

Fully Decentralized Multi-Agent RL

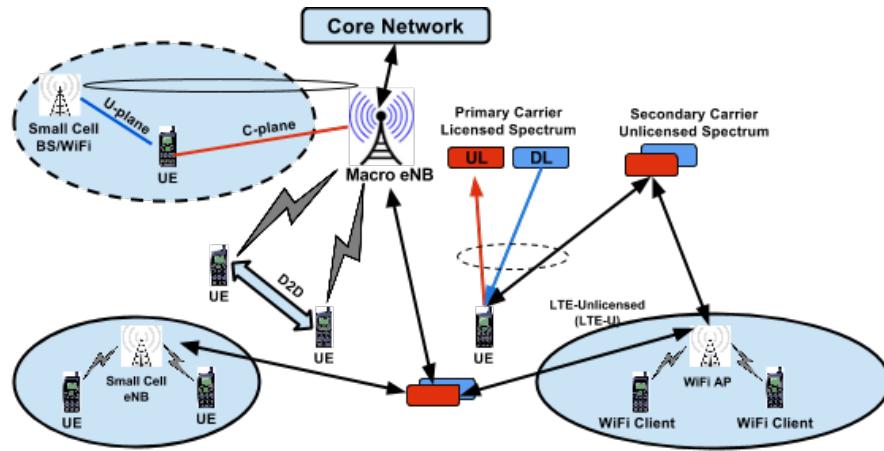
Zongqing Lu

zongqing.lu@pku.edu.cn

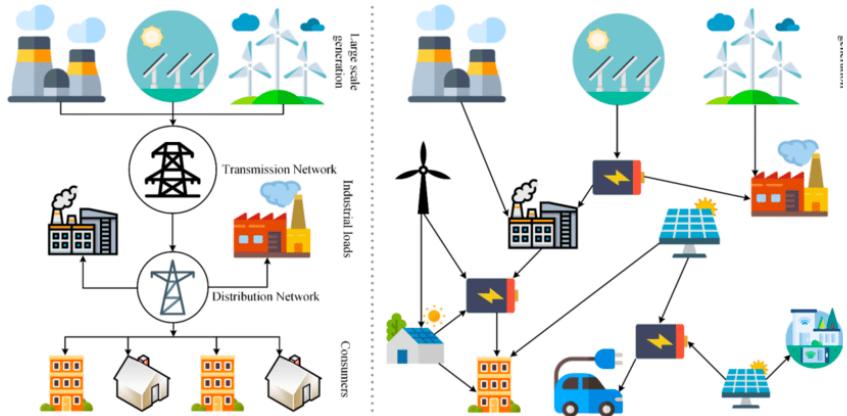




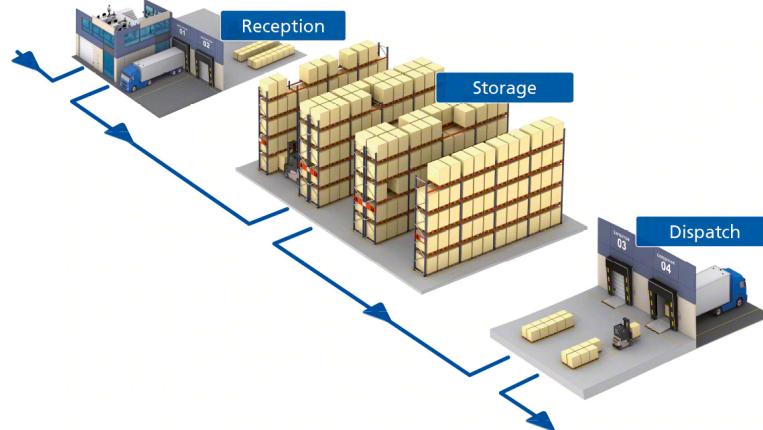
How to control traffic signals in a region to maximize *throughput*?



How to control 5G base stations to optimally *tradeoff coverage and power*?



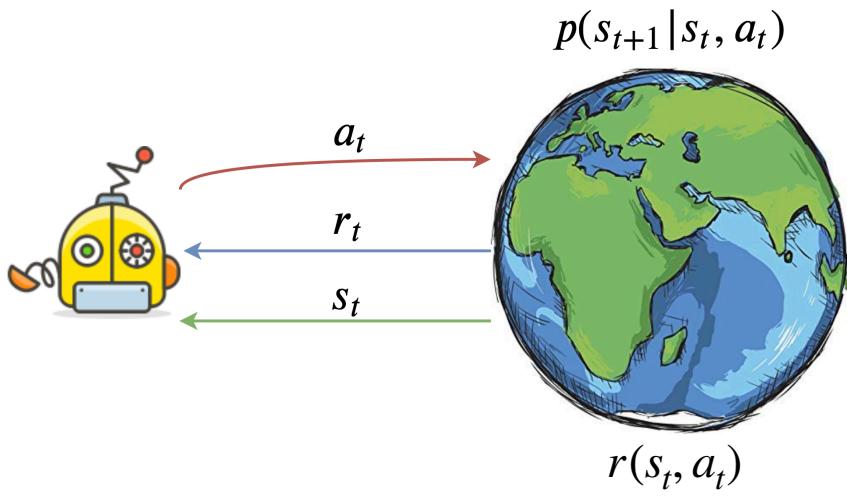
How to control power generators to meet the market demand?



How to optimize supply chain/logistics?

RL and MARL

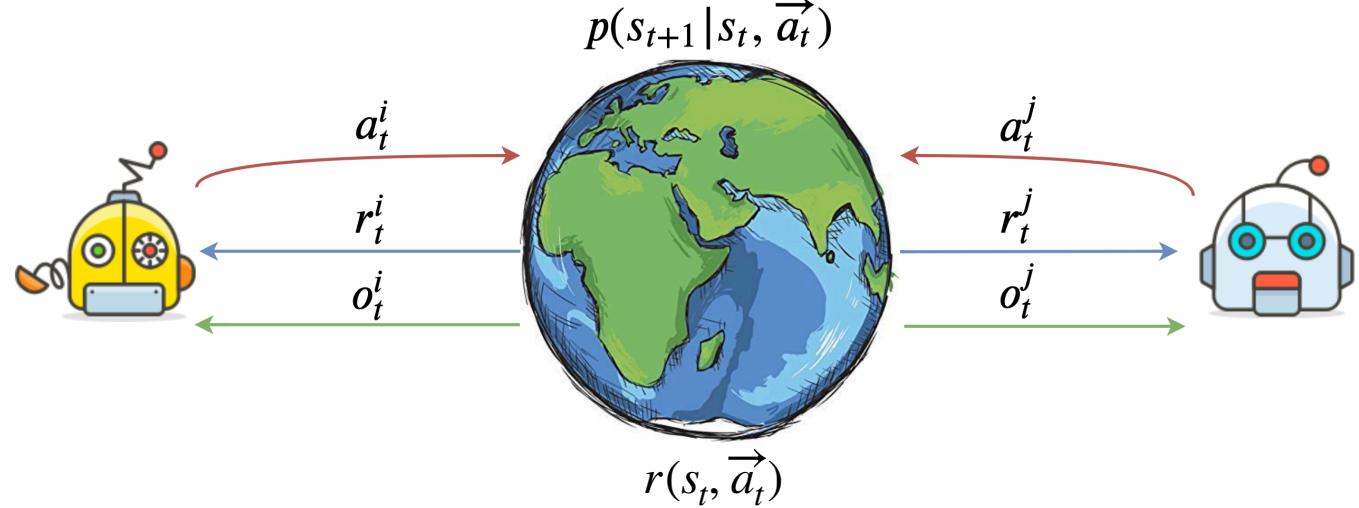
RL



$$\arg \max_{\pi} \mathbb{E}_{\pi} \sum_{t=0}^{\infty} \gamma^t r_t$$

- Q-learning
- Policy gradient (actor-critic)
- Model-based RL

MARL

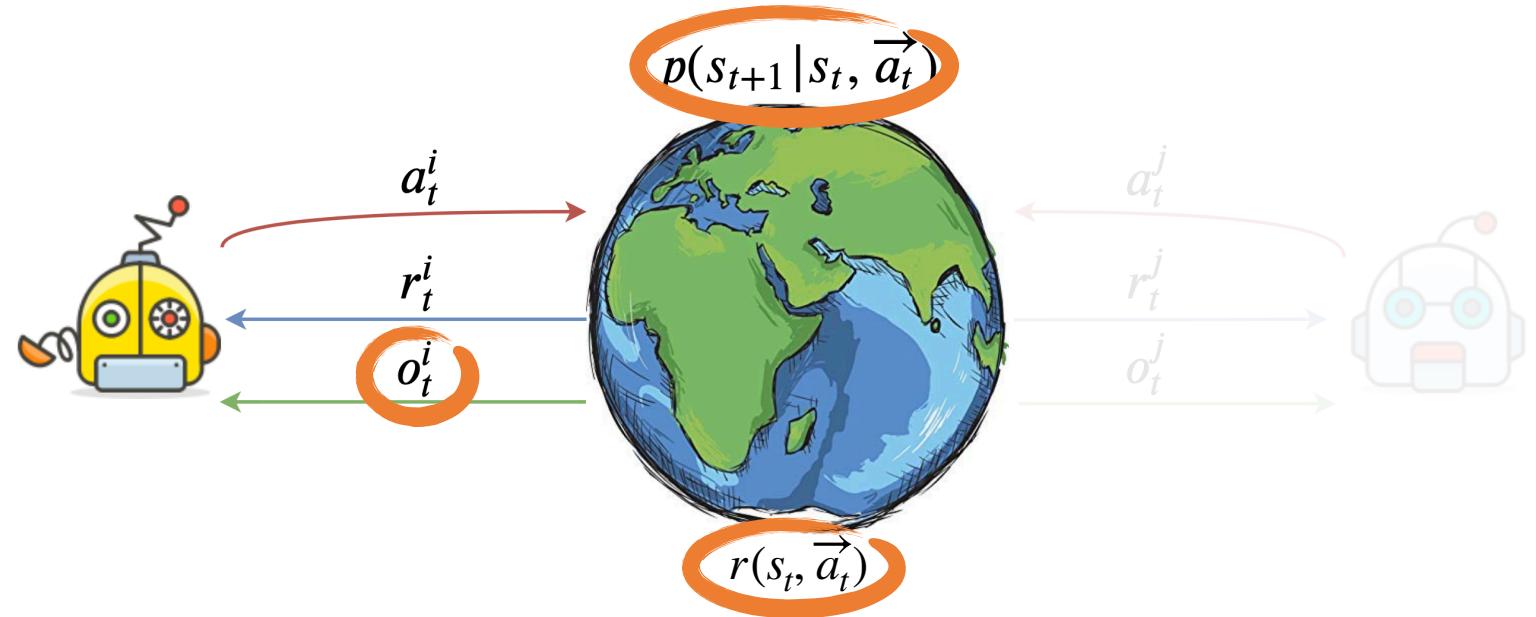


$$r^i = r^j \quad \arg \max_{\pi} \mathbb{E}_{\pi} \sum_{t=0}^{\infty} \gamma^t r_t$$

