



# 强化学习原理及应用 Reinforcement Learning (RL): Theories & Applications

*DCS6289 Spring 2022*

Yucong Zhang (张宇聪)

School of Computer Science and Engineering  
Sun Yat-Sen University

# Lecture 11 : Multi-Agent RL

24<sup>th</sup> May. 2022

- ❑ Learning cooperation
  - ❑ MAPPO
  - ❑ EOI

## □ MAPPO

- On-policy and Off-policy
  - The experiences are collected using the latest learned policy, and then using that experience to improve the policy.
  - Off-policy learning allows the use of older samples(collected using the older policies) in the calculation.
  - On-policy RL algorithms are significantly less sample efficient than off-policy methods.
- Main contributions
  - Multi-Agent PPO(MAPPO), with minimal hyperparameter tuning and without any domain-specific algorithmic changes or architectures, achieves final performances comparable to various off-policy methods
  - MAPPO obtains these strong results while often using a comparable number of samples to many off-policy methods
  - Five implementation and algorithmic factors that govern the practical performance of MAPPO

## □ MAPPO

### ➤ Component

- Following the CTDE structure: following the algorithmic structure of the single-agent PPO algorithm by learning a policy  $\pi_\theta$  and a value function  $V_\phi(s)$  which can take extra global information.
- Generalized Advantage Estimation(GAE)
- Advantage normalization
- Observation normalization
- Gradient clipping
- Value clipping
- Layer normalization
- ReLU activation with orthogonal initialization
- Large batch size
- Five concrete implementation details which are insightful and particularly critical to MAPPO's practical performance: **value normalization, value function inputs, training data usage, policy and value clipping, death masking.**

## □ MAPPO

### ➤ Value Normalization

- To stabilize value learning
- Standardize the targets of the value function by using running estimates of the average and standard deviation of the value targets
- Conclusion: using value normalization never hurts training and often significantly improves the final performance of MAPPO

### ➤ Training Data Usage

- A major trick in PPO is the use of importance sampling to perform off-policy corrections, allowing for sample reuse.
- MAPPO's performance degrades when samples are re-used too often.
- Conclusion: avoid using too many training epochs and do not split data into mini-batches.

