

Hierarchical Multi-Agent Reinforcement Learning with Options

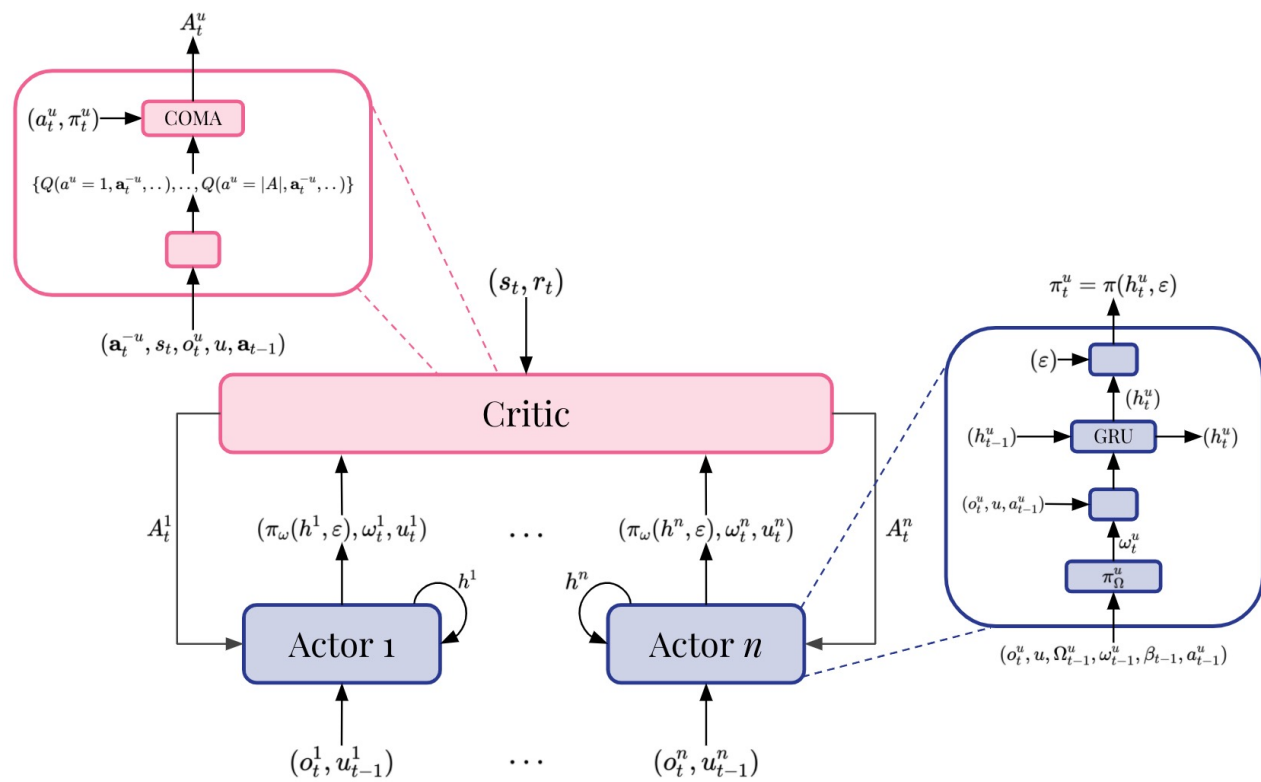
Student: Ang Wan Qi

Supervisor: Asst Prof Lana Obraztsova

Project Objectives

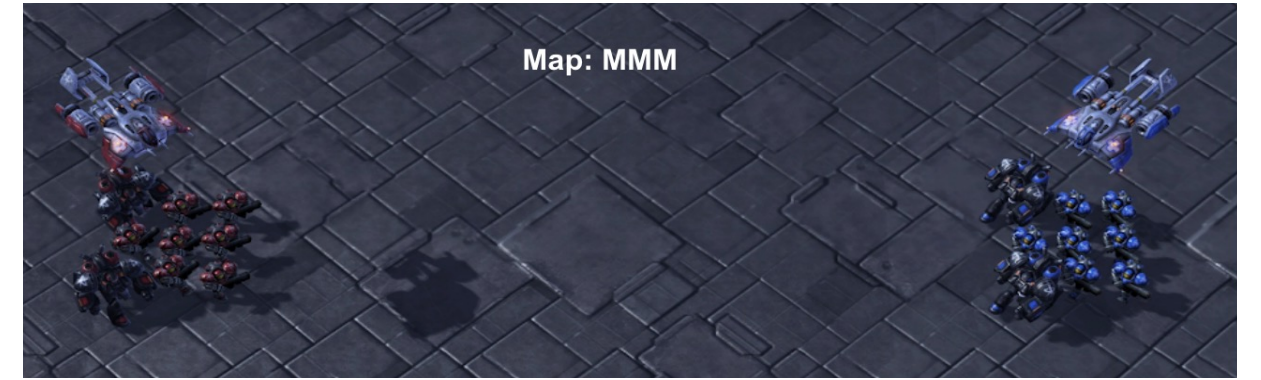
Cooperative multi-agent reinforcement learning (MARL) demands that multiple agents work together towards a common goal. Agents in these tasks often select and complete actions at different times and require asynchronous action execution. However, current methods have actions modeled as primitive operations and synchronized action execution. Hence, this project aims to explore hierarchical reinforcement learning (HRL) methods such as the options framework to achieve temporal abstraction of actions. A novel algorithm named COMA-OC is proposed to train the decentralized agents in StarCraft II's micromanagement scenarios against an enemy army.