Cooperative MARL

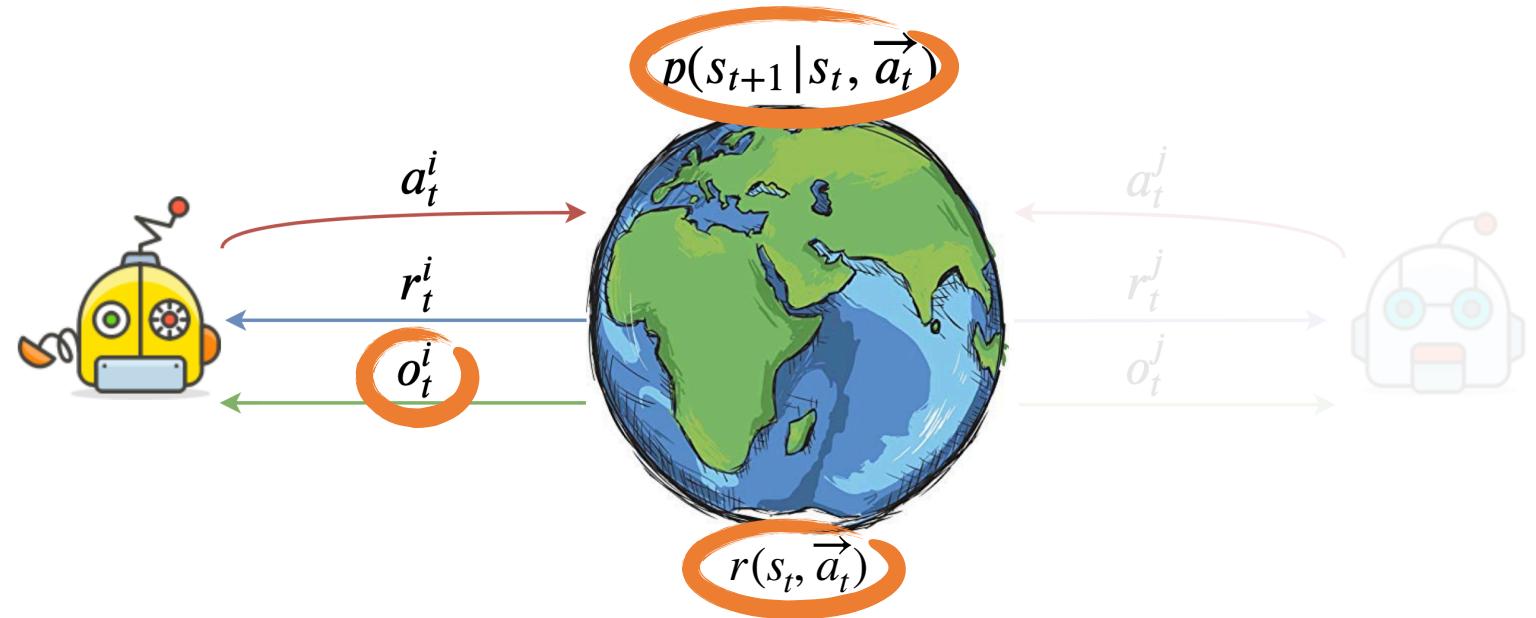
Challenges in cooperative MARL

- Non-stationarity
- Partial observation
- Coordination

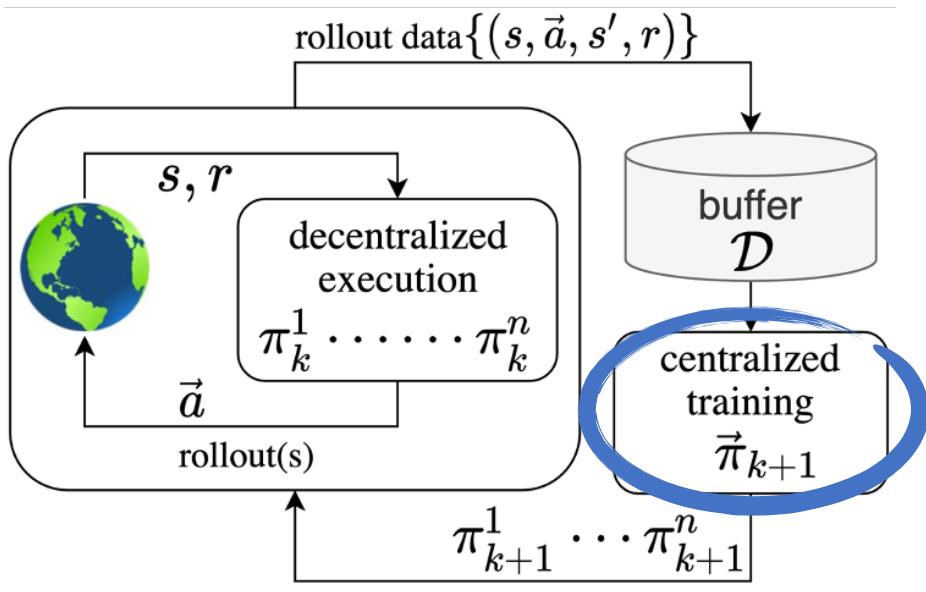


Challenges in cooperative MARL

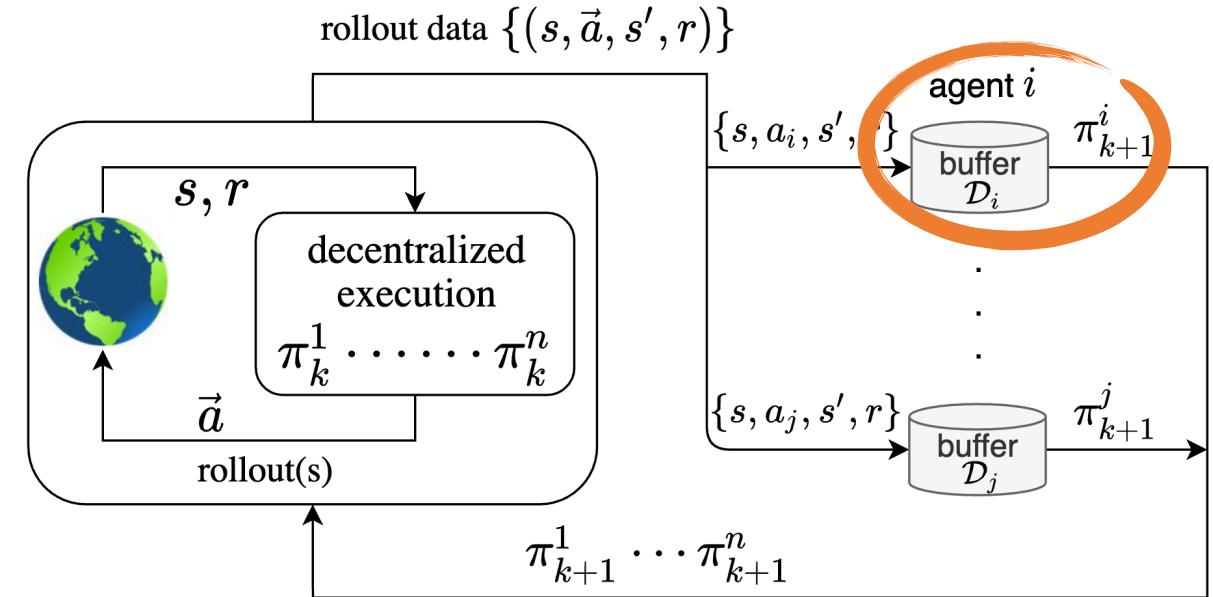
- Non-stationarity
- Partial observation
- Coordination



Learning paradigms



Centralized Training with Decentralized Execution (CTDE)



Decentralized Learning

Why should we care decentralized learning?

- Isn't CTDE enough for MARL?
 - In many cases, there is no centralized entity for training
- Decentralized learning
 - Much better scalability and flexibility
 - A old problem, back to the found of MARL
 - *The way how humans learn!*

Fully decentralized multi-agent RL

- Each agent learns on its **own action**, *without communication or parameter-sharing*
- **Non-stationarity**

$$p(s'|s, \mathbf{a}) \rightarrow p(s'|s, a_i) = \sum_{a_{-i}} \pi_{-i}(a_{-i} | s) p(s'|s, a_i, a_{-i})$$

$$r(s, \mathbf{a}) \rightarrow r(s, a_i) = \sum_{a_{-i}} \pi_{-i}(a_{-i} | s) r(s, a_i, a_{-i})$$

As π_{-i} continuously updates during learning, the environment becomes **non-stationary**
Simple methods, e.g., independent Q-learning/PPO, have **no convergence guarantee!**

Outline

- 1. Fully decentralized Q-learning**
- 2. Fully decentralized policy optimization**

- ✓ Su et al., A Minimalist Approach to Decentralized Multi-Agent Reinforcement Learning, 2022
- ✓ Jiang and Lu, I2Q: A fully decentralized Q-learning algorithm, 2022
- ✓ Su and Lu, Fully decentralized policy optimization, 2022
- ✓ Luo, Jiang, and Lu, Model-based decentralized policy optimization, 2022

IQL

```
1: repeat
2:   all agents interact in the environment
3:   for  $i \leftarrow 1, n$  do
4:     agent  $i$  updates by Q-learning
5:   end for
6: until terminate
```

✗ no convergence guarantee

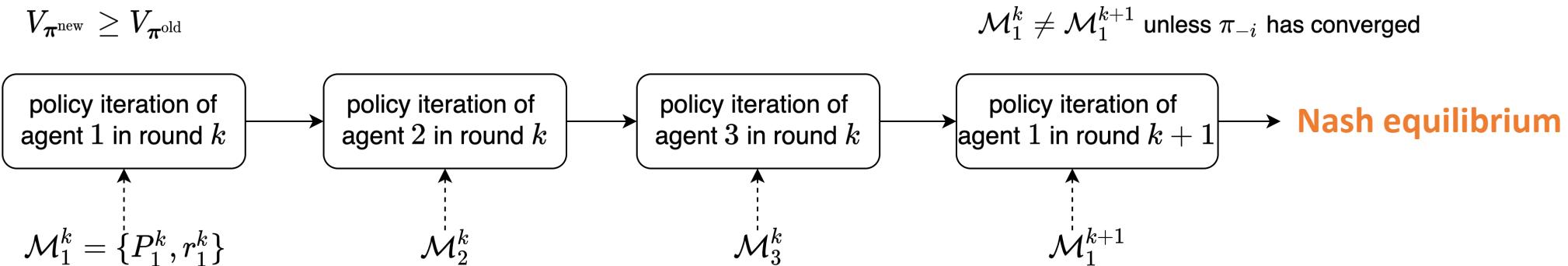
MA2QL

```
1: repeat
2:   for  $i \leftarrow 1, n$  do
3:     all agents interact in the environment
4:     agent  $i$  updates by Q-learning
5:   end for
6: until terminate
```

Will this guarantee convergence?

