

Appendix

1 DDPG Experiments

In addition to SAC, we also provide experiments on DDPG[1]. Instead of the probabilistic policy used in SAC, DDPG has a deterministic policy. The average performances are shown in 1. Though DDPG-based algorithms have worse performance than SAC, CLOS, ER, GEM and APD still have a satisfying result, which is consistent to that in SAC.

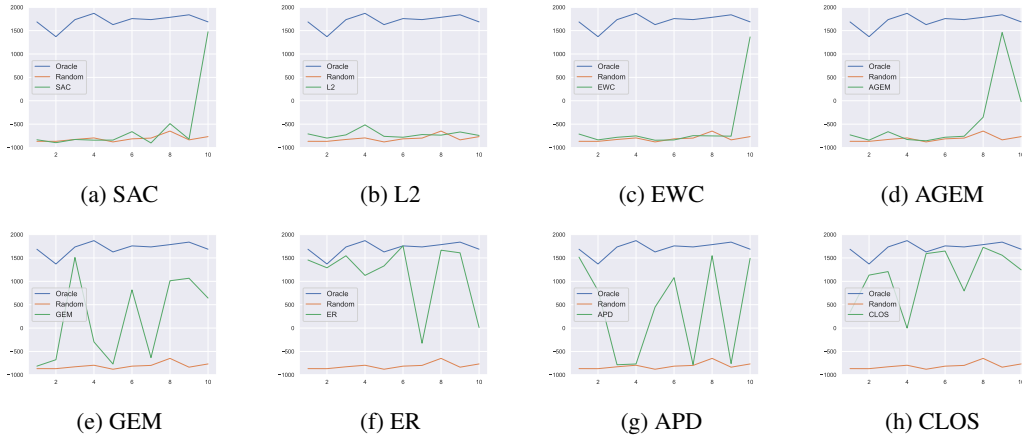


Figure 1: Performance of lifelong learning methods with the RL framework. The x-axis denotes task IDs. The y-axis denotes cumulative rewards. Each lifelong reinforcement learning algorithm consists of a lifelong learning algorithm and DDPG as the RL algorithm. The figure shows the performance of the policy on each task after finishing learning the last one. Oracle and random agent is included in each figure for comparison.

2 Reinforcement Lifelong Learning Framework

The pseudo code is shown in Algorithm. 1. The descriptive figure of it is shown in Fig. 2.

References

- [1] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. In *ICLR (Poster)*, 2016.

Algorithm 1 Reinforcement Lifelong Learning Framework

```

1: INPUT  $K$ : the number of tasks needed to be learnt;
2:  $D \leftarrow \{\} * K$  # Initialize  $K$  empty replay buffer;
3:  $C \leftarrow \{c_1, c_2, \dots, c_K\}$  # Initialize the critic set;
4: for  $k \in \{1, 2, \dots, K\}$  do
5:    $(s'_t, a'_t, s'_{t+1}, r'_t, d'_t) \sim RollOut(\theta, k)$ ;
6:    $D_k \leftarrow D_k \cup (s'_t, a'_t, s'_{t+1}, r'_t, d'_t)$ ;
7:    $(s_t^k, a_t^k, s_{t+1}^k, r_t^k, d_t^k) \sim D_k$  # Sample a data from the replay buffer;
8:    $l_c \leftarrow (Q_{\phi^k}(s_t^k, a_t^k) - (r_t^k + \gamma(1 - d_t^k)Q_{\phi^k}(s_{t+1}^k, \pi_\theta(s_t^k))))^2$  # Compute the critic loss;
9:    $g_c^k \leftarrow \nabla_{\phi^k} l_c$ ;
10:   $\phi^k \leftarrow \phi^k - \beta g_c^k$ ; # Update the critic parameters;
11:   $l_a, g_a^k \leftarrow -c_k(\pi_\theta(s_t^k), k), -\nabla c_k(\pi_\theta(s_t^k), k)$  # Compute the actor loss;
12:   $l'_a, g_a^{k'} \leftarrow LifelongLearning(l_a, g_a^k)$  # Update the Actor Loss and gradients according
    to the lifelong learning algorithm
13:   $\theta_c \leftarrow \theta_c - \alpha g_a^{k'}$  # Update the shared module in the actor
14:   $\theta_s \leftarrow \theta_s - \alpha g_a^k$  # Update the task-specific module in the actor.
15: end for
  
```

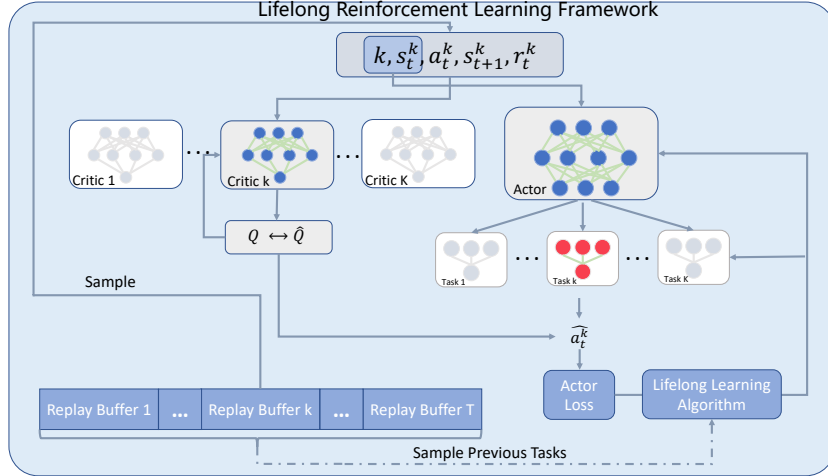


Figure 2: Reinforcement Lifelong Learning Framework. The framework combines an Actor-Critic algorithm with a lifelong learning algorithm. The actor is decomposed into shared parts and task-specific parts. Each task is assigned an individual critic. We keep the replay memory of previous tasks and sample data from them when needed in the lifelong learning process.