COMA-OC Architecture

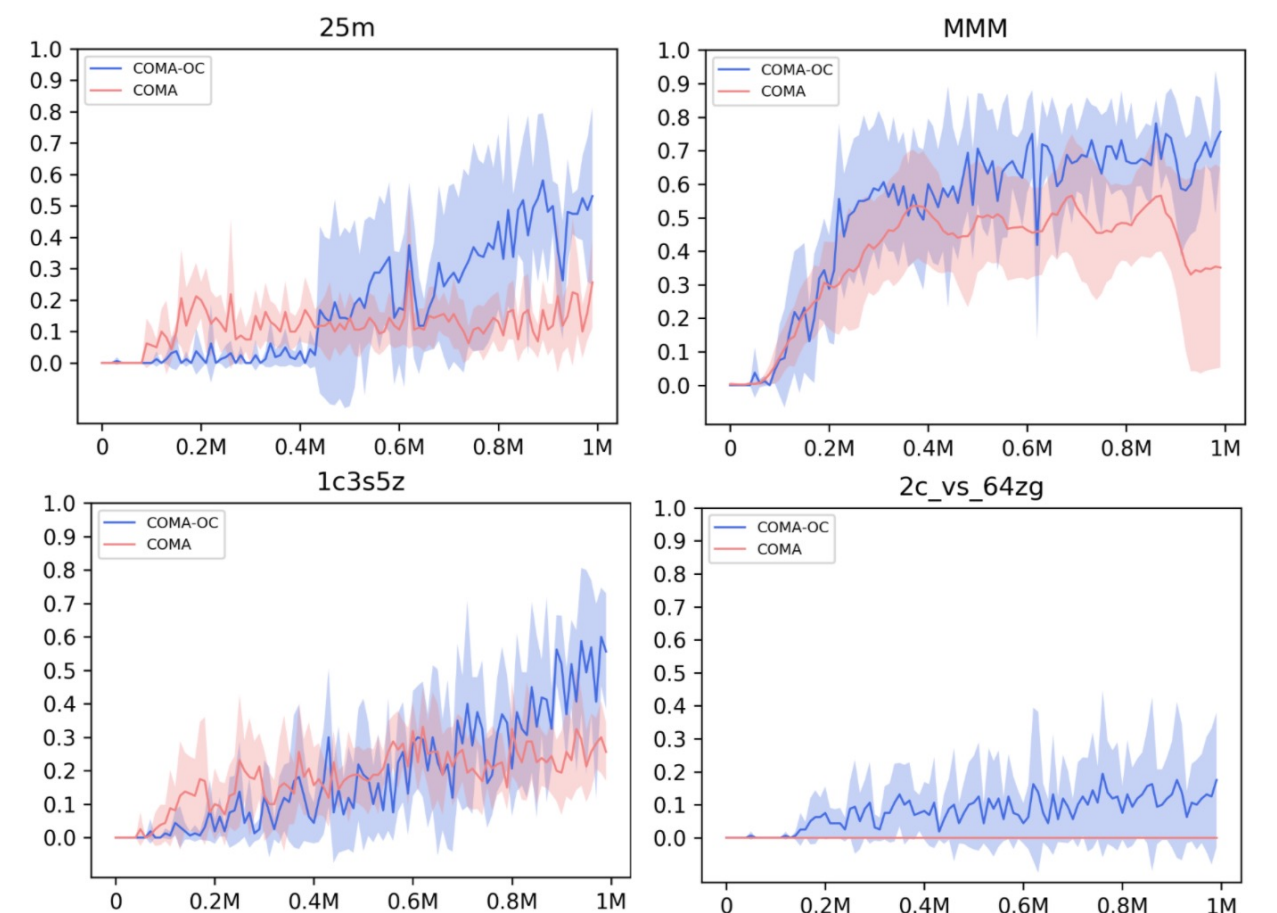


Evaluation Environment

COMA-OC is evaluated in StarCraft II where 2 armies battle each other in a scenario.



Results



COMA-OC Pseudocode

Algorithm 1 COMA-OC Algorithm

```

1: Initialise  $\theta_1^c, \hat{\theta}_1^c, \theta_1^\pi, \theta_1^\beta$ 
2: for each training episode  $e$  do
3:   Empty buffer
4:   for  $e_c = 1$  to  $\frac{\text{BatchSize}}{n}$  do
5:      $s_1$  = initial state,  $t = 0$ ,  $h_0^u = \mathbf{0}$  for each agent  $u$ 
6:     Choose  $\omega$  according to an  $\epsilon$ -soft policy over options  $\pi_\Omega(s_1)$ 
7:     while  $s_t \neq \text{terminal}$  and  $t < T$  do
8:        $t = t + 1$ 
9:       for each agent  $u$  do
10:        if  $\omega$  terminates in  $s_t$  according to  $\beta_{\omega, \theta^\beta}$  then
11:          Choose new  $\omega$  according to  $\epsilon\text{-soft}\pi_\Omega(s_t)$ 
12:           $h_t^u = \text{Actor}(o_t^u, h_{t-1}^u, a_{t-1}^u, u, a, \theta_i)$ 
13:          Sample  $a_t^u$  from  $\pi_{\omega, \theta^\pi}(h_t^u, \epsilon(e))$ 
14:          Get reward  $r_t$  and next state  $s_{t+1}$ 
15:          Add episode to buffer
16:   Collate episodes in buffer into single batch
17:   for  $t = 1$  to  $T$  do // processing all agents in parallel via single batch
18:     Batch unroll RNN using states, actions and rewards
19:     Calculate TD( $\lambda$ ) targets  $y_t^u$  using  $\hat{\theta}_i^c$  // targets for actions
20:     Calculate update target  $g_t$  for options:
21:     if  $s$  is non-terminal then
22:        $g_t^u = r + \gamma(1 - \beta_{\omega, \theta^\beta}(s))Q_\Omega(s, \omega) + \gamma\beta_{\omega, \theta^\beta}(s)\max_{\omega'} Q_\Omega(s, \omega')$ 
23:   for  $t = T$  down to 1 do
24:      $\Delta Q_1^u(t) = y_t^u - Q_1(s_t^u, \mathbf{a})$  // for actions
25:      $\Delta Q_2^u(t) = g_t^u - Q_2(s_t^u, \mathbf{a})$  // for options
26:      $\Delta \theta^c = \nabla_{\theta^c}((\Delta Q_1^u(t))^2 + (\Delta Q_2^u(t))^2)$  // calculate critic gradient
27:      $\theta_{i+1}^c = \theta_i^c - \alpha \Delta \theta^c$  // update critic weights
28:     Every C steps reset  $\hat{\theta}_i^c = \theta_i^c$ 
29:   for  $t = T$  down to 1 do
30:      $A^u(s_t^u, \mathbf{a}) = Q(s_t^u, \mathbf{a}) - \sum_a Q(s_t^u, a, \mathbf{a}^{-u})\pi(a|h_t^u)$  // calculate COMA
31:      $\Delta \theta^\pi = \Delta \theta^\pi + \nabla_{\theta^\pi} \log \pi(a|h_t^u) A^u(s_t^u, \mathbf{a})$  // accumulate actor gradients
32:      $\Delta \theta^\beta = \Delta \theta^\beta + \nabla_{\theta^\beta} \beta(s_t^u) A^u(s_t^u, \mathbf{a})$  // accumulate beta gradients
33:      $\theta_{i+1}^\pi = \theta_i^\pi + \alpha \Delta \theta^\pi$  // update actor weights (pi loss)
34:      $\theta_{i+1}^\beta = \theta_i^\beta + \alpha \Delta \theta^\beta$  // update actor weights (beta loss)

```

Conclusion and Future Work

Results showed that learning options leads to significant improvements in performance in various scenarios as compared to the baseline method. One possible improvement is to incorporate communication tools between agents to enhance cooperation.