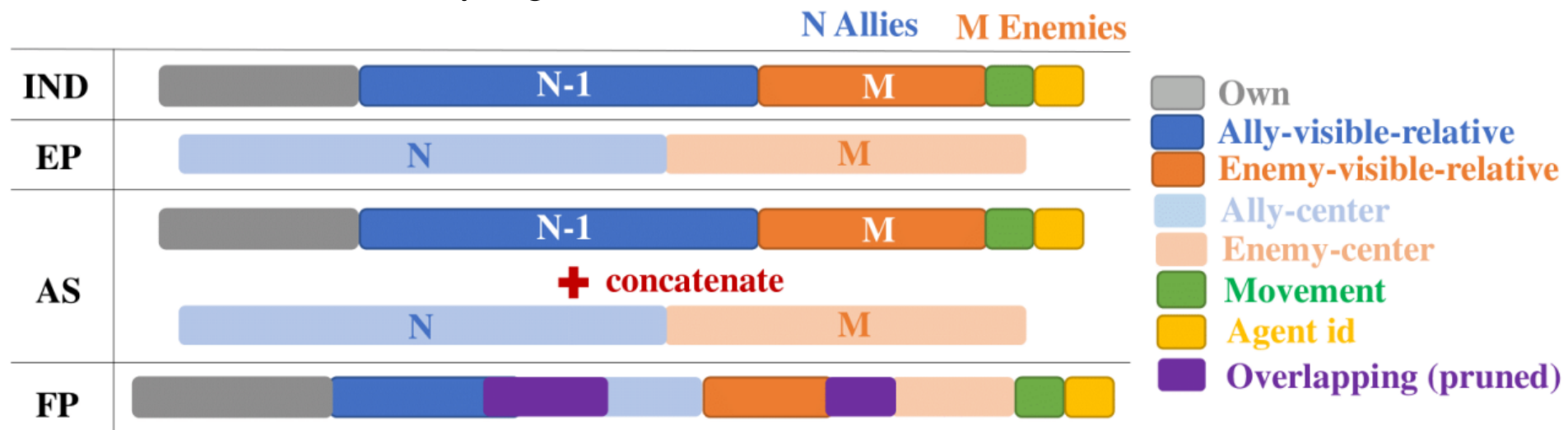
## □ MAPPO

- Input Representation to Value Function
  - Two common implementations for a centralized value function: **concatenation of all local observations**, an **environment-provided global state**. Other methods: **agent-specific global state representation**.
  - **Concatenation of all local observations(CL)** :
    - the value **input dimension can be extremely large** when the number of agent or the dimension of local observation is large.
    - CL also may **not contain sufficient global information** to reduce a POMDP to an MDP.
  - **Environment-provided global state(EP)**: Only contains information common to all agents and can **omit important local information**
  - **Agent-Specific Global State(AS)**: concatenate the environment state and the local observation
  - **Feature-Pruned Agent-Specific Global State(FP)**: to address the overlap between local and global state, removing all the duplicated features in AS.

# Multi-Agent RL

## MAPPO

- Input Representation to Value Function
  - IND: using decentralized inputs
  - EP: using environment-provided global state
  - AS: agent-specific global state which concatenates EP and IND
  - FP: removing the duplicated features in AS
  - Conclusion: include agent-specific features in the global state and check that these features do not make the state dimension substantially higher





## □ MAPPO

### ➤ PPO Clipping

- Clipping the importance ratio and value losses
- Constrain the policy and value functions from drastically changing between iterations
- $\epsilon$  hyperparameter: lower  $\epsilon$  values slows learning speed and higher  $\epsilon$  values result in larger variance and larger volatility in the performance
- Conclusion: tuning the clipping ratio  $\epsilon$  as a trade-off between training stability and fast convergence

### ➤ Death Masking

- Using these informative states for dead agents during value learning amplifies the bias of the learned value function
- Using an agent-specific constant vector, i.e., a zero vector with the agent's ID, as the input to the value function after an agent dies.
- Conclusion: use zero states with agent ID as the value input for dead agents.

# Multi-Agent RL

## □ MAPPO

### Algorithm 1 Recurrent-MAPPO

---

Initialize  $\theta$ , the parameters for policy  $\pi$  and  $\phi$ , the parameters for critic  $V$ , using Orthogonal initialization (Hu et al., 2020)  
Set learning rate  $\alpha$   
**while**  $step \leq step_{max}$  **do**  
  set data buffer  $D = \{\}$   
  **for**  $i = 1$  **to**  $batch\_size$  **do**  
     $\tau = []$  empty list  
    initialize  $h_{0,\pi}^{(1)}, \dots, h_{0,\pi}^{(n)}$  actor RNN states  
    initialize  $h_{0,V}^{(1)}, \dots, h_{0,V}^{(n)}$  critic RNN states  
    **for**  $t = 1$  **to**  $T$  **do**  
      **for all agents**  $a$  **do**  
         $p_t^{(a)}, h_{t,\pi}^{(a)} = \pi(o_t^{(a)}, h_{t-1,\pi}^{(a)}; \theta)$   
         $u_t^{(a)} \sim p_t^{(a)}$   
         $v_t^{(a)}, h_{t,V}^{(a)} = V(s_t^{(a)}, h_{t-1,V}^{(a)}; \phi)$   
      **end for**  
      Execute actions  $u_t$ , observe  $r_t, s_{t+1}, o_{t+1}$   
       $\tau += [s_t, o_t, h_{t,\pi}, h_{t,V}, u_t, r_t, s_{t+1}, o_{t+1}]$   
      **end for**  
      Compute advantage estimate  $\hat{A}$  via GAE on  $\tau$ , using PopArt  
      Compute reward-to-go  $\hat{R}$  on  $\tau$  and normalize with PopArt  
      Split trajectory  $\tau$  into chunks of length  $L$   
      **for**  $l = 0, 1, \dots, T/L$  **do**  
         $D = D \cup (\tau[l : l + T], \hat{A}[l : l + L], \hat{R}[l : l + L])$   
      **end for**  
    **end for**  
    **for** mini-batch  $k = 1, \dots, K$  **do**  
       $b \leftarrow$  random mini-batch from  $D$  with all agent data  
      **for each data chunk**  $c$  in the mini-batch  $b$  **do**  
        update RNN hidden states for  $\pi$  and  $V$  from first hidden state in data chunk  
      **end for**  
    **end for**  
    Adam update  $\theta$  on  $L(\theta)$  with data  $b$   
    Adam update  $\phi$  on  $L(\phi)$  with data  $b$   
  **end while**

