# Multi-Agent Learning from Learners
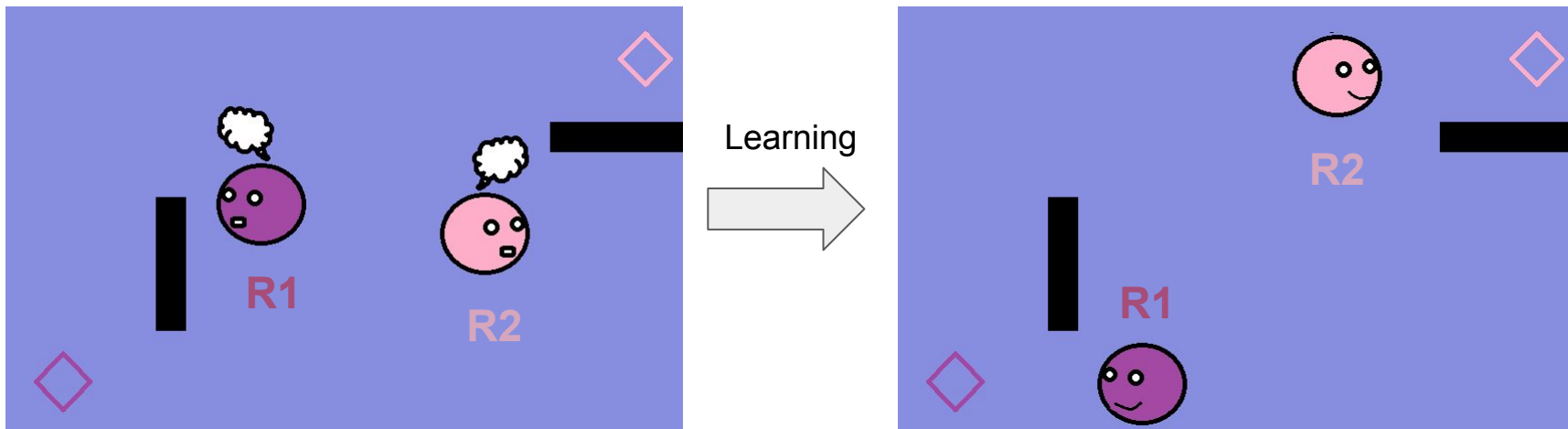
Mine Melodi Caliskan, Francesco Chini, Setareh Maghsudi
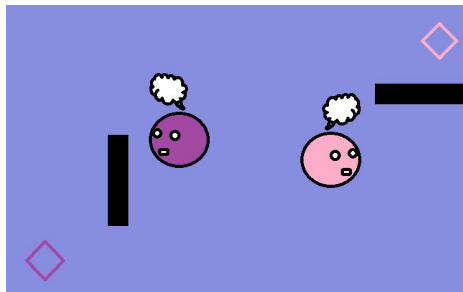
# Introduction

We study the "Learning from a Learner" problem in multi-agent setting
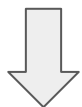
➔ **Goal:** Infer the reward functions of other agents that you interact with who are **not experts** but are **still learning**



