# MASER: Multi-Agent Reinforcement Learning with Subgoals Generated from Experience Replay Buffer

Jeewon Jeon

Advisor : Youngchul Sung

15 July, 2022

**Smart Information Systems Research Lab**
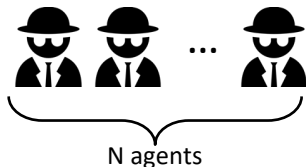**Dept. of Electrical Engineering, KAIST**

# Content

1. Introduction & Contribution

2. Proposed Algorithm

3. Experiment and Result

4. Conclusion

# Introduction

Multi agent reinforcement learning with sparse rewards

- Joint action space grows exponentially with the number of agents.

- Typically, agents receive a global reward.

- Reward comes only under certain circumstances, e.g., success/failure.

- We approach multi-agent sparse RL with sub-goals and
-     propose a method to determine sub-goals by exploiting experience replay buffer.

4 actions

N agents

- Total joint action space : $4^N$
- $4^N$ contributes only 1 sparse & global reward
- Which agent and actions contributes more?

# Contributions

Our method has <span style="color:red">three</span> contributions :

- **Generating and assigning subgoals**: MASER finds <span style="color:blue">subgoals</span> for agents from <span style="color:blue">the experience replay buffer</span>. This eliminates the necessity of predesigning good subgoals based on domain knowledge.

- **Giving individual rewards**: MASER designs <span style="color:blue">individual rewards</span> for <span style="color:blue">local agents</span> to reach their subgoals while maximizing the joint return.

- **Actionable distance relevant to Q-learning** : To determine the intrinsic reward based on the Euclidean distance in the transformed domain, MASER uses representational transform based on <span style="color:blue">actionable distance relevant to Q-learning</span> derived from Amari 0-divergence.
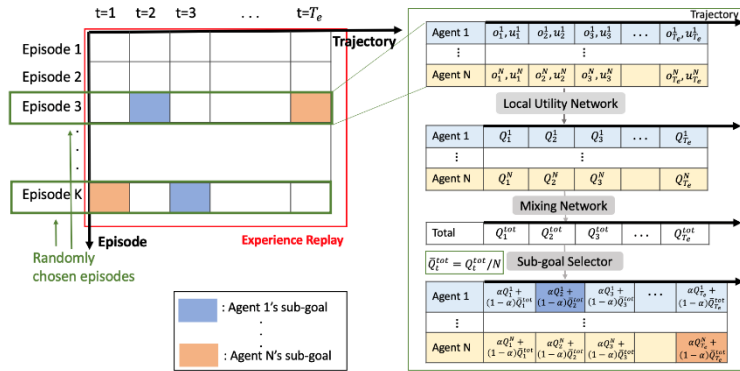
# Proposed Algorithm : Generating Subgoals



Figure1. Generating Subgoals

$$t_*^i = argmax_t[\alpha Q^i\left(o_t^i, argmax_u Q^i(o_t^i, u)\right) + (1 - \alpha)Q^{tot}(\boldsymbol{o}_t, \boldsymbol{u}_t)]$$

Local Q-value $\qquad$ Global Q-value

Generating different subgoals for each agent $\qquad$ Consider other agents' status

Subgoal for agent $i : o_g^i = o_{t_*^i}^i$

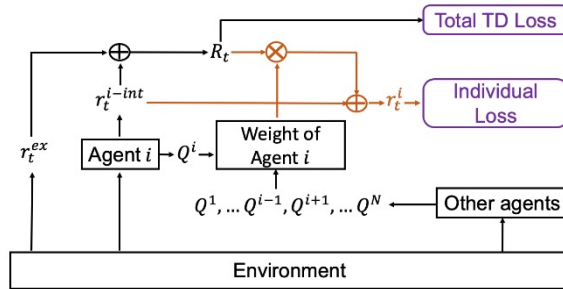# Proposed Algorithm : Overall Reward Design



Figure2. Overall reward design diagram

Both mixing and local utility parameters are learned to achieve the subgoals as well as maximizing the overall extrinsic reward

Individual intrinsic reward

$$r_t^{i-int} = -\left\| \boxed{\phi^i}(o_t^i) - \phi^i(o_g^i) \right\|_2 \longrightarrow \text{Agents to reach their subgoals}$$

Actionable representational transform

Proxy reward

$$R_t = \boxed{r_t^{ex}} + \lambda \frac{1}{N} \sum_{i=1}^{N} r_t^{i-int} \longrightarrow \text{To update mixing network } \theta$$

Extrinsic reward (sparse)

Individual reward

$$r_t^i = \boxed{softmax(\max_u Q^i(o_t^i, u))} \cdot R_t + \lambda r_t^{i-int}$$

Contribution of Agent $i$ to overall extrinsic reward

$\longrightarrow$ To update utility parameter $\theta\_i$

# Proposed Algorithm : Q-Function-Based Representation Learning

Actionable distance

$$D_Q\left(o_t^i, o_g^i\right) = 1 \; - \; \frac{< Q^i\left(o_t^i, \cdot\right), Q^i\left(o_g^i, \cdot\right) >}{\left\|Q^i\left(o_t^i, \cdot\right)\right\| \times \left\|Q^i\left(o_g^i, \cdot\right)\right\|}$$

$\longrightarrow$ 1 – cosine similarity between $Q^i\left(o_t^i, \cdot\right), Q^i(o_g^i, \cdot)$