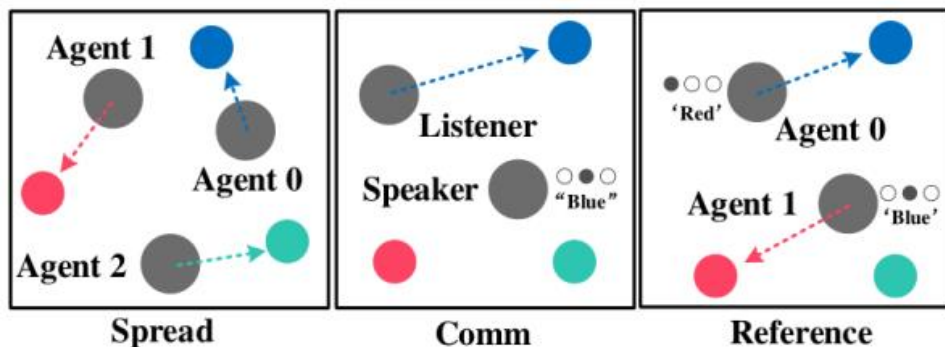
---

# Multi-Agent RL

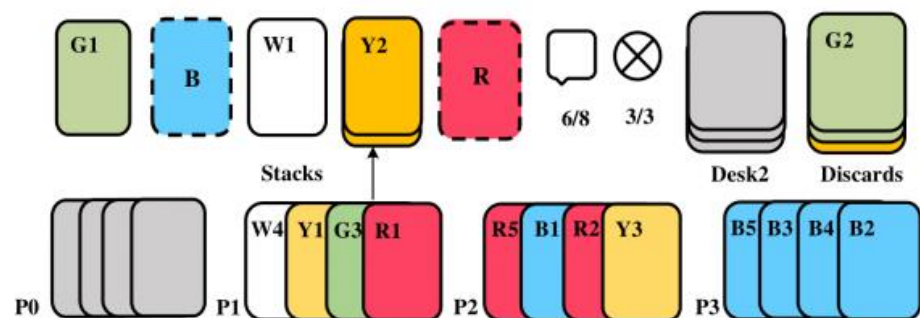
## MAPPO

### ➤ Main Results

#### ➤ Three multi-agent testbeds: MPE, SC2, Hanabi



(a) MPE scenarios



(b) 4-player Hanabi-Full



(c) SMAC corridor

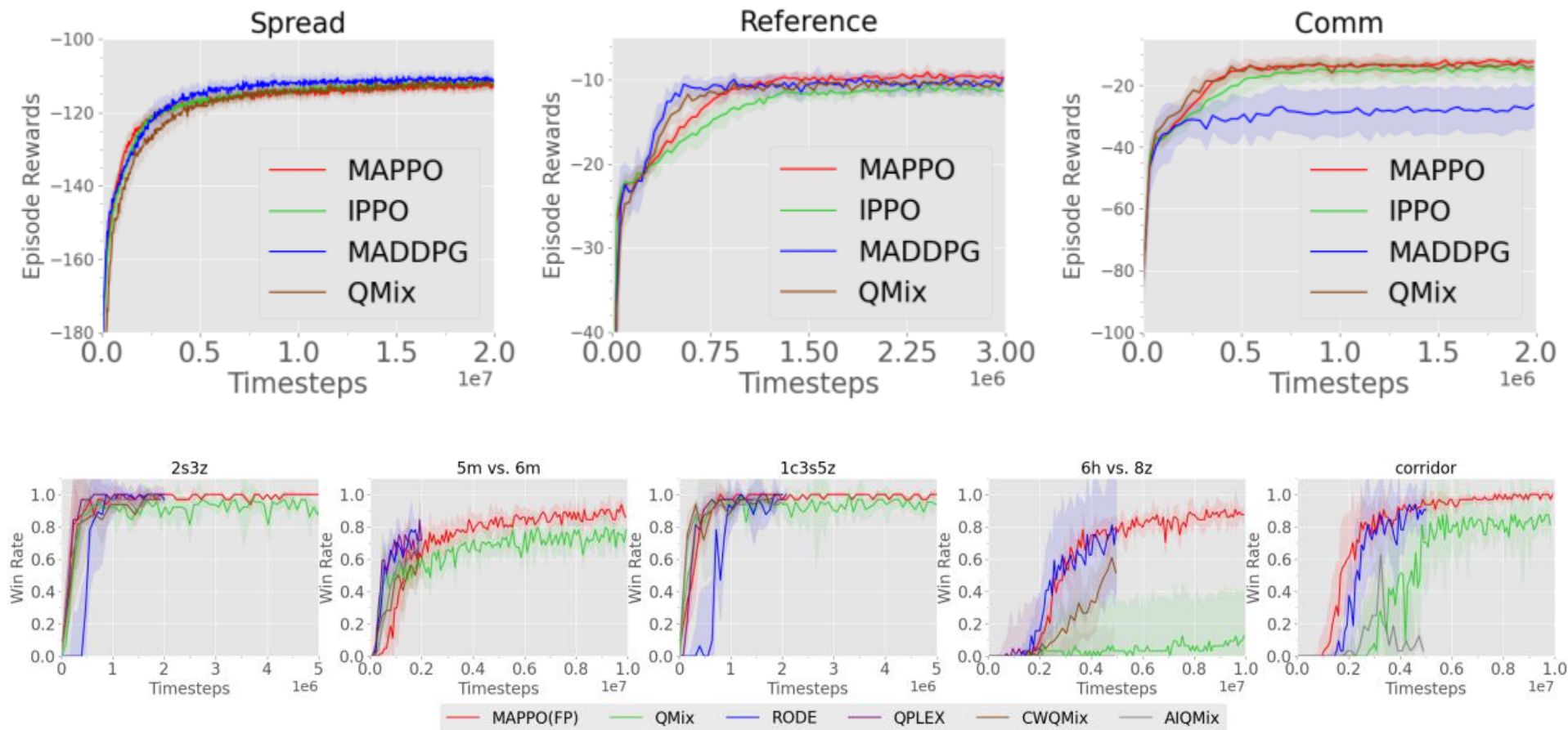


(d) SMAC 2c\_vs\_64zg

# Multi-Agent RL

## MAPPO

### ➤ Main Results





# Multi-Agent RL

## MAPPO

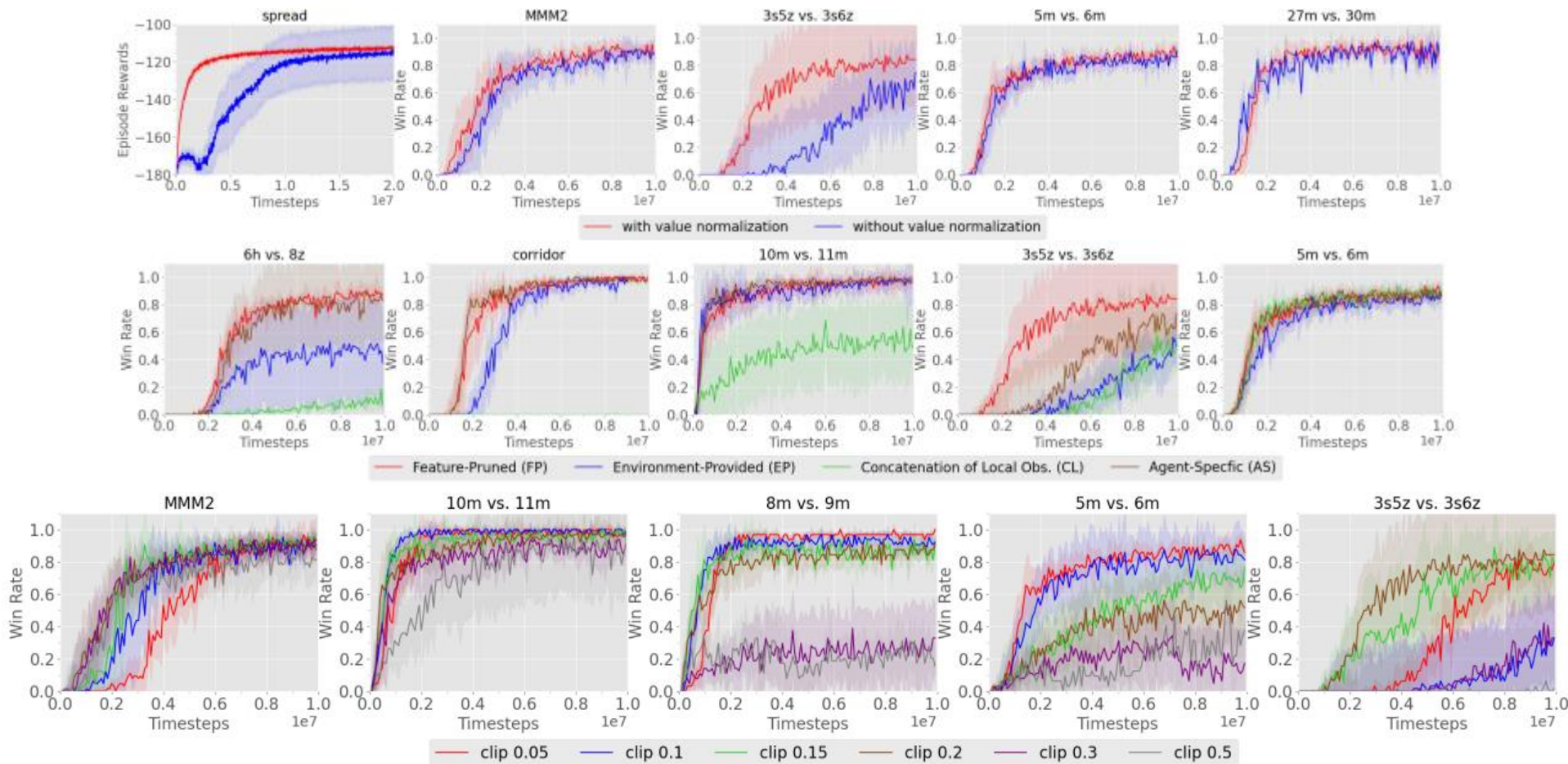
### ➤ Main Results

Map	MAPPO(FP)	MAPPO(AS)	IPPO	QMIX	RODE*	MAPPO*(FP)	MAPPO*(AS)
2m vs_1z	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>95.3</b> (5.2)	/	<u>100.0</u> (0.0)	<u>100.0</u> (0.0)
3m	<b>100.0</b> (0.0)	<b>100.0</b> (1.5)	<b>100.0</b> (0.0)	96.9(1.3)	/	<u>100.0</u> (0.0)	<u>100.0</u> (1.5)
