

**Algorithm S2** Asynchronous n-step Q-learning - pseudocode for each actor-learner thread.

---

// Assume global shared parameter vector  $\theta$ .  
// Assume global shared target parameter vector  $\theta^-$ .  
// Assume global shared counter  $T = 0$ .  
Initialize thread step counter  $t \leftarrow 1$   
Initialize target network parameters  $\theta^- \leftarrow \theta$   
Initialize thread-specific parameters  $\theta' = \theta$   
Initialize network gradients  $d\theta \leftarrow 0$   
**repeat**