

---

# Multi-Agent Reinforcement Learning is A Sequence Modeling Problem

---

Muning Wen<sup>1,2</sup>, Jakub Grudzien Kuba<sup>3</sup>, Runji Lin<sup>4</sup>,  
Weinan Zhang<sup>1</sup>, Ying Wen<sup>1</sup>, Jun Wang<sup>2,5</sup>, Yaodong Yang<sup>6</sup>

<sup>1</sup>Shanghai Jiao Tong University, <sup>2</sup>Digital Brain Lab,

<sup>3</sup>University of Oxford, <sup>4</sup>Institute of Automation, Chinese Academy of Science,

<sup>5</sup>University College London, <sup>6</sup>Institute for AI, Peking University

Prequels:

- Settling the Variance of Multi-Agent Policy Gradient ([NeurIPS 2021](#)).
- Trust Region Policy Optimisation in Multi-Agent Reinforcement Learning ([ICLR 2022](#)).

# Motivation of Multi-Agent Transformer

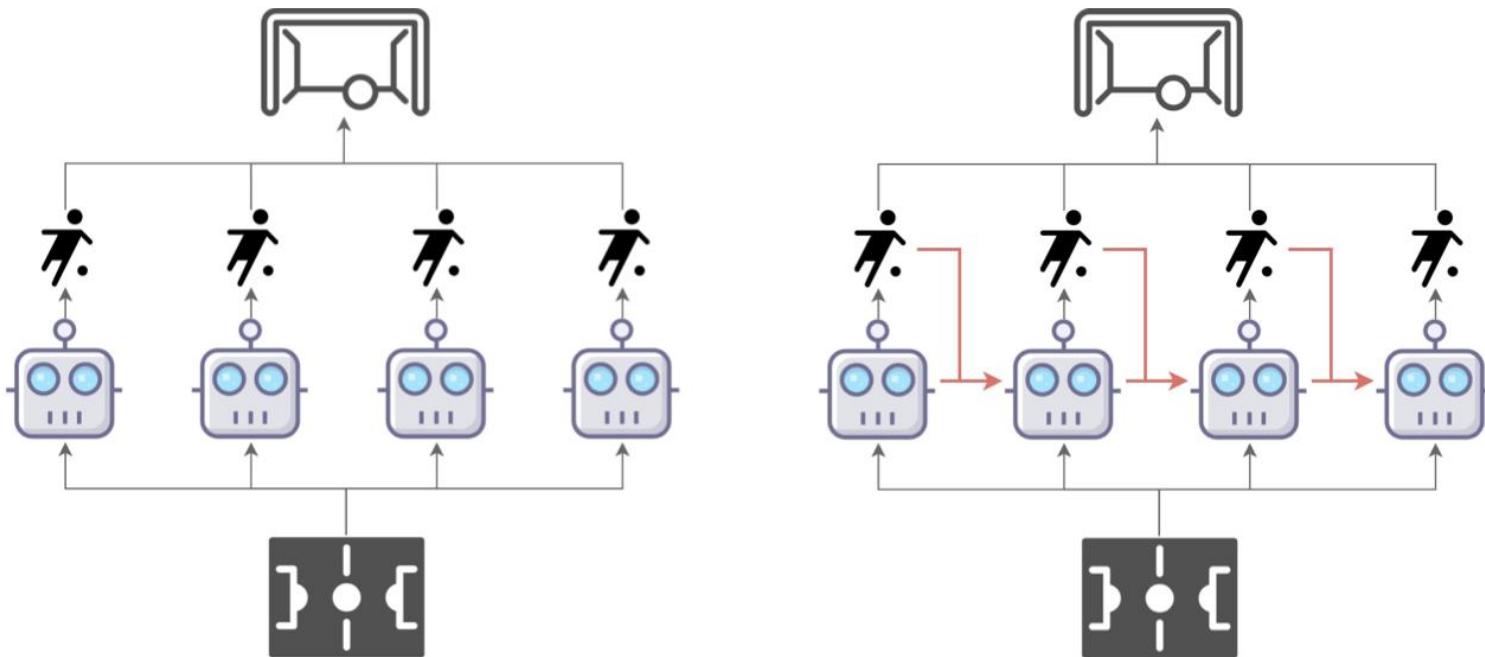


- How to abstract multi-agent decision making as a sequence modeling problem and leverage the prosperous development of the sequence models
- Building the bridge between MARL and sequence models so that the modeling power of modern sequence models, e.g. the Transformer, can be unleashed for MARL.
- Guarantee the monotonic performance improvement of joint policy with acceptable sample efficiency and training speed.
- More information could be found in our project website (<https://sites.google.com/view/multi-agent-transformer>) or github(<https://github.com/PKU-MARL/Multi-Agent-Transformer>).

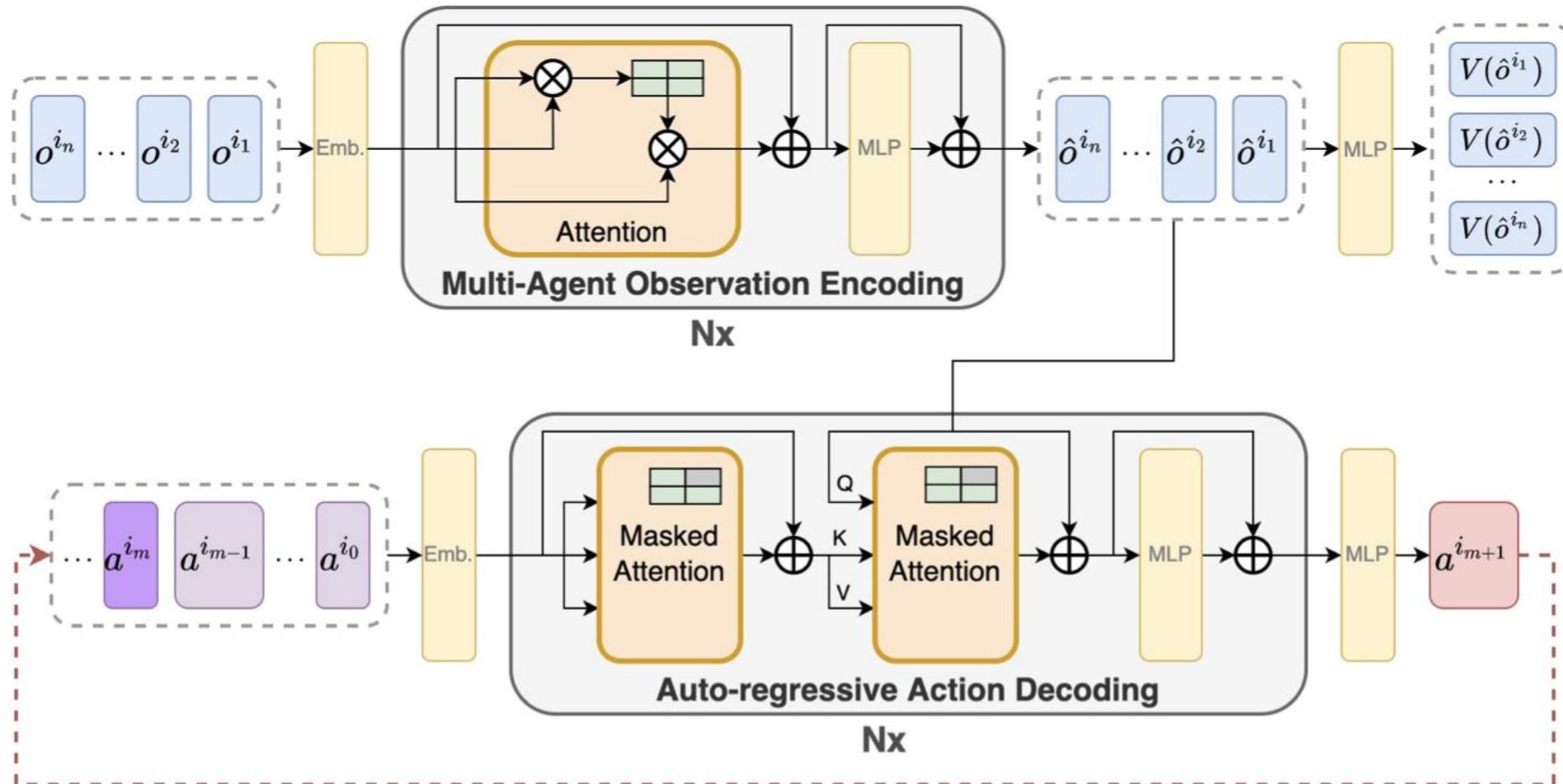
# The Multi-Agent Sequential Decision Paradigm