2svs1sc	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>100.0</b> (1.5)	96.9(2.9)	<u>100.0</u> (0.0)	<u>100.0</u> (0.0)	<u>100.0</u> (0.0)
2s3z	<b>100.0</b> (0.7)	<b>100.0</b> (1.5)	<b>100.0</b> (0.0)	95.3(2.5)	<u>100.0</u> (0.0)	96.9(1.5)	96.9(1.5)
3svs3z	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>96.9</b> (12.5)	/	<u>100.0</u> (0.0)	<u>100.0</u> (0.0)
3svs4z	<b>100.0</b> (1.3)	<b>98.4</b> (1.6)	<b>99.2</b> (1.5)	<b>97.7</b> (1.7)	/	<u>100.0</u> (2.1)	<u>100.0</u> (1.5)
so many baneling	<b>100.0</b> (0.0)	<b>100.0</b> (0.7)	<b>100.0</b> (1.5)	96.9(2.3)	/	<u>100.0</u> (1.5)	96.9(1.5)
8m	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>100.0</b> (0.7)	97.7(1.9)	/	<u>100.0</u> (0.0)	<u>100.0</u> (0.0)
MMM	<b>96.9</b> (0.6)	93.8(1.5)	<b>96.9</b> (0.0)	<b>95.3</b> (2.5)	/	<u>93.8</u> (2.6)	<u>96.9</u> (1.5)