A simplest fix to non-stationarity

- Fix the policies of other agents while an agent is learning



Multi-agent alternate policy iteration

✓ convergence guarantee

✗ one-step policy improvement may take very long to converge

$$\pi_i^{*,k} = \arg \max_{\pi_i} \mathbb{E}_{\pi_{-i}^k} [Q_{\pi_i, \pi_{-i}^k}(s, a_i, a_{-i})]$$



Q-iteration

ε -convergent Q-iteration

$$\pi_i^{*,k} = \arg \max_{\pi_i} \mathbb{E}_{\pi_{-i}^k} [Q_{\pi_i, \pi_{-i}^k}(s, a_i, a_{-i})]$$

✗ although off-policy at each turn, still asymptotic convergence

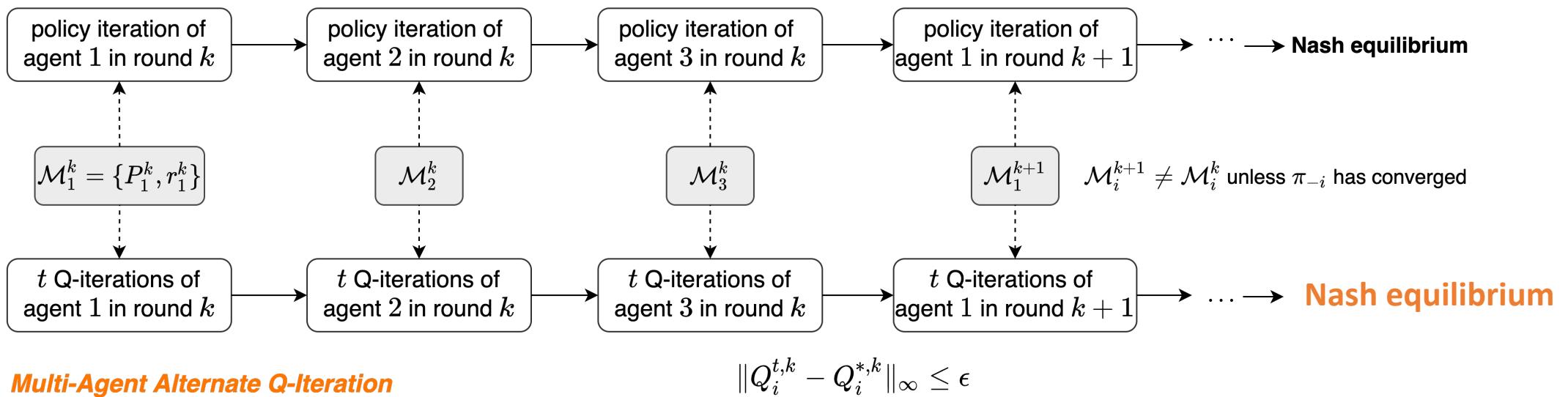


As each agent updates the same Q-function at different turns, can we *truncate* Q-iteration at each turn?

$$\|Q_i^{t,k} - Q_i^{*,k}\|_\infty \leq \varepsilon, \quad \text{when } t \geq \frac{\log((1-\gamma)\varepsilon) - \log(2R + 2\varepsilon)}{\log \gamma}$$

Multi-agent alternate Q-iteration

Multi-Agent Alternate Policy Iteration



If all agents in turn take Q-iteration to $\|Q_i^k - Q_i^{*,k}\|_\infty \leq \epsilon$, then their joint policy sequence $\{\pi^k\}$ converges to a Nash equilibrium, where $\pi_i^k(s) = \arg \max_{a_i} Q_i^k(s, a_i)$.

MA2QL vs. IQL

IQL

```
1: repeat
2:   all agents interact in the environment
3:   for  $i \leftarrow 1, n$  do
4:     agent  $i$  updates by Q-learning
5:   end for
6: until terminate
```

✗ no convergence guarantee

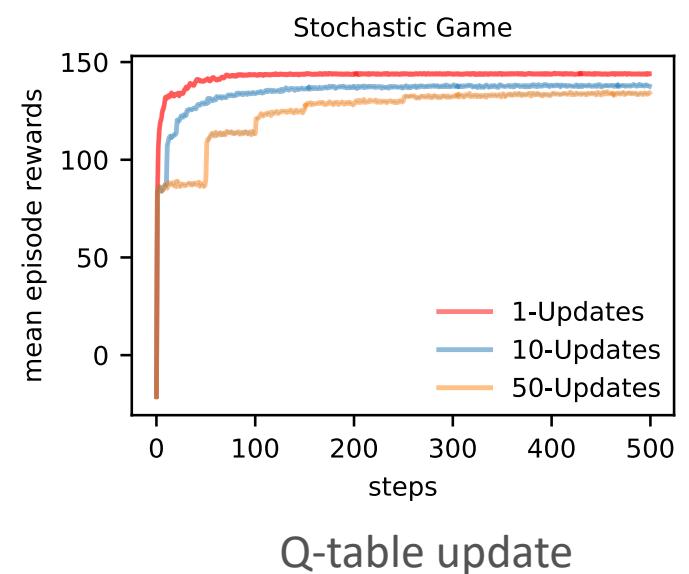
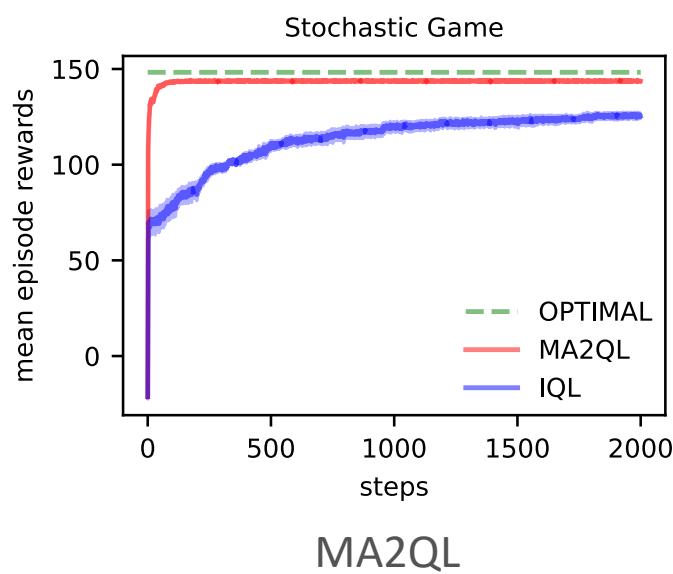
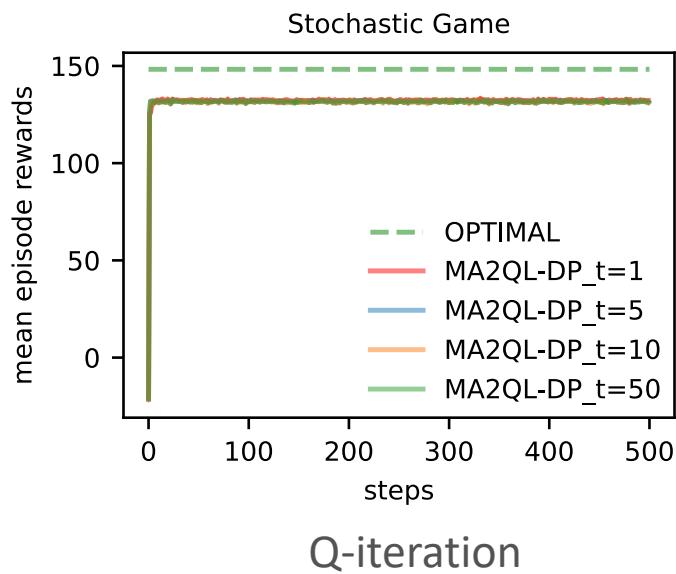
MA2QL

```
1: repeat
2:   for  $i \leftarrow 1, n$  do
3:     all agents interact in the environment
4:     agent  $i$  updates by Q-learning
5:   end for
6: until terminate
```

