

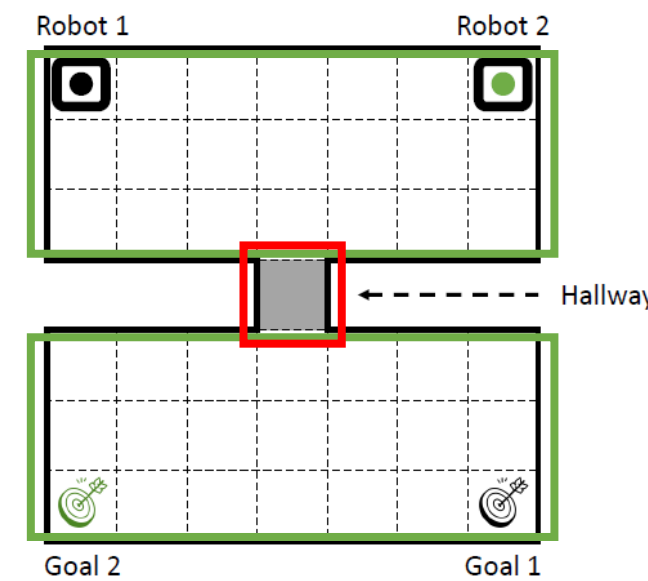
# Multi-Agent Sparse Interaction Modeling is An Anomaly Detection Problem

Chao Li<sup>1</sup>, Shaokang Dong<sup>1</sup>, Shangdong Yang<sup>2,3</sup>, Hongye Cao<sup>1</sup>, Wenbin Li<sup>1</sup>, Yang Gao<sup>1</sup><sup>1</sup> State Key Laboratory for Novel Software Technology, Nanjing University <sup>2</sup> School of Computer Science, Nanjing University of Posts and Telecommunications<sup>3</sup> Guangxi Key Lab of Multi-Source Information Mining & Security, Guangxi Normal University

chaoli1996@smail.nju.edu.cn

## 1. Motivation

- Multi-Agent Reinforcement Learning (MARL) heavily suffers from **sample inefficiency** problem
- Algorithm perspective: the trial-and-error paradigm inherent in RL;
- Task perspective: policy search over enormous state-joint action space.



### Sparse Interaction Property

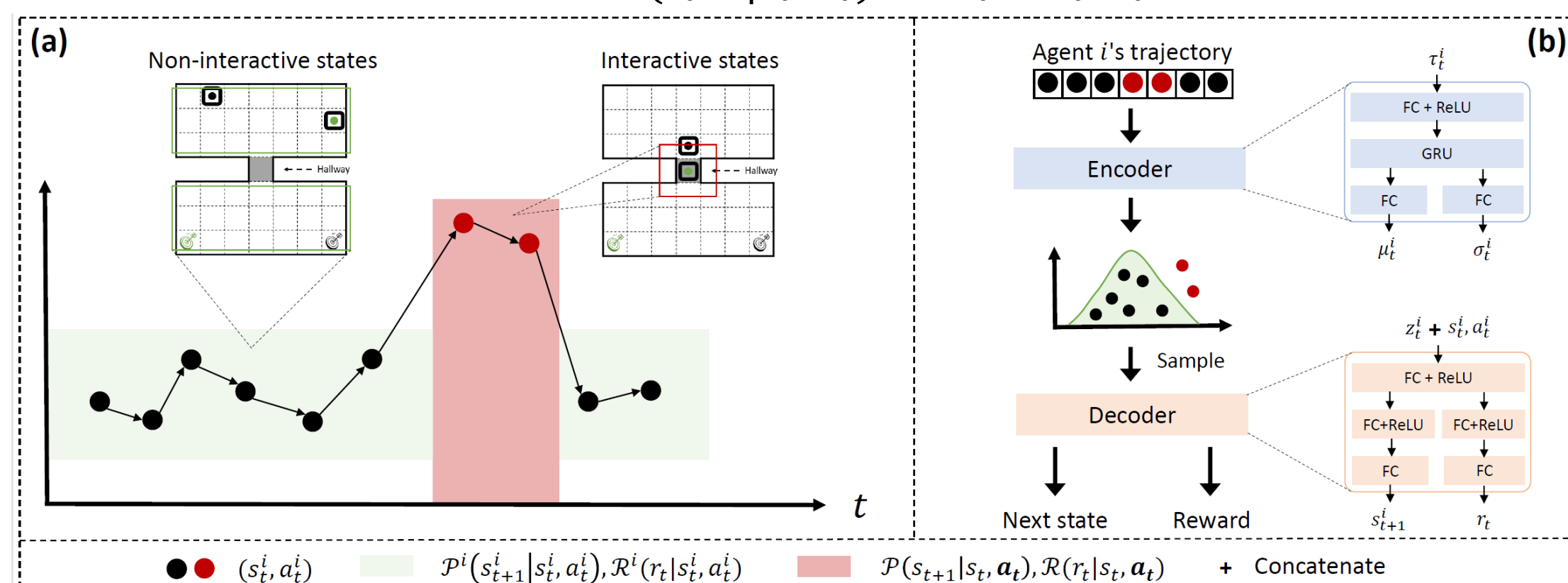
- On **rare** interactive states where agents need to coordinate, the multi-agent task continues;
- On **commonplace** non-interactive states where agents can act independently to achieve desired outcomes, task reduces to a single-agent task.