1c3s5z	<b>100.0</b> (0.0)	96.9(2.6)	<b>100.0</b> (0.0)	96.1(1.7)	<u>100.0</u> (0.0)	<u>100.0</u> (0.0)	96.9(2.6)
bane vs bane	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	<u>100.0</u> (46.4)	<u>100.0</u> (0.0)	<u>100.0</u> (0.0)
3svs5z	<b>100.0</b> (0.6)	<b>99.2</b> (1.4)	<b>100.0</b> (0.0)	<b>98.4</b> (2.4)	78.9(4.2)	<u>98.4</u> (5.5)	<u>100.0</u> (1.2)
2cvs64zg	<b>100.0</b> (0.0)	<b>100.0</b> (0.0)	98.4(1.3)	92.2(4.0)	<u>100.0</u> (0.0)	<u>96.9</u> (3.1)	95.3(3.5)
8mvs9m	<b>96.9</b> (0.6)	<b>96.9</b> (0.6)	<b>96.9</b> (0.7)	92.2(2.0)	/	<u>84.4</u> (5.1)	<u>87.5</u> (2.1)
25m	<b>100.0</b> (1.5)	<b>100.0</b> (4.0)	<b>100.0</b> (0.0)	85.9(7.1)	/	<u>96.9</u> (3.1)	<u>93.8</u> (2.9)
5mvs6m	<b>89.1</b> (2.5)	<b>88.3</b> (1.2)	<b>87.5</b> (2.3)	75.8(3.7)	<u>71.1</u> (9.2)	<u>65.6</u> (14.1)	<u>68.8</u> (8.2)