Jacq, A., Geist, M., Paiva, A., and Pietquin, O. Learning from a learner. In International Conference on Machine Learning, pp. 2990–2999. PMLR, 2019
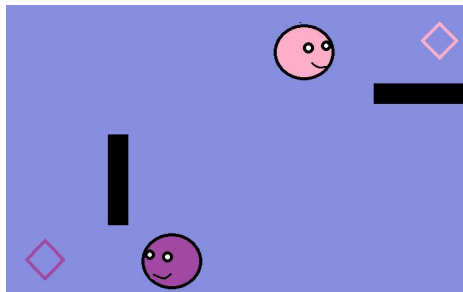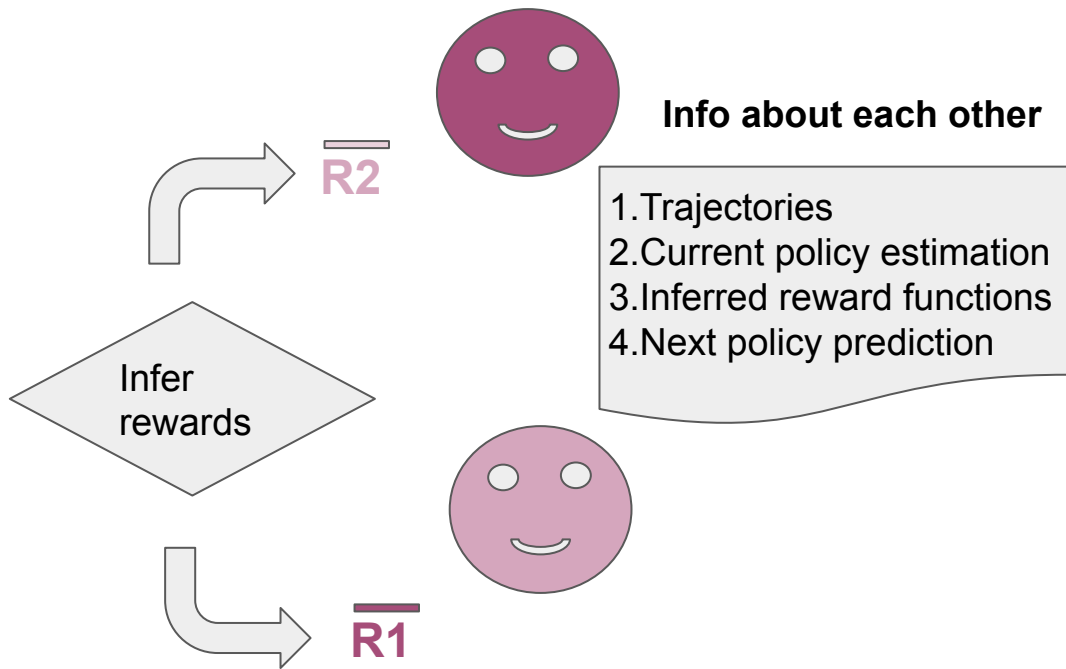
# Introduction



h=0

Policy Improvements

h=20

Infer rewards

R2

R1

**Info about each other**

1. Trajectories
2. Current policy estimation
3. Inferred reward functions
4. Next policy prediction

# Introduction

**Potential applications:**

- **Autonomous cars:** Cars from different companies might have different reward functions e.g safety or energy efficiency, shared environment and no equilibrium
- ➔ Predict behaviour using recovered reward functions
- **Fairness:** The agents might use the information about other agents' rewards in order to learn altruistic behaviours
- **Decentralization:** Use the information about the reward function to decentralize MARL algorithms that requires reward information e.g. Nash Q-learning

# Problem Setting

- $N$ agents acting together in the same environment
- Each agent $i$ is trying to maximize its own reward $R_i$ (general-sum)
- Agents can only observe the state s the actions $a_1, \ldots, a_N$ performed by the other agents and their own reward $R^i(s, a_1, \ldots, a_N)$

→ Assume agents are optimizing entropy-regularized objective (individually):

$$\mathcal{J}(\pi^i) = \mathbb{E}_{\pi^i, \boldsymbol{\pi}^{-i}} \left[ \sum_{t \geq 0} \gamma^t \left( R^i(s_t, \boldsymbol{a}_t) + \alpha \mathcal{H} \left( \pi^i(\cdot | s_t) \right) \right) \right]$$

# Modeling other agents while optimizing your own policy

1. **Policy Improvements:** Multi-Agent Soft Policy Iteration (MA-SPI)

- Evaluate

$$\widetilde{Q}_{\text{soft}}^{\pi^i}(s, a^i) = \tilde{R}^i(s, a^i) + \gamma \mathop{\mathbb{E}}_{\boldsymbol{\pi}} \left[ \widetilde{Q}_{\text{soft}}^{\pi^i}(s', a_{\text{new}}^i) + \alpha \mathcal{H}\left(\pi^i(\cdot|s')\right) \right]$$
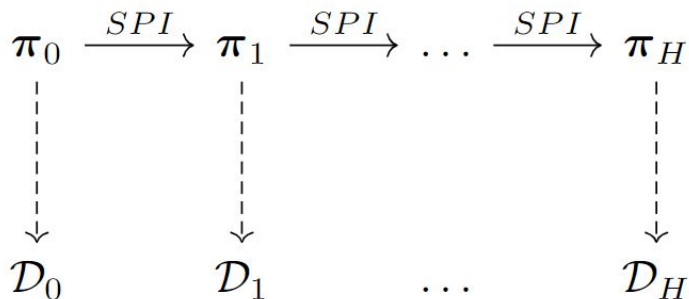
- Improve

$$\pi_{\text{new}}^i(a^i|s) \propto \exp\left(\tfrac{1}{\alpha} \widetilde{Q}_{\text{soft}}^{\pi^i}(s, a^i)\right)$$