✓ guaranteed convergence to Nash equilibrium

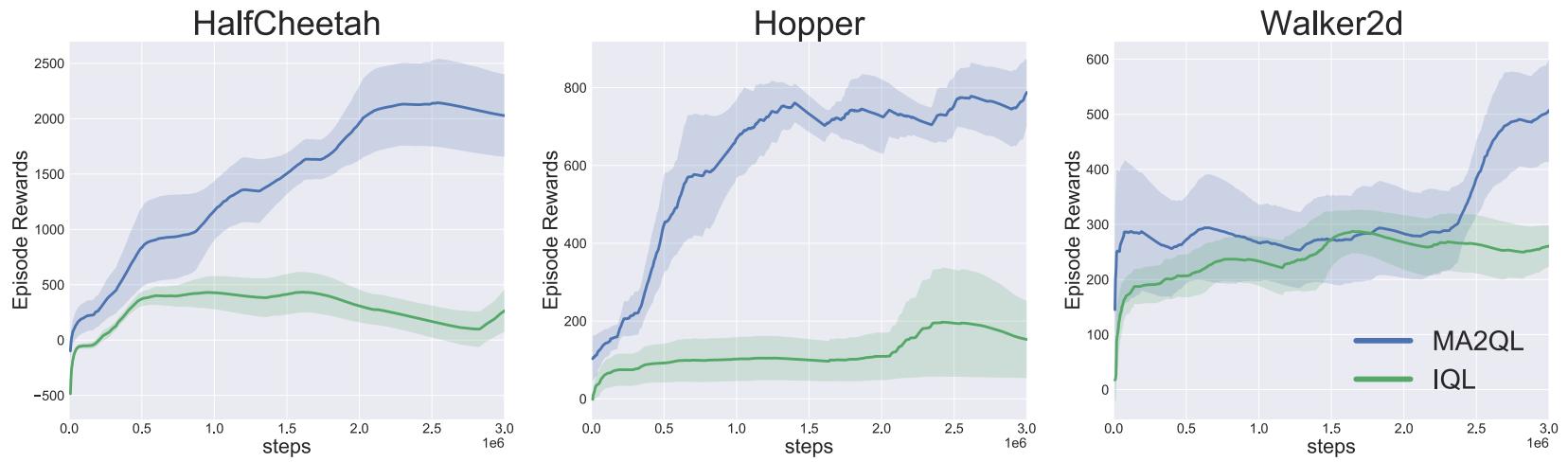
MA2QL vs. IQL

- A didactic game (tabular case)

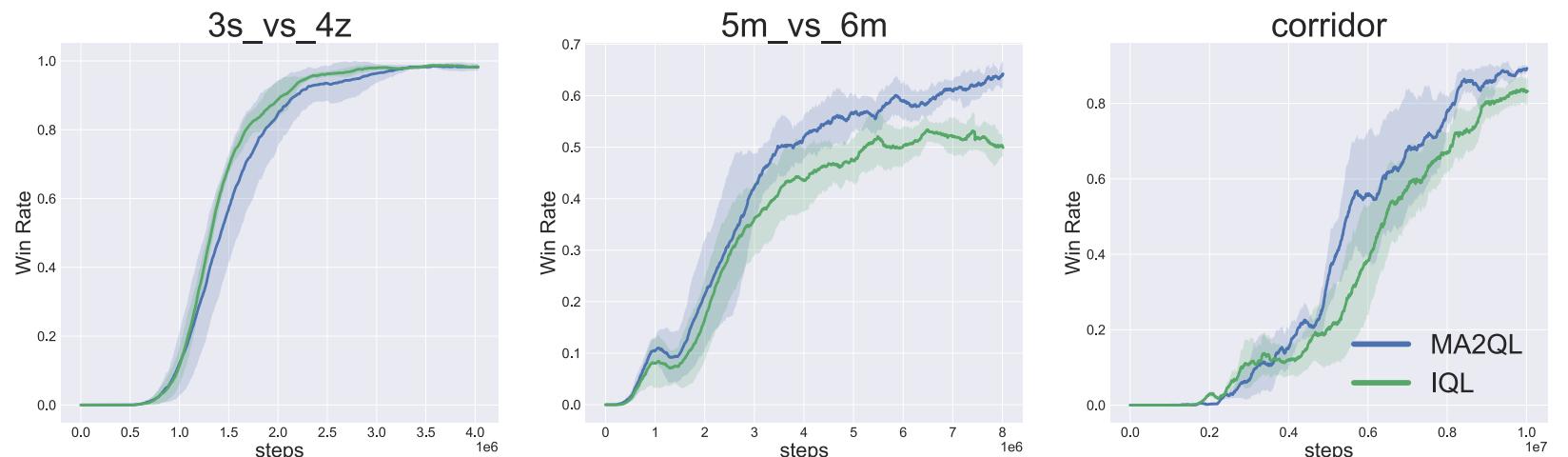


MA2QL vs. IQL

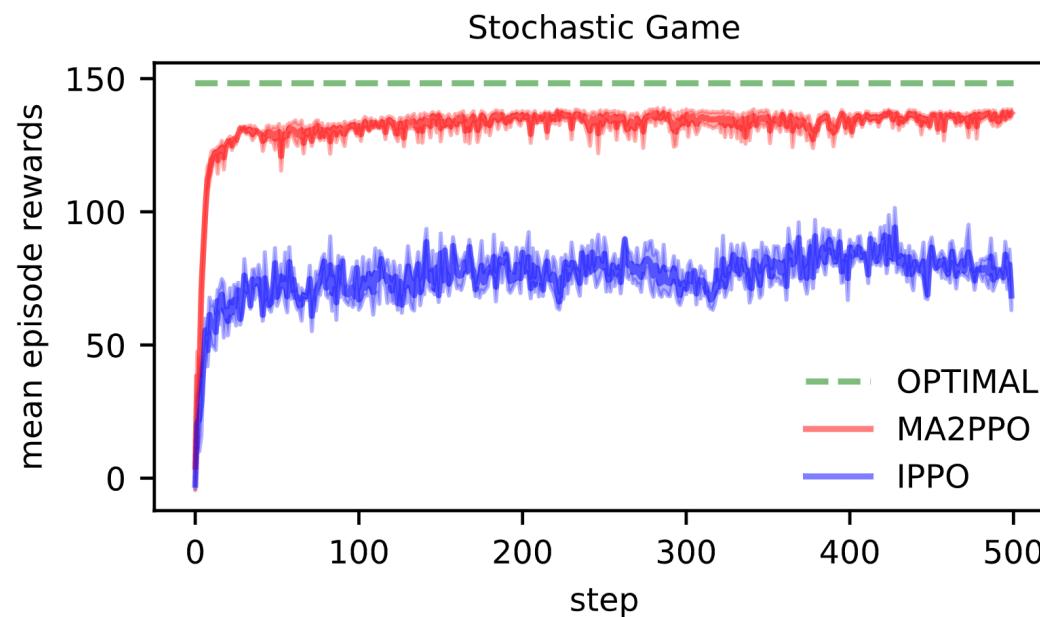
Muti-agent MuJoCo



StarCraft



Alternate PPO vs PPO



MA2QL

MA2QL

```
1: repeat
2:   for  $i \leftarrow 1, n$  do
3:     all agents interact in the environment
4:     agent  $i$  updates by Q-learning
5:   end for
6: until terminate
```

✓ guaranteed convergence to Nash equilibrium

✗ Nash equilibrium does not guarantee the optimality

Can we have a performance guarantee?

Revisit $P_i(s' | s, a_i)$

- Is there an optimal transition function P_i^* ?
 - Optimal in terms of **what**?

Suppose there is only one optimal policy $\pi^*(s) = \arg \max_a Q^*(s, a) \rightarrow \pi^*(s) = \{\pi_1^*, \dots, \pi_N^*\}$



$$\pi_{-i}^*(s, a_i) = \arg \max_{a_{-i}} Q^*(s, a_i, a_{-i})$$



$$P_i^*(s' | s, a_i) = P(s' | s, a_i, \pi_{-i}^*(s, a_i))$$



$$\max_{a_i} Q_i^*(s, a_i) = \max_{a_i} \max_{a_{-i}} Q^*(s, a_i, a_{-i}) = Q^*(s, \pi^*(s))$$

The global optimum under P_i^*

Can we construct P_i^* ?

- Yes, in deterministic environments!
- QSS learning $Q(s, a) \rightarrow Q(s, s')$

$$Q_i^{\text{ss}}(s, s') = r + \gamma \max_{s'' \in \mathcal{N}(s')} Q_i^{\text{ss}}(s', s'') \rightarrow \max_{s'} Q_i^{\text{ss}}(s, s') = \max_a Q^*(s, a)$$

$$s'^* = \arg \max_{s' \in \mathcal{N}(s, a_i)} Q_i^{\text{ss}}(s, s') \rightarrow P_i^*$$

In deterministic environments, all agents will converge to ***the optimal policies***, if each agent i performs Q-learning on the transition function $s'^* = \arg \max_{s' \in \mathcal{N}(s, a_i)} Q_i^{\text{ss}}(s, s')$.

Ideal independent Q-learning (I2Q)

- Each agent has three modules

① $f(s, a_i)$ approximates $P_i^\star(s' | s, a_i)$ by maximizing $\mathbb{E}_{(s, a_i, s') \sim \mathcal{D}_i} [\lambda Q_i^{\text{ss}}(s, f_i(s, a_i)) - (f_i(s, a_i) - s')^2]$

② $Q(s, s')$ minimizes $\mathbb{E}_{(s, a_i, s', r) \sim \mathcal{D}_i} \left[\left(Q_i^{\text{ss}}(s, s') - r - \gamma \bar{Q}_i^{\text{ss}}(s', f_i(s', a_i'^*)) \right)^2 \right]$, $a_i'^* = \arg \max_{a_i'} Q_i(s', a_i')$

③ $Q(s, a_i)$ minimizes $\mathbb{E}_{(s, a_i, r) \sim \mathcal{D}_i} \left[\left(Q_i(s, a_i) - r - \gamma \max_{a_i'} \bar{Q}_i(f_i(s, a_i), a_i') \right)^2 \right]$