3s5z	<b>96.9</b> (0.7)	<b>96.9</b> (1.9)	<b>96.9</b> (1.5)	88.3(2.9)	<u>93.8</u> (2.0)	71.9(11.8)	53.1(15.4)
10mvs11m	<b>96.9</b> (4.8)	<b>96.9</b> (1.2)	<b>93.0</b> (7.4)	<b>95.3</b> (1.0)	<u>95.3</u> (2.2)	81.2(8.3)	<u>89.1</u> (5.5)
MMM2	<b>90.6</b> (2.8)	<b>87.5</b> (5.1)	<b>86.7</b> (7.3)	<b>87.5</b> (2.6)	<u>89.8</u> (6.7)	51.6(21.9)	28.1(29.6)
3s5zvs3s6z	<b>84.4</b> (34.0)	63.3(19.2)	<b>82.8</b> (19.1)	<b>82.8</b> (5.3)	<u>96.8</u> (25.11)	<u>75.0</u> (36.3)	18.8(37.4)
27mvs30m	<b>93.8</b> (2.4)	85.9(3.8)	69.5(11.8)	39.1(9.8)	<u>96.8</u> (1.5)	<u>93.8</u> (3.8)	<u>89.1</u> (6.5)
6hvs8z	<b>88.3</b> (3.7)	<b>85.9</b> (30.9)	<b>84.4</b> (33.3)	9.4(2.0)	<u>78.1</u> (37.0)	<u>78.1</u> (5.6)	<u>81.2</u> (31.8)
corridor	<b>100.0</b> (1.2)	<b>98.4</b> (0.8)	<b>98.4</b> (3.1)	84.4(2.5)	<u>65.6</u> (32.1)	<u>93.8</u> (3.5)	<u>93.8</u> (2.8)

