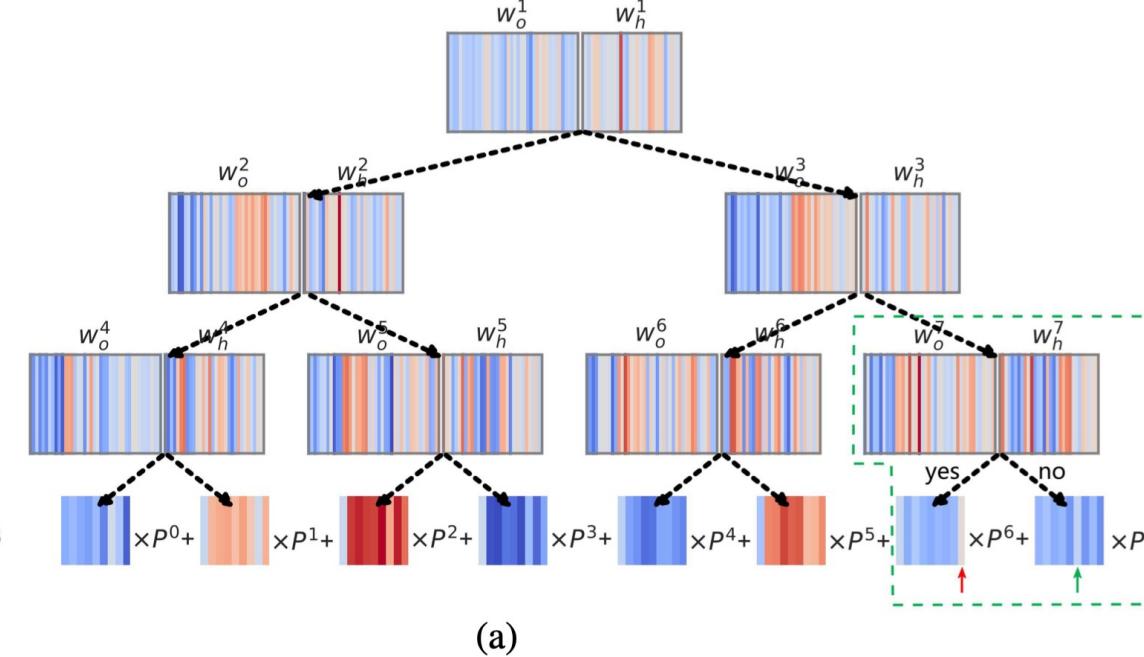
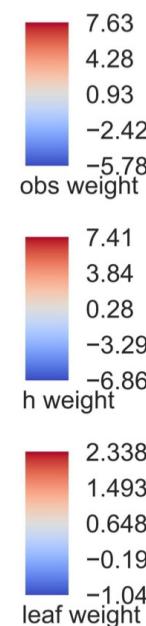


➤ 模型参数比较

- 线性模型、参数较少

Method	3m	2s3z	5m_vs_6m	3s5z	8m_vs_9m	MMM2	6h_vs_8z
VDN	28,297	31,883	30,412	35,534	32,911	39,250	32,206
QMIX	37,738	62,892	55,789	111,951	96,304	173,651	72,847
QTRAN	70,911	84,437	80,518	101,492	94,513	120,320	88,436
MIXRTs (ours)	20,880	34,448	28,752	48,560	38,592	62,736	35,440

Primary Results: Interpretability



➤ 解释软决策树的结构

- 更红的颜色意味着获得更高的特征权重
- 在绿色方框中，位置17、14和29处有着更红的颜色，表示敌人是否可见、自身的生命值等特征，易发现智能体更倾向于攻击敌人

Primary Results: Interpretability



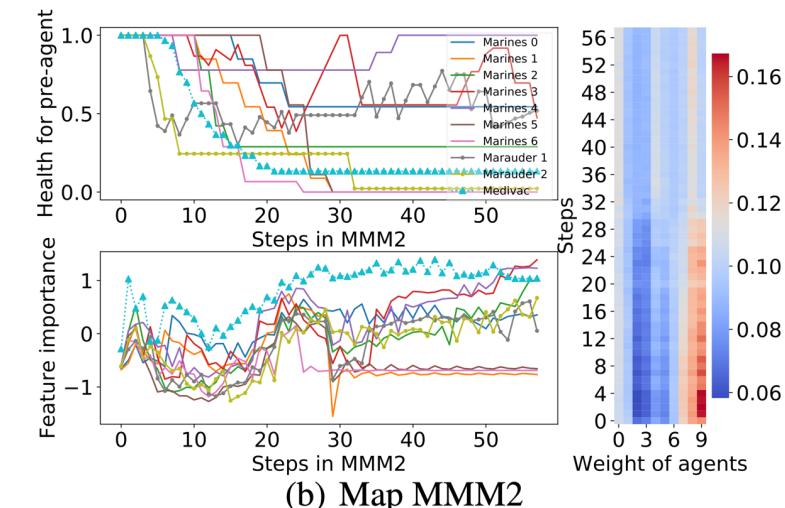
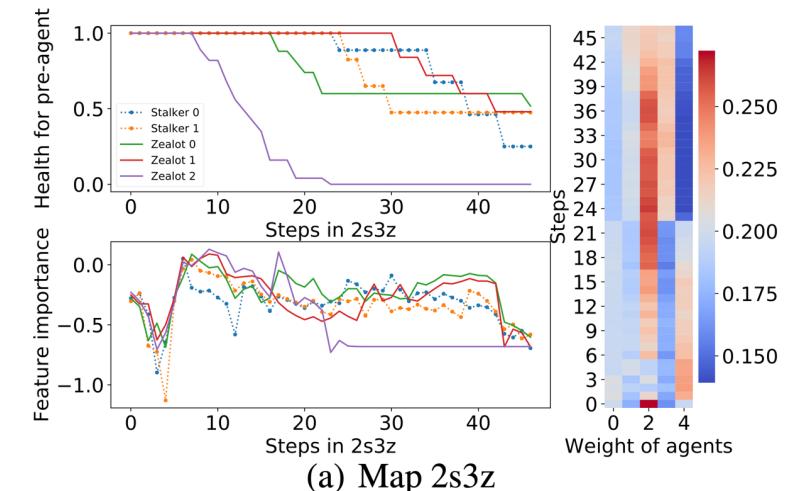
■ 案例分析

- 对特征的重要度解释

- 特征重要度的定义

$$I(o_i^t) = \sum_j P^j(o_i^t, h_i^{t-1}) w_o^j$$

- 特征重要度在某一回合中：以健康属性为例
 - 相同的智能体有着相同趋势
 - 死亡的智能体具有较低的重要度
 - 医疗船在战斗中具有较高的重要度
 -





Thank You.

Zhi Wang (王志)

<https://heyuanmingong.github.io>

Email: zhiwang@nju.edu.cn

Nanjing University, China

2025-04-27