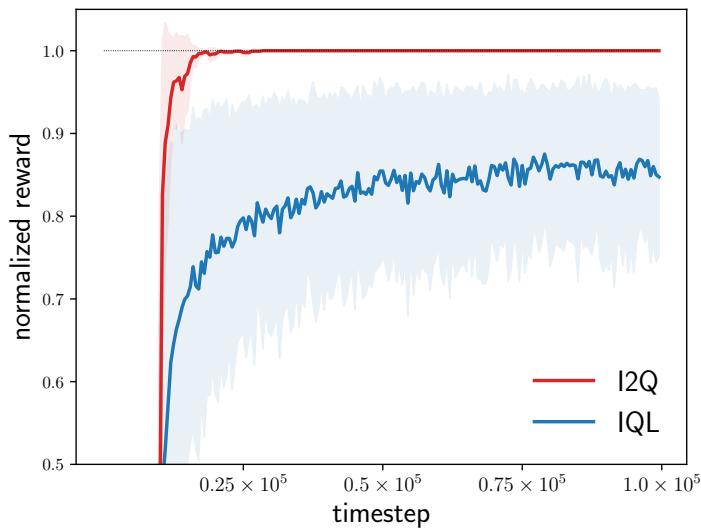
Learning algorithm

Algorithm 1. I2Q for each agent i

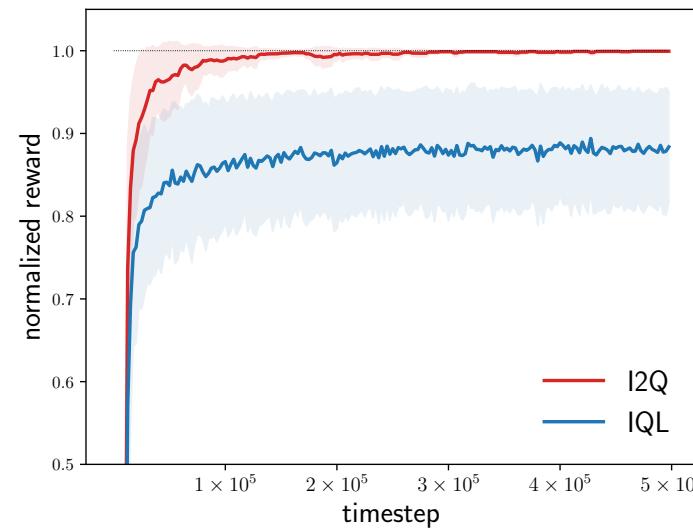
- 1: Initialize transition model f_i , Q-networks Q_i and Q_i^{ss} , and the target networks \bar{Q}_i and \bar{Q}_i^{ss} .
 - 2: Initialize the replay buffer \mathcal{D}_i .
 - 3: **for** $t = 1, \dots, \text{max_iteration}$ **do**
 - 4: All agents interact in the environment and store experiences (s, a_i, s', r) in replay buffer \mathcal{D}_i
 - 5: Sample a mini-batch from \mathcal{D}_i
 - 6: Update f_i by maximizing ①
 - 7: Update Q_i^{ss} by minimizing ②
 - 8: Update Q_i by minimizing ③
 - 9: Update the target networks \bar{Q}_i and \bar{Q}_i^{ss}
 - 10: **end for**
-

I2Q

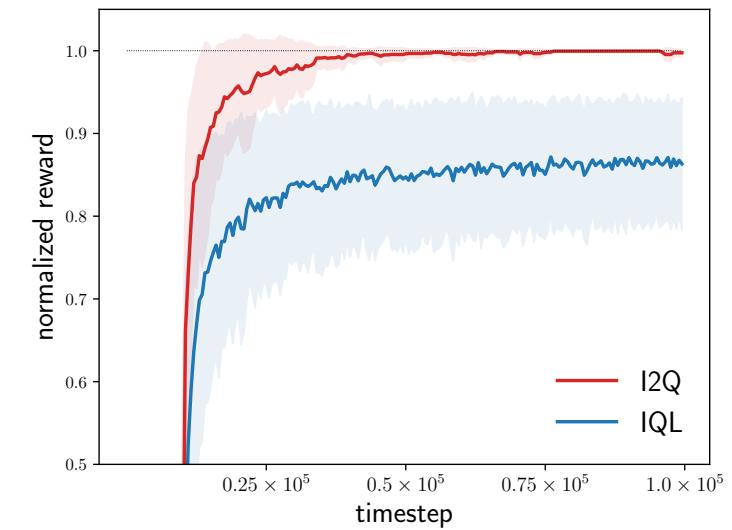
- L -stage matrix game (deterministic tabular case)



$L = 3$



$L = 4$

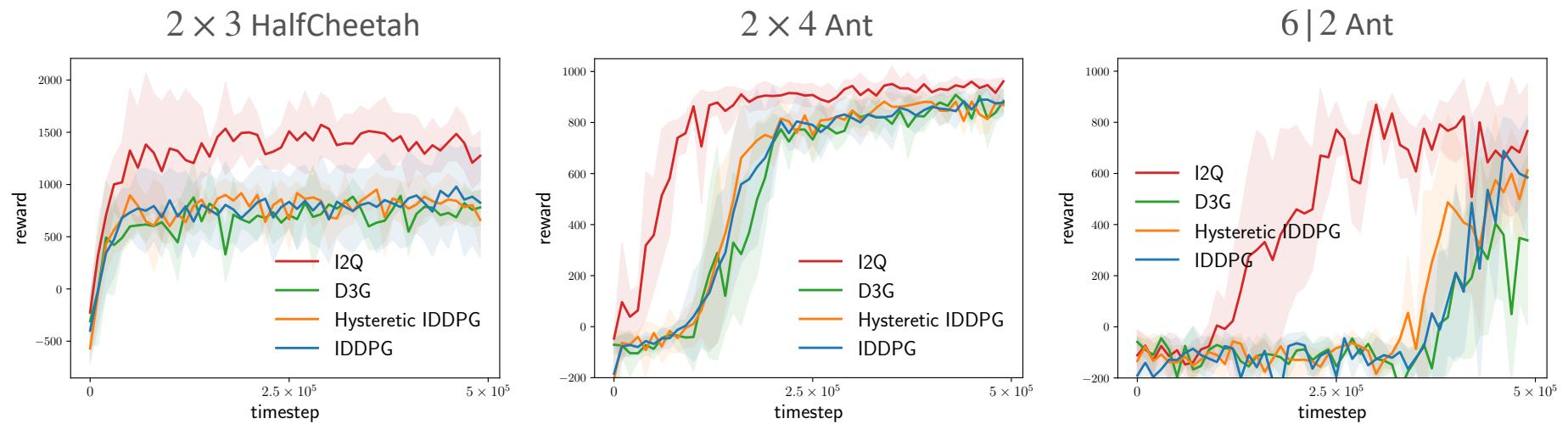


$L = 5$

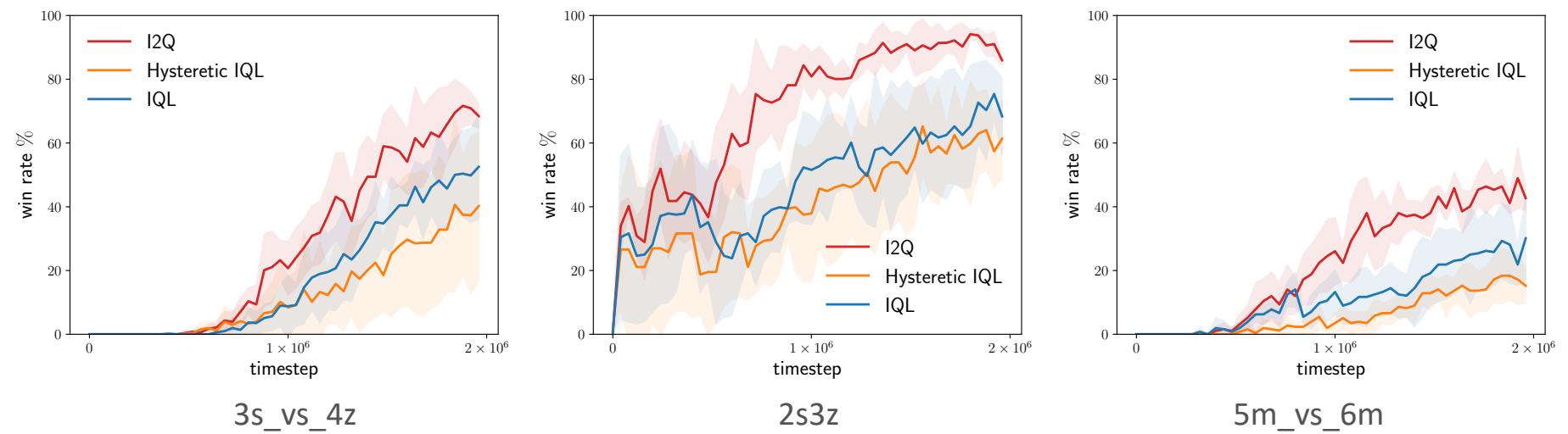
✓ converged to the global optimum

I2Q vs. IQL

Muti-agent MuJoCo
partially observable



StarCraft
stochastic environment



Fully decentralized Q-learning

- Remarks

- ✓ MA2QL guarantees Nash equilibrium for learning Q-tables
- ✓ I2Q guarantees the global optimum in deterministic environments for learning Q-tables
- ✗ Both are not for learning Q-networks

Can we have a performance guarantee even using neural networks?

Outline

- 1. Fully decentralized Q-learning**
- 2. Fully decentralized policy optimization**

- ✓ Su et al., A Minimalist Approach to Decentralized Multi-Agent Reinforcement Learning, 2022
- ✓ Jiang and Lu, I2Q: A fully decentralized Q-learning algorithm, 2022
- ✓ Su and Lu, Fully decentralized policy optimization, 2022
- ✓ Luo, Jiang, and Lu, Model-based decentralized policy optimization, 2022

Policy optimization?

- Independent PPO (IPPO) shows strong performance in several benchmarks
 - de Witt et al., 2020
 - Papoudakis et al., 2021
- **Does IPPO have a convergence guarantee?**

Does IPPO have a convergence guarantee?

- Fully decentralized critic

$$Q_{\pi^{-i}}^{\pi^i}(s, a_i) = r_{\pi^{-i}}(s, a_i) + \gamma \mathbb{E}_{a_{-i} \sim \pi^{-i}, s' \sim P(\cdot | s, a_i, a_{-i}), a'_i \sim \pi^i} [Q_{\pi^{-i}}^{\pi^i}(s', a'_i)]$$



$$Q_{\pi^{-i}}^{\pi^i}(s, a_i) = \mathbb{E}_{\pi^{-i}}[Q^{\boldsymbol{\pi}}(s, a_i, a_{-i})]$$

$$V_{\pi^{-i}}^{\pi^i}(s) = \mathbb{E}_{\pi^{-i}}[V^{\boldsymbol{\pi}}(s)] = V^{\boldsymbol{\pi}}(s)$$

What does IPPO optimize?

- Joint TRPO objective

$$J(\boldsymbol{\pi}_{\text{new}}) - J(\boldsymbol{\pi}_{\text{old}}) \geq \mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^{\text{joint}}(\boldsymbol{\pi}_{\text{new}}) - C \cdot D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}} \| \boldsymbol{\pi}_{\text{new}})$$

where $\mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^{\text{joint}}(\boldsymbol{\pi}_{\text{new}}) = \sum_s \rho_{\text{old}}(s) \sum_a \boldsymbol{\pi}_{\text{new}}(\mathbf{a} | s) A_{\text{old}}(s, \mathbf{a})$

- IPPO objective

$$\mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^i(\boldsymbol{\pi}_{\text{new}}^i) = \sum_s \rho_{\text{old}}(s) \sum_{a_i} \boldsymbol{\pi}_{\text{new}}^i(a_i | s) A_{\text{old}}^i(s, a_i)$$

where $A_{\text{old}}^i(s, a_i) = \mathbb{E}_{\boldsymbol{\pi}_{\text{old}}^{-i}}[A_{\text{old}}(s, a_i, a_{-i})]$

X optimizing this objective does not necessarily guarantee the improvement of the joint policy

What should we optimize?