# Multi-Agent RL

## MAPPO

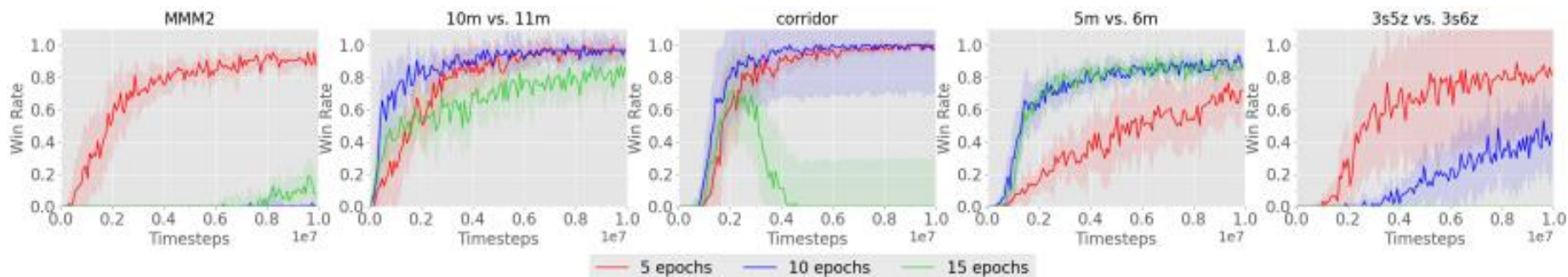
### ➤ Main Results--Ablation



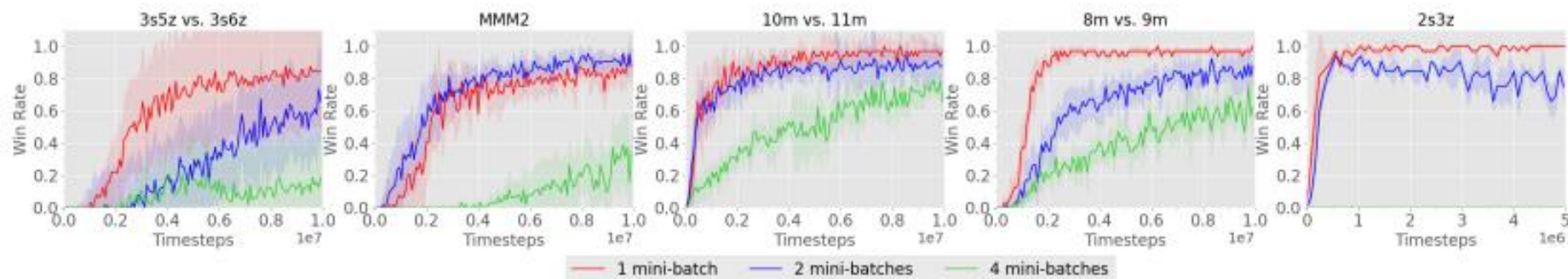
# Multi-Agent RL

## MAPPO

### ➤ Main Results--Ablation



(a) effect of different training epochs.



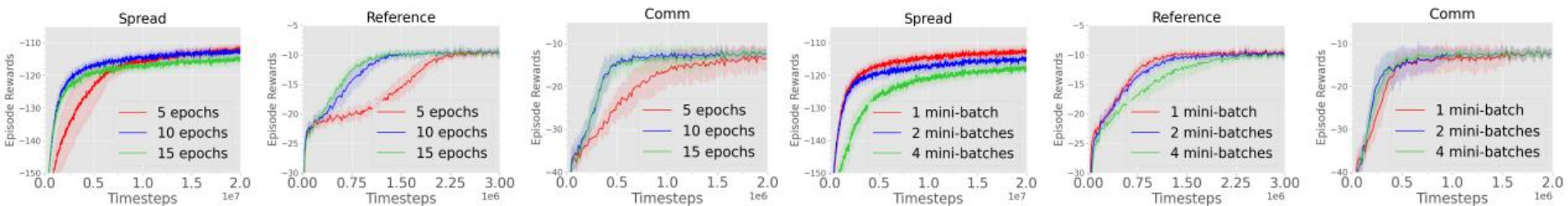
(b) effect of different mini-batch numbers.



# Multi-Agent RL

## MAPPO

### ➤ Main Results--Ablation



(a) effect of different training epochs.

(b) effect of different mini-batch numbers.