Loss function

$$L_D\left(\phi^i\right) = E_{o_t^i}\left[\left\|\phi^i\left(o_t^i\right) - \phi^i\left(o_g^i\right)\right\|_2 - D_Q\left(o_t^i, o_g^i\right)\right]^2$$

$\longrightarrow$ By minimizing loss function, $\phi^i$ is learned to represent actionable distance

# Proposed Algorithm : Overall Flow



Figure3. Overflow

- Reach subgoals by giving intrinsic reward

- Novel distance function with cosine similarity

- Episodic correction after reaching subgoals.
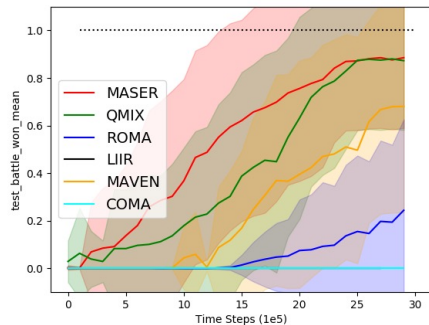
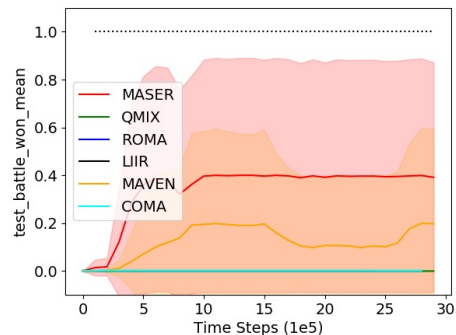# Results (Experiments on StarCraft 2 with Sparse Rewards)
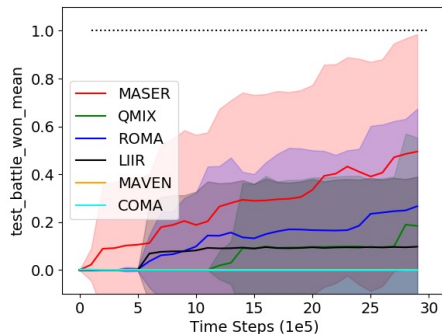

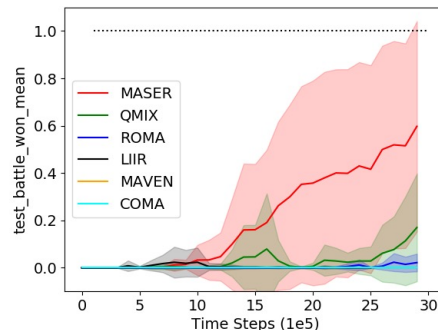
Figure4(a). 3m

Figure4(b). 2m_vs_1z

Figure4(c). 8m
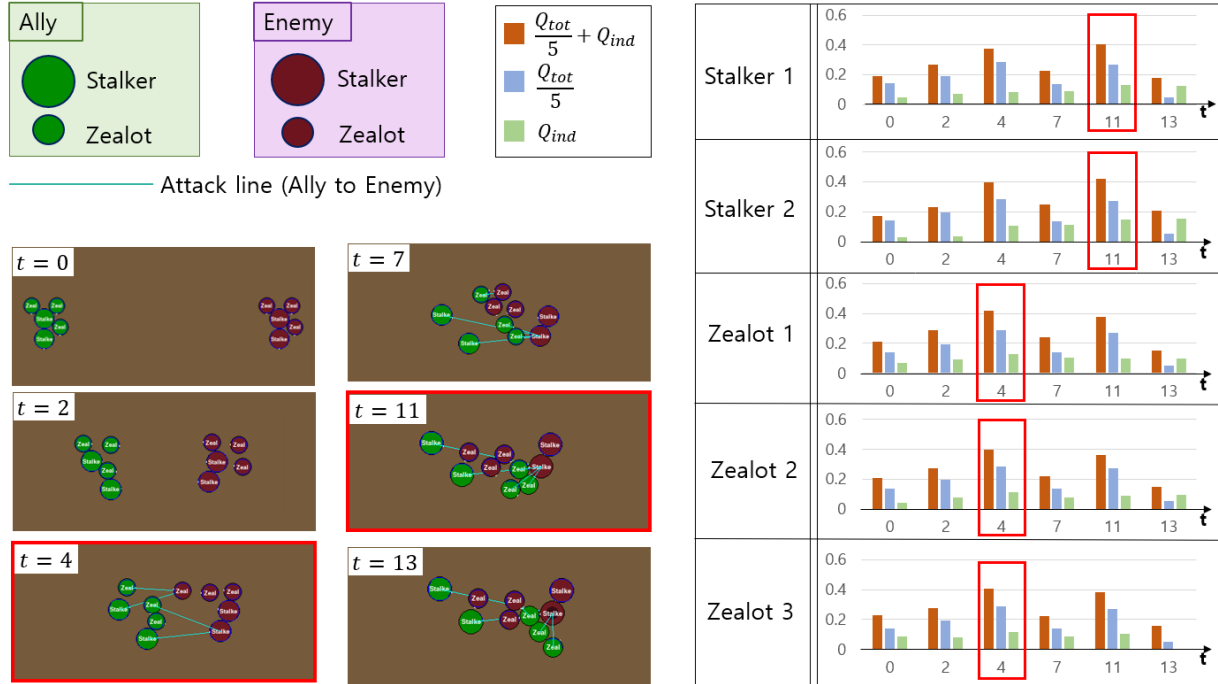
Figure4(d). 2s3z

# Example of Generated Subgoals



Figure5. Visualization of suboglas on 2s3z

# Thank you!