# Modeling other agents while optimizing your own policy

**2. Recovering Reward Functions:**

- Estimate policies of the other agents from trajectories

$$\boldsymbol{\pi}_0 \xrightarrow{SPI} \boldsymbol{\pi}_1 \xrightarrow{SPI} \ldots \xrightarrow{SPI} \boldsymbol{\pi}_H$$

$$\mathcal{D}_0 \qquad \mathcal{D}_1 \qquad \ldots \qquad \mathcal{D}_H$$

- Infer reward functions

$$\mathbb{E}_{a^{-i} \sim \boldsymbol{\pi}^{-i}} \left[ \overline{R^i}(s, \boldsymbol{a}^{-i}, a^i) \right] = \alpha \ln \pi_{\text{new}}^i(a^i|s) + \alpha\gamma \mathop{\mathbb{E}}_{\substack{\boldsymbol{a}^{-i} \sim \boldsymbol{\pi}^{-i} \\ s' \sim P(\cdot|s, \boldsymbol{a}^{-i}, a^i)}} \left[ D_{\text{KL}}(\pi^i(\cdot|s') \| \pi_{\text{new}}^i(\cdot|s')) \right]$$
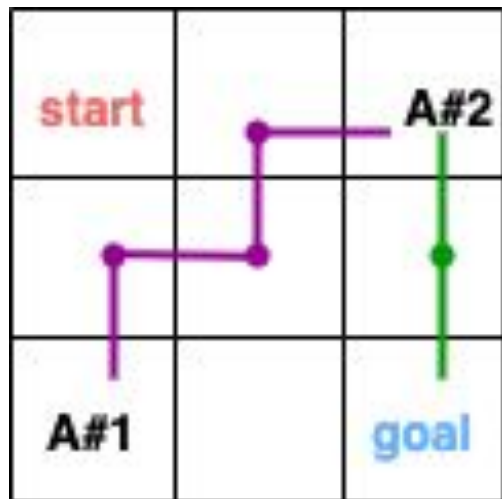
# Error Bounds

In the paper, we provide error bounds

- For the reward recovery in terms of policy estimations
- For the predicted policy improvement in terms of recovered rewards

→ These are novel contributions even in the single-agent case.

# Experiments



$$M_{\text{hom}}^i = -d(\text{agent}_i, \text{goal}) + d(\text{agent}_i, \text{agent}_j) \text{ for } i = 1, 2$$

$$M_{\text{het}}^i = \begin{cases} -d(\text{agent}_i, \text{goal}) - d(\text{agent}_i, \text{agent}_j) & i = 1 \\ \\ -d(\text{agent}_i, \text{goal}) + d(\text{agent}_i, \text{agent}_j) & i = 2 \end{cases}$$

# Results

| Metric | $M_{\text{hom}}$ | $M_{\text{het}}$ |
|--------|------------------|------------------|
| PCC #1 | $0.48 \pm 0.06$ | $0.45 \pm 0.04$ |
| PCC #2 | $0.59 \pm 0.02$ | $0.42 \pm 0.02$ |
| $\hat{P}$ | $\mathbf{0.54} \pm 0.03$ | $\mathbf{0.44} \pm 0.01$ |
| SCC #1 | $0.44 \pm 0.14$ | $0.51 \pm 0.02$ |
| SCC #2 | $0.60 \pm 0.04$ | $0.43 \pm 0.03$ |
| $\hat{S}$ | $\mathbf{0.52} \pm 0.06$ | $\mathbf{0.47} \pm 0.01$ |

Wed 26 Jul 2 p.m. HST — 3:30 p.m. HST
Exhibit Hall 1 #606