**Theorem 1** (Multi-Agent Advantage Decomposition [13]). *Let  $i_{1:n}$  be a permutation of agents. Then, for any joint observation  $\mathbf{o} \in \mathcal{O}$  and joint action  $\mathbf{a} = \mathbf{a}^{i_{1:n}} \in \mathcal{A}$ , the following equation always holds with no further assumption needed,*

$$A_{\pi}^{i_{1:n}}(\mathbf{o}, \mathbf{a}^{i_{1:n}}) = \sum_{m=1}^n A_{\pi}^{i_m}(\mathbf{o}, \mathbf{a}^{i_{1:m-1}}, a^{i_m}).$$



# Multi-Agent Transformer

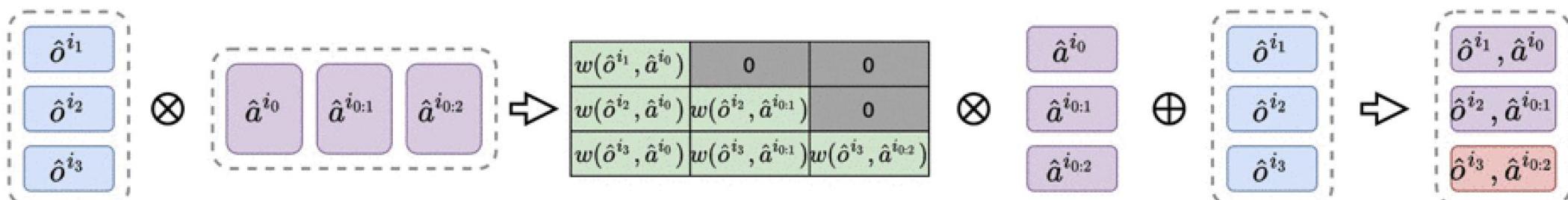


$$L_{\text{Encoder}}(\phi) = \frac{1}{Tn} \sum_{m=1}^n \sum_{t=0}^{T-1} \left[ R(\mathbf{o}_t, \mathbf{a}_t) + \gamma V_{\bar{\phi}}(\hat{\mathbf{o}}_{t+1}^{i_m}) - V_{\phi}(\hat{\mathbf{o}}_t^{i_m}) \right]^2, \quad (4)$$

$$L_{\text{Decoder}}(\theta) = -\frac{1}{Tn} \sum_{m=1}^n \sum_{t=0}^{T-1} \min \left( \mathbf{r}_t^{i_m}(\theta) \hat{A}_t, \text{clip}(\mathbf{r}_t^{i_m}(\theta), 1 \pm \epsilon) \hat{A}_t \right), \quad (5)$$

$$\mathbf{r}_t^{i_m}(\theta) = \frac{\pi_{\theta}^{i_m}(\mathbf{a}_t^{i_m} | \hat{\mathbf{o}}_t^{i_{1:n}}, \hat{\mathbf{a}}_t^{i_{1:m-1}})}{\pi_{\theta_{\text{old}}}^{i_m}(\mathbf{a}_t^{i_m} | \hat{\mathbf{o}}_t^{i_{1:n}}, \hat{\mathbf{a}}_t^{i_{1:m-1}})},$$

An example of how the attention work in decoder.