- This is the only objective for fully decentralized learning

$$\mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^i(\boldsymbol{\pi}_{\text{new}}^i) = \sum_s \rho_{\text{old}}(s) \sum_{a_i} \pi_{\text{new}}^i(a_i | s) A_{\text{old}}^i(s, a_i)$$

- A novel lower bound

$$J(\boldsymbol{\pi}_{\text{new}}) - J(\boldsymbol{\pi}_{\text{old}}) \geq \mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^{\text{joint}}(\boldsymbol{\pi}_{\text{new}}) - C \cdot D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}} \| \boldsymbol{\pi}_{\text{new}})$$



$$J(\boldsymbol{\pi}_{\text{new}}) - J(\boldsymbol{\pi}_{\text{old}}) \geq \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^i(\boldsymbol{\pi}_{\text{new}}^i) - \tilde{M} \cdot \sum_{i=1}^N \sqrt{D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}}^i \| \boldsymbol{\pi}_{\text{new}}^i)} - C \cdot \sum_{i=1}^N D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}}^i \| \boldsymbol{\pi}_{\text{new}}^i)$$

Fully decentralized surrogate

$$J(\boldsymbol{\pi}_{\text{new}}) - J(\boldsymbol{\pi}_{\text{old}}) \geq \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^i(\boldsymbol{\pi}_{\text{new}}^i) - \tilde{M} \cdot \sum_{i=1}^N \sqrt{D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}}^i \| \boldsymbol{\pi}_{\text{new}}^i)} - C \cdot \sum_{i=1}^N D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}}^i \| \boldsymbol{\pi}_{\text{new}}^i)$$



$$\boldsymbol{\pi}_{\text{new}}^i = \arg \max_{\boldsymbol{\pi}^i} \left(\frac{1}{N} \mathcal{L}_{\boldsymbol{\pi}_{\text{old}}}^i(\boldsymbol{\pi}^i) - \tilde{M} \cdot \sqrt{D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}}^i \| \boldsymbol{\pi}^i)} - C \cdot D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\text{old}}^i \| \boldsymbol{\pi}^i) \right)$$



- ✓ the joint policy of agents improves monotonically and converges to *suboptimum*
- ✗ the step size of the policy update will be small, similar to TRPO

Fully decentralized policy optimization (DPO)

- The practical algorithm

$$\pi_{\text{new}}^i = \arg \max_{\pi^i} \left(\frac{1}{N} \mathcal{L}_{\pi_{\text{old}}}^i(\pi^i) - \tilde{M} \cdot \sqrt{D_{\text{KL}}^{\max}(\pi_{\text{old}}^i \| \pi^i)} - C \cdot D_{\text{KL}}^{\max}(\pi_{\text{old}}^i \| \pi^i) \right)$$

 In practice

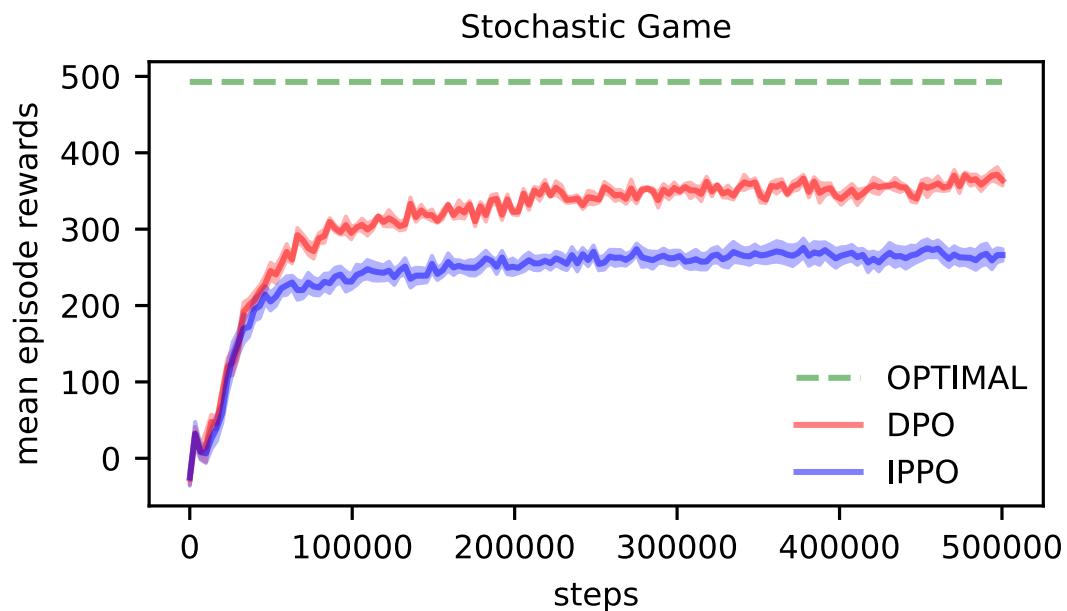
$$\pi_{\text{new}}^i = \arg \max_{\pi^i} \left(\frac{1}{N} \mathcal{L}_{\pi_{\text{old}}}^i(\pi^i) - \beta_1^i \sqrt{D_{\text{KL}}^{\text{avg}}(\pi_{\text{old}}^i \| \pi^i)} - \beta_2^i D_{\text{KL}}^{\text{avg}}(\pi_{\text{old}}^i \| \pi^i) \right)$$

If $D_{\text{KL}}^{\text{avg}}(\pi_{\text{old}}^i \| \pi_{\text{new}}^i) > d_{\text{target}} * \delta$, then $\beta_j^i \leftarrow \beta_j^i * \omega \quad \forall j \in \{1,2\}$

If $D_{\text{KL}}^{\text{avg}}(\pi_{\text{old}}^i \| \pi_{\text{new}}^i) < d_{\text{target}} / \delta$, then $\beta_j^i \leftarrow \beta_j^i / \omega \quad \forall j \in \{1,2\}$

DPO vs. IPPO

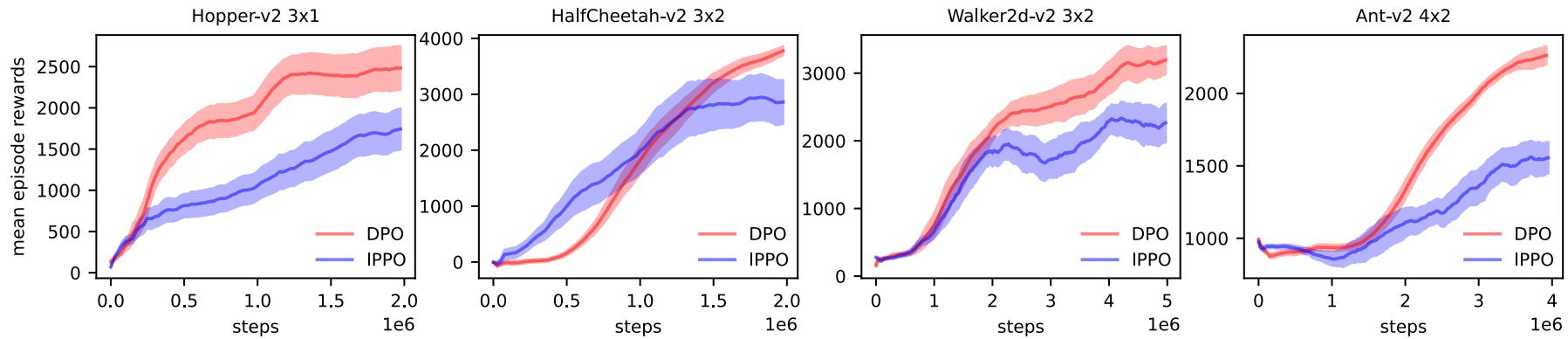
- Tabular case



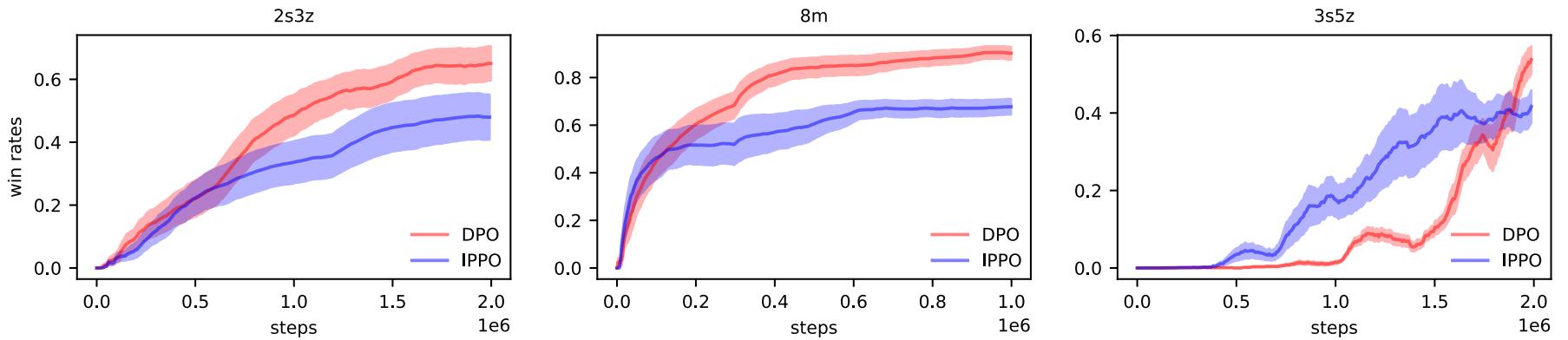
✓ better sub-optimum empirically
✗ still a large gap to the optimum

DPO vs. IPPO

Muti-agent MuJoCo



StarCraft



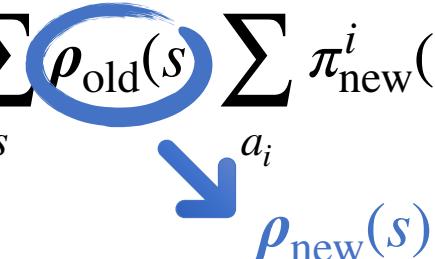
DPO

- Remarks

✓ DPO guarantees monotonic improvement and sub-optimum ?