- Our focus: improve MARL from the task perspective
- Model the sparse interaction structure among agents;
- Utilize this structure to instruct per-agent policy learning.

## 2. Methodology

### 2.1 Model the sparse interaction by dynamics

- In per-agent local trajectory, **rare** interactive states (**outliers**) adhere to the dynamics of the multi-agent task:  $P(s_{t+1}|s_t, \mathbf{a}_t), R(r_t|s_t, \mathbf{a}_t)$ ;
- In contrast, **common** non-interactive states adhere to the dynamics of the single-agent task:  $P^i(s_{t+1}^i|s_t^i, \mathbf{a}_t^i), R^i(r_t^i|s_t^i, \mathbf{a}_t^i)$ .



### 2.2 Dynamics Distribution Modeling

- We characterize the trajectory dynamics by a latent variable, and the objective is defined as:  $\min \mathcal{D}_{\text{KL}}(q(z^i | \tau_i^i) \| p(z^i | \tau_i^i))$
$$\mathcal{D}_{\text{KL}}(q(z^i | \tau_i^i) \| p(z^i | \tau_i^i)) = \mathcal{D}_{\text{KL}}(q(z^i | \tau_i^i) \| p(z^i)) - \mathbb{E}_{z^i \sim q(z^i | \tau_i^i)} \log p(\tau_i^i | z^i) + \log p(\tau_i^i).$$

$$\max \mathbb{E}_{z^i \sim q(z^i | \tau_i^i)} \log p(\tau_i^i | z^i) + \mathcal{D}_{\text{KL}}(q(z^i | \tau_i^i) \| p(z^i))$$

$$\max \mathbb{E}_{z^i \sim q(z^i | \tau_i^i)} \sum_{t=0}^{T-1} \log p(s_{t+1}^i, r_t^i | s_t^i, a_t^i, z^i) - \mathcal{D}_{\text{KL}}(q(z^i | \tau_i^i) \| p(z^i)).$$

Encoder    Decoder (Reconstruction Likelihood)    KL Divergence between posterior and prior

- We achieve this optimization using a VAE-like network.

### 2.3 Interactive State Discovery

- For interactive states, the reconstruction likelihood is poor
- We define the prediction discrepancy as follows:

$$Dis(s_t^i, a_t^i) = (s_{t+1}^i - \hat{s}_{t+1}^i)^2 + (r_t^i - \hat{r}_t^i)^2.$$

- A large value indicates that current state-action is interactive.

### 2.4 Interaction-Instructed Exploration

- We use the prediction discrepancy as the intrinsic reward:
$$r^i(s_{t+1}^i, a_t^i) = Dis(s_t^i, a_t^i) = (s_{t+1}^i - \hat{s}_{t+1}^i)^2 + (r_t^i - \hat{r}_t^i)^2.$$

**Intuition:** This intrinsic reward encourages agents to explore more on interactive states, enhancing their coordinated behaviors and thus accelerating coordinated policy learning.

- Consequently, the MARL algorithm aims to maximize the total reward of all agents, which is defined as follows:

$$r(s_t, \mathbf{a}_t) = r_t + \alpha \sum_i r^i(s_{t+1}^i, a_t^i),$$

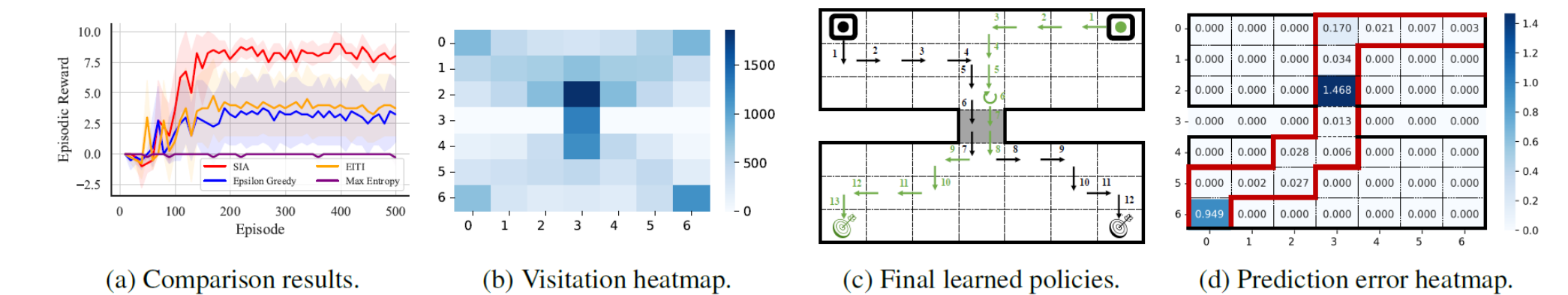
- where  $r_t$  denotes the extrinsic task reward.

## 3. Experiments

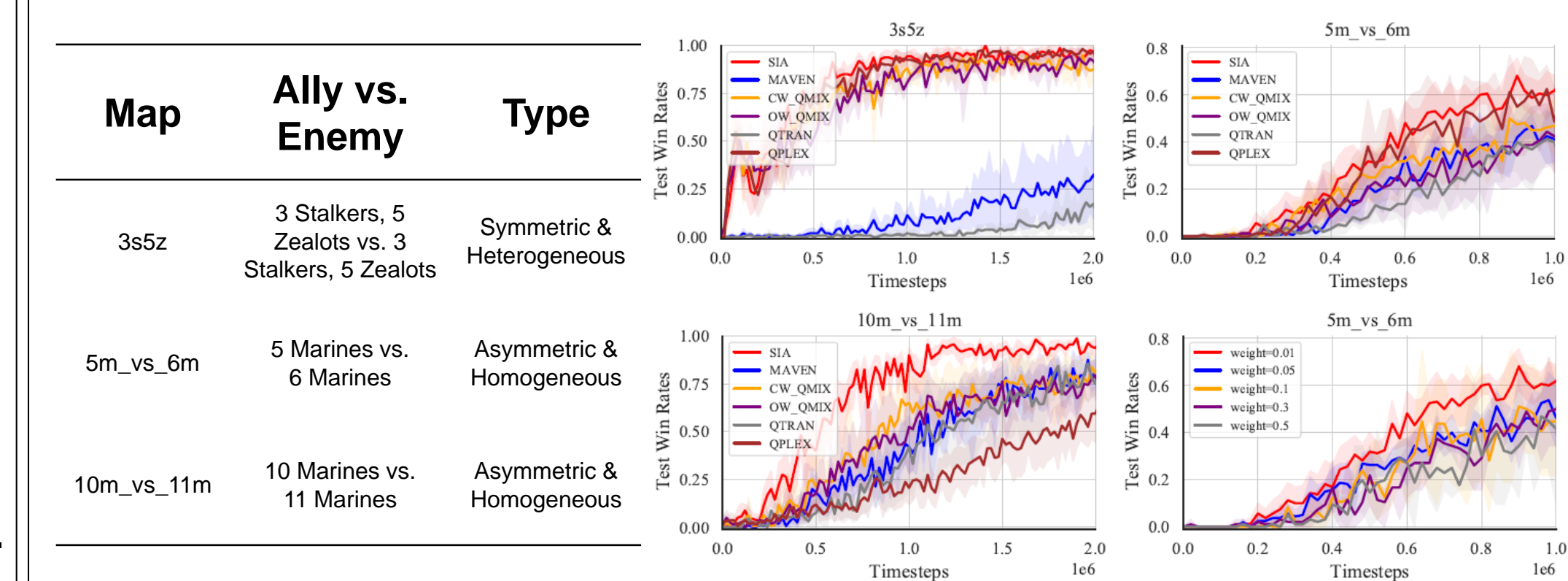
### Benchmarks

- Hall Way (A Didactic Example);
- StarCraft Multi-Agent Challenge;

### Results



**Summary:** SIA successfully identifies interactive states, and the interaction instructed exploration encourages more exploration on them, leading to superior performance.



**Summary:** The superior performance on complex tasks further verifies its effectiveness.

## 4. Conclusion

- Modeling the interaction structure among agents and utilizing it to improve MARL is promising. More ways are worth further exploration;
- In future, we would focus on the nearly decomposable property of multi-agent tasks to enhance multi-agent coordination on large-scale scenarios!