# Single-task Experiments

---

Table 1: Performance evaluations of win rate and standard deviation on the SMAC benchmark.

Task	Difficulty	MAT	MAT-Dec	MAPPO	HAPPO	QMIX	Steps
1c3s5z	Easy	<b>100.0</b> <sub>(2.4)</sub>	<b>100.0</b> <sub>(0.4)</sub>	<b>100.0</b> <sub>(2.2)</sub>	97.5 <sub>(1.8)</sub>	96.9 <sub>(1.5)</sub>	2e6
MMM	Easy	<b>100.0</b> <sub>(2.2)</sub>	98.1 <sub>(2.1)</sub>	95.6 <sub>(4.5)</sub>	81.2 <sub>(22.9)</sub>	91.2 <sub>(3.2)</sub>	2e6
2c vs 64zg	Hard	<b>100.0</b> <sub>(1.3)</sub>	95.9 <sub>(2.3)</sub>	<b>100.0</b> <sub>(2.7)</sub>	90.0 <sub>(4.8)</sub>	90.3 <sub>(4.0)</sub>	5e6
3s vs 5z	Hard	<b>100.0</b> <sub>(1.7)</sub>	<b>100.0</b> <sub>(1.3)</sub>	<b>100.0</b> <sub>(2.5)</sub>	91.9 <sub>(5.3)</sub>	92.3 <sub>(4.4)</sub>	5e6
3s5z	Hard	<b>100.0</b> <sub>(1.9)</sub>	<b>100.0</b> <sub>(3.3)</sub>	72.5 <sub>(26.5)</sub>	90.0 <sub>(3.5)</sub>	84.3 <sub>(5.4)</sub>	3e6
5m vs 6m	Hard	<b>90.6</b> <sub>(4.4)</sub>	83.1 <sub>(4.6)</sub>	88.2 <sub>(6.2)</sub>	73.8 <sub>(4.4)</sub>	75.8 <sub>(3.7)</sub>	5e6
8m vs 9m	Hard	<b>100.0</b> <sub>(3.1)</sub>	95.0 <sub>(4.6)</sub>	93.8 <sub>(3.5)</sub>	86.2 <sub>(4.4)</sub>	92.6 <sub>(4.0)</sub>	5e6
10m vs 11m	Hard	<b>100.0</b> <sub>(1.4)</sub>	<b>100.0</b> <sub>(2.0)</sub>	96.3 <sub>(5.8)</sub>	77.5 <sub>(9.7)</sub>	95.8 <sub>(6.1)</sub>	5e6
25m	Hard	<b>100.0</b> <sub>(1.3)</sub>	86.9 <sub>(5.6)</sub>	<b>100.0</b> <sub>(2.7)</sub>	0.6 <sub>(0.8)</sub>	90.2 <sub>(9.8)</sub>	2e6
27m vs 30m	Super Hard	<b>100.0</b> <sub>(0.7)</sub>	95.3 <sub>(2.2)</sub>	93.1 <sub>(3.2)</sub>	0.0 <sub>(0.0)</sub>	39.2 <sub>(8.8)</sub>	1e7
MMM2	Super Hard	<b>93.8</b> <sub>(2.6)</sub>	91.2 <sub>(5.3)</sub>	81.8 <sub>(10.1)</sub>	0.3 <sub>(0.4)</sub>	88.3 <sub>(2.4)</sub>	1e7
6h vs 8z	Super Hard	<b>98.8</b> <sub>(1.3)</sub>	93.8 <sub>(4.7)</sub>	88.4 <sub>(5.7)</sub>	0.0 <sub>(0.0)</sub>	9.7 <sub>(3.1)</sub>	1e7
3s5z vs 3s6z	Super Hard	<b>96.5</b> <sub>(1.3)</sub>	85.3 <sub>(7.5)</sub>	84.3 <sub>(19.4)</sub>	82.8 <sub>(21.2)</sub>	68.8 <sub>(21.2)</sub>	2e7