✗ The assumption of stationary discounted state distribution *may not* hold in multi-agent settings

$$\mathcal{L}_{\pi_{\text{old}}}^i(\pi_{\text{new}}^i) = \sum_s \rho_{\text{old}}(s) \sum_{a_i} \pi_{\text{new}}^i(a_i | s) A_{\text{old}}^i(s, a_i)$$

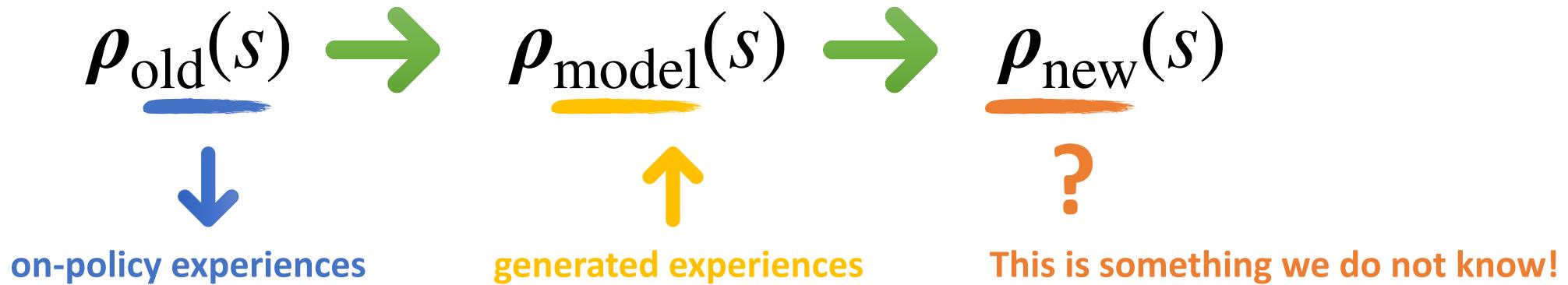


$$\rho_{\text{new}}(s)$$

Can we bridge the gap between $\rho_{\text{old}}(s)$ and $\rho_{\text{new}}(s)$?



Can we bridge the gap between $\rho_{\text{old}}(s)$ and $\rho_{\text{new}}(s)$?



Bridging the gap by learning from experiences
generated by the environment model?

What model can we have?

- Fully decentralized model with latent variable

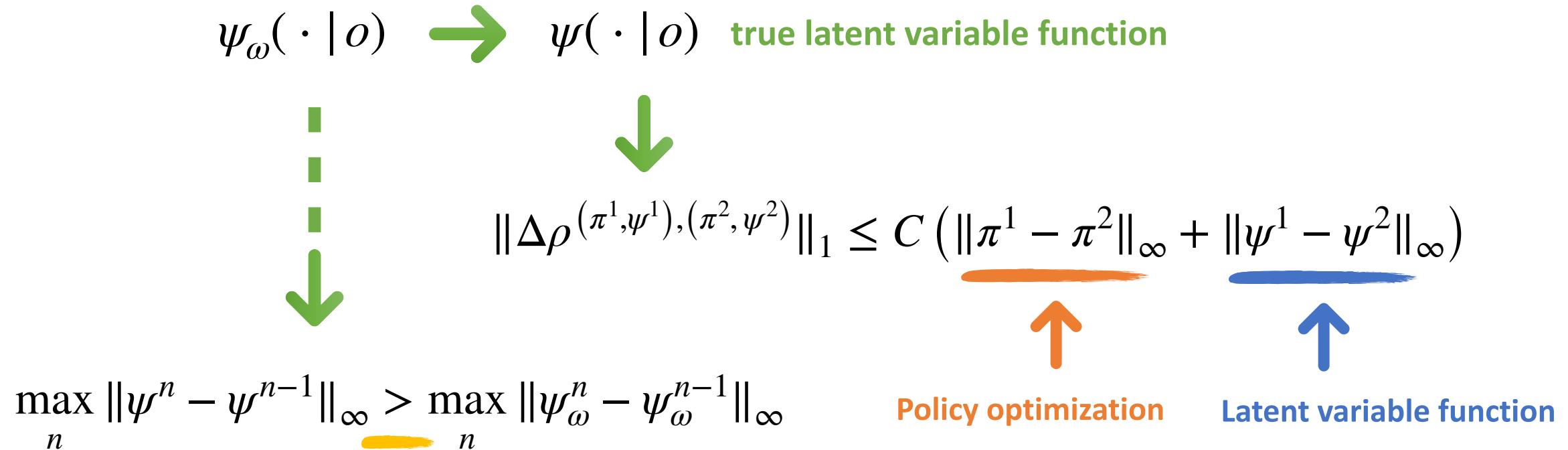
$$P_i(o'_i | o_i, a_i, z_i)$$

$$R(o_i, a_i, z_i)$$

z_i captures the inaccessible information of state and other agents' policies

$\psi_\omega(\cdot | o) \rightarrow \psi(\cdot | o)$ true latent variable function

Does such a simple model work?



More stable policy optimization on model

Can we guarantee policy improvement?

$$\left| \eta(\pi, \psi^{n+1}) - \eta^{model}(\pi, \psi_\omega^n) \right| \leq C(\epsilon_\theta, \epsilon_\pi, \epsilon_\omega, \epsilon_\psi)$$

- model learning** → ϵ_θ : model error
- policy optimization** → ϵ_π : policy divergence
- model learning** → ϵ_ω : learned latent variable function error
- ϵ_ψ : true latent variable function divergence

?

Can we guarantee policy improvement?

$$\hat{\epsilon}_\omega \triangleq \max_o D_{\text{TV}}(\psi^{n+1} \| \psi_\omega^n)$$



$$\left| \eta(\pi, \psi^{n+1}) - \eta^{\text{model}}(\pi, \psi_\omega^n) \right| \leq C(\epsilon_\theta, \epsilon_\pi, \hat{\epsilon}_\omega)$$



Constraining $\hat{\epsilon}_\omega$ controls the return bound

How to constrain $\hat{\epsilon}_\omega$?

- Fully decentralized model with latent variable prediction

$$\hat{\epsilon}_\omega \triangleq \max_o D_{\text{TV}}(\psi^{n+1} \| \psi_\omega^n)$$



$$\max_o D_{\text{TV}}(\psi^l(o) \| f(\psi^1(o), \dots, \psi^{l-1}(o)))$$

Model-based decentralized policy optimization

- ① Model learning → $\max_{\theta, \omega_l, \phi} \sum_{j=1}^l \left(\mathbb{E}_{(o, a, o', r) \sim \mathcal{D}^j, z \sim \psi_{\omega_j}(o)} \left[P_\theta(o' | o, a, z) - (R_\phi(o, a, z) - r)^2 \right] \right)$
- ② Latent variable prediction → $\max_{\zeta} \mathbb{E}_{o \sim \mathcal{D}^l, z^1 \sim \psi_{\omega_1}(o), \dots, z^l \sim \psi_{\omega_l}(o)} \left[f_\zeta(z^l | z^1, \dots, z^{l-1}) \right]$
- ③ Model rollout → $z \sim f_\zeta \left(\cdot | z^2 \sim \psi_{\omega_2}(o), \dots, z^l \sim \psi_{\omega_l}(o) \right), \hat{o}' \sim P_\theta(o' | o_a, z), \hat{r} = R_\phi(o, a, z)$
- ④ Policy optimization → DPO/IPPO on branched model rollout

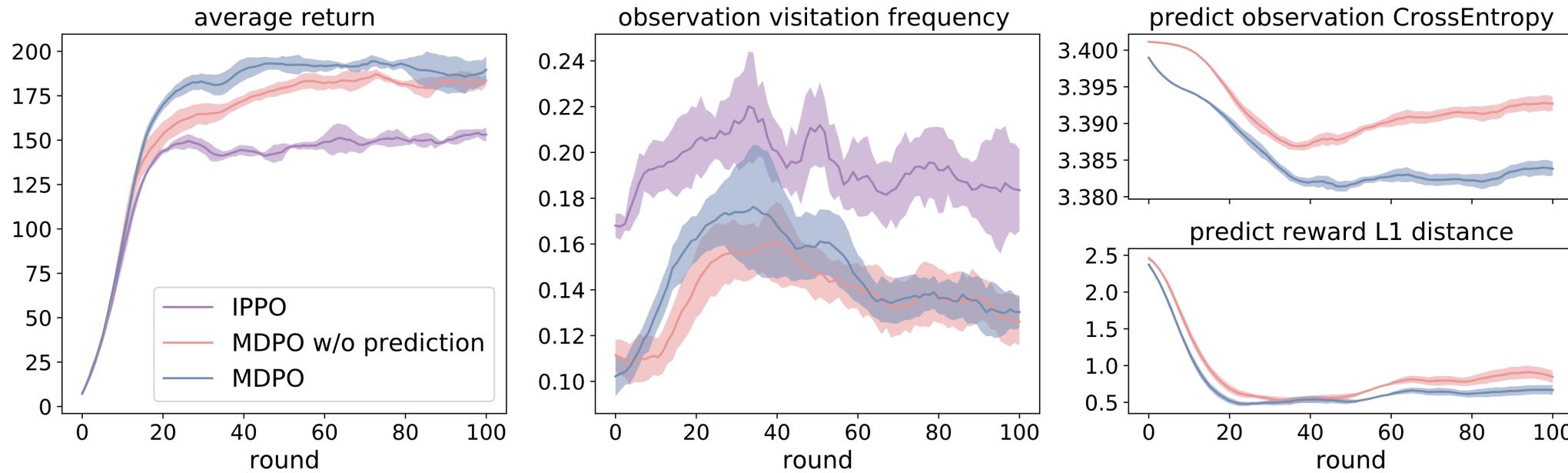
Learning algorithm

Algorithm 1. MDPO

- 1: **Initiate** $\mathcal{D}_{env} = \{\mathcal{D}^1, \dots, \mathcal{D}^l\}, \pi, P_\theta, R_\phi, \Psi = \{\psi_{\omega_1}, \dots, \psi_{\omega_l}\}, f_\zeta$.
- 2: **repeat**
- 3: policy rollout in environment and obtain \mathcal{D}^l
- 4: optimize P_θ, R_ϕ and ψ_{ω_l} on \mathcal{D}_{env} with ①
- 5: optimize prediction function f_ζ on \mathcal{D}^l with ②
- 6: obtain branched model rollout $\mathcal{D}_{rollout}$ based on \mathcal{D}^l using $P_\theta, R_\phi, \pi, \Psi$, and f_ζ with ③
- 7: optimize policy π using $\mathcal{D}_{rollout}$ by IPPO/DPO ④
- 8: **for** $j \leftarrow 1, \dots, l - 1$ **do**
- 9: $\mathcal{D}^j \leftarrow \mathcal{D}^{j+1}, \psi_{\omega_j} \leftarrow \psi_{\omega_{j+1}}$
- 10: **end for**
- 11: **until** terminate

MDPO, MDPO w/o Prediction, and IPPO

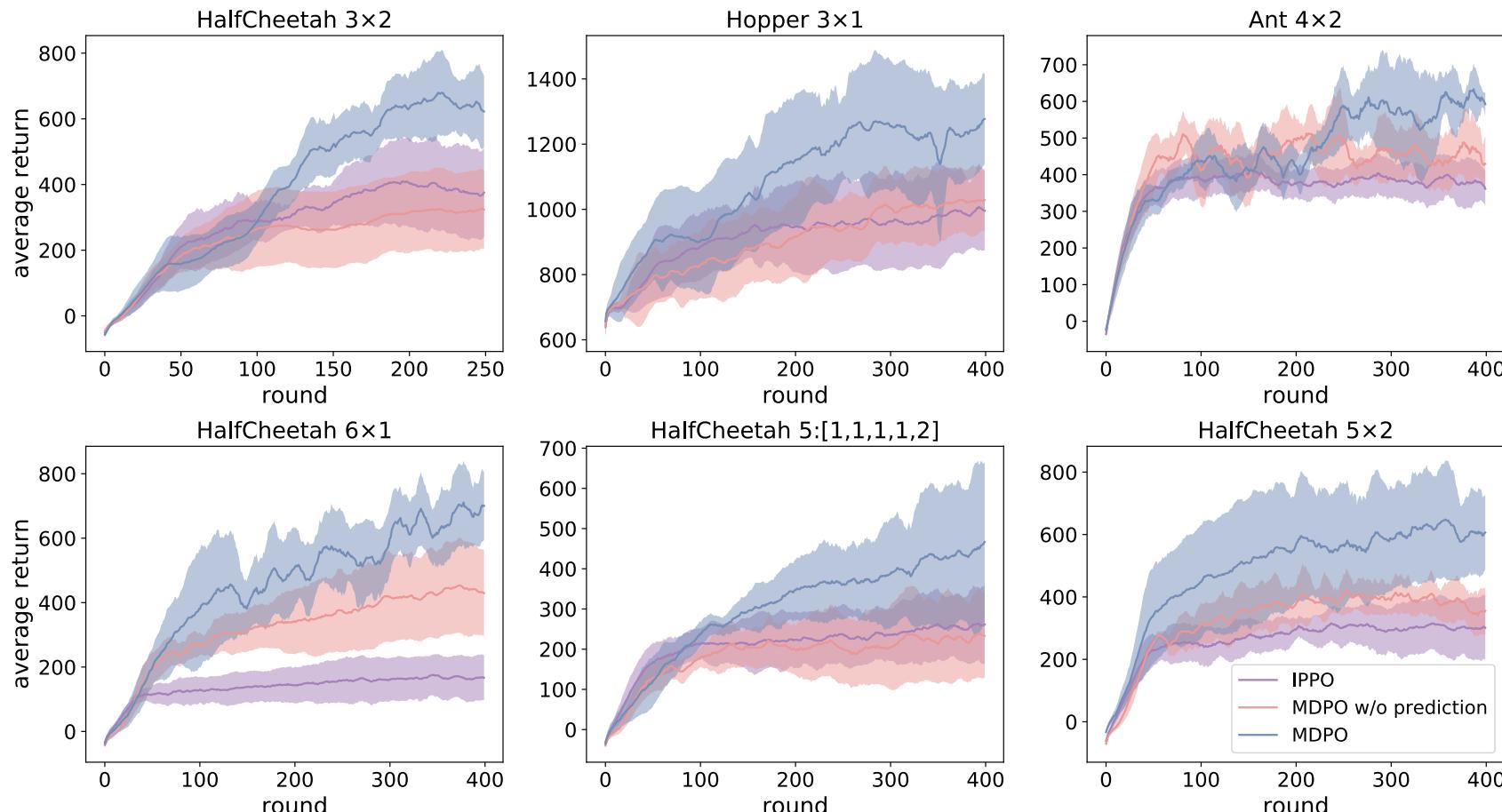
- A cooperative stochastic game



- ✓ MDPO indeed induces more stable visitation frequency
- ✓ Latent variable prediction better controls $\hat{\epsilon}_\omega$

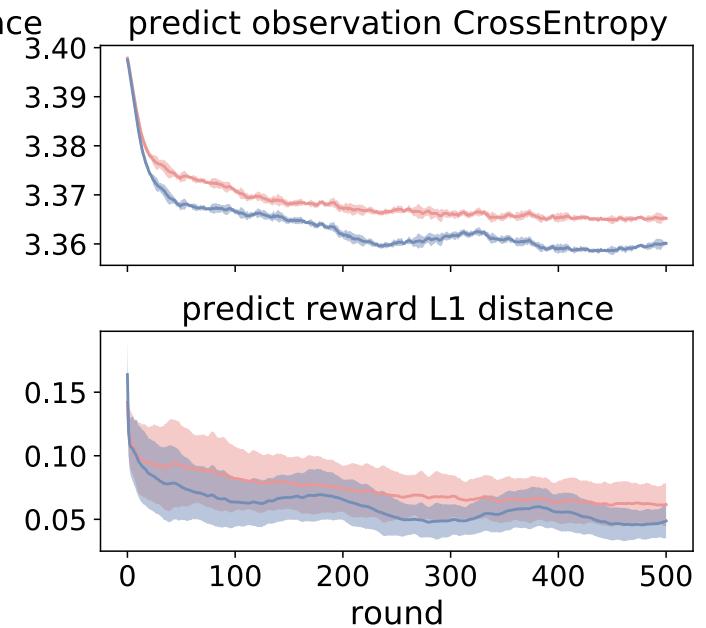
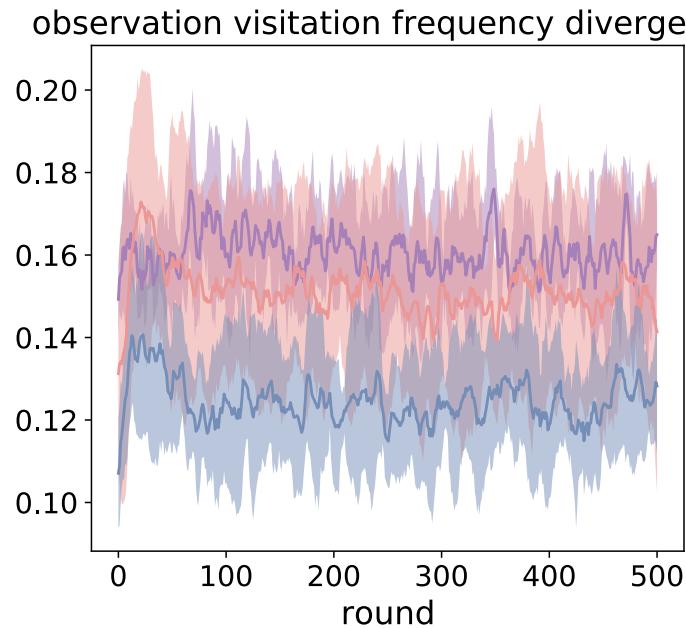
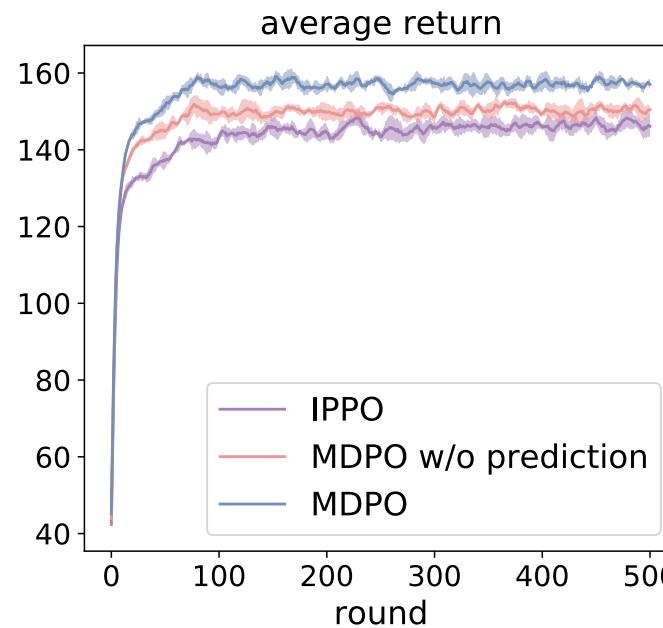
MDPO, MDPO w/o Prediction, and IPPO

- Multi-agent MuJoCo



MDPO, MDPO w/o Prediction, and IPPO

- Does MDPO help in single-agent non-stationary environments?



✓ MDPO also helps in single-agent non-stationary environments

Outline

1. Fully decentralized Q-learning
2. Fully decentralized policy optimization
3. ???!!!

- ✓ Su et al., A Minimalist Approach to Decentralized Multi-Agent Reinforcement Learning, 2022
- ✓ Jiang and Lu, I2Q: A fully decentralized Q-learning algorithm, 2022
- ✓ Su and Lu, Fully decentralized policy optimization, 2022
- ✓ Luo, Jiang, and Lu, Model-based decentralized policy optimization, 2022

Takeaways

- ✓ Fully decentralized MARL is widely applicable
- ✓ Fully decentralized MARL can be much better than IQL and IPPO

- ✗ These methods are preliminary
- ✗ More efforts are needed

Fully decentralized learning can be a unified paradigm for RL and cooperative MARL!

That is all!
Questions?

zongqing.lu@pku.edu.cn