# Single-task Experiments

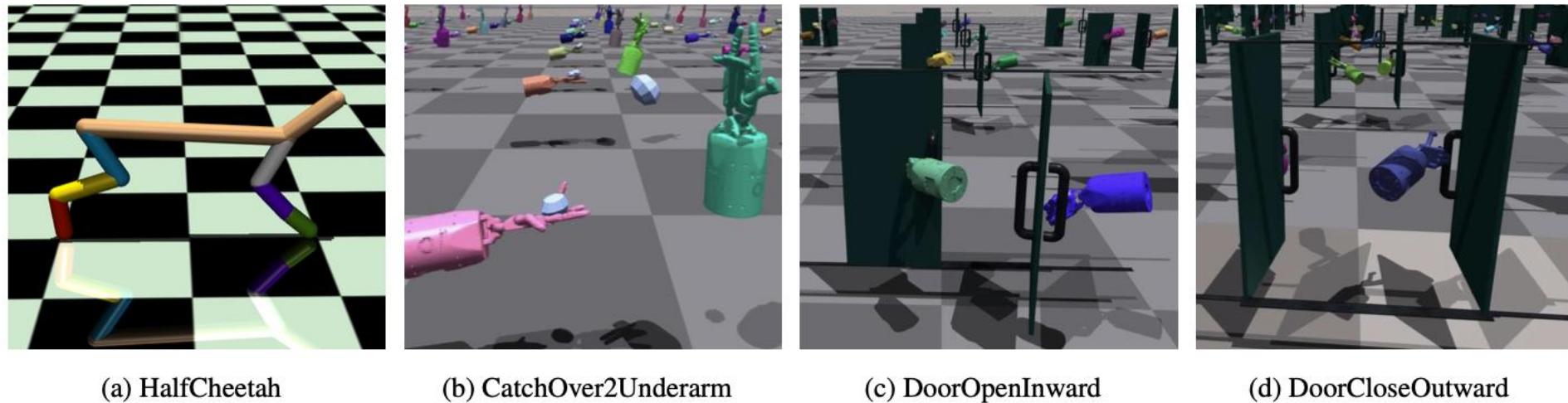


Figure 3: Demonstrations of the Bi-DexHands and the HalfCheetah environments.

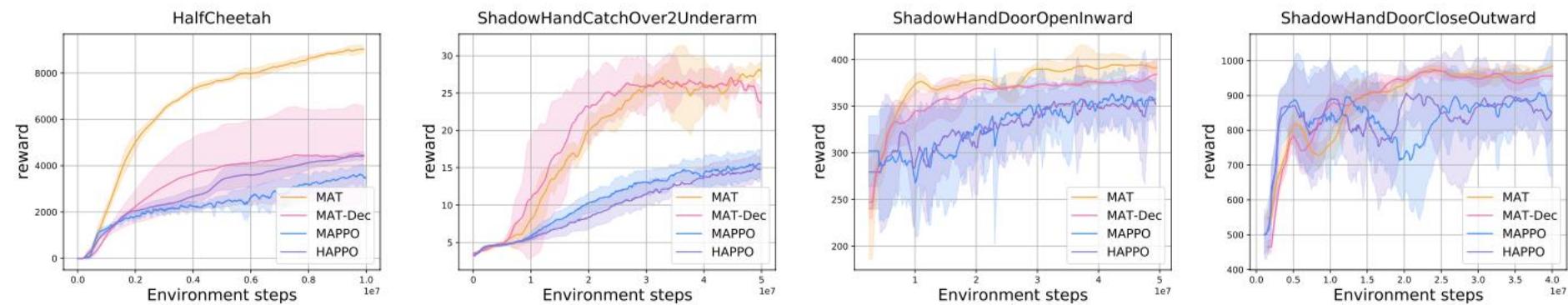


Figure 4: Performance comparisons on the Multi-Agent MuJoCo and the Bi-DexHands benchmarks.

---

# Thank you for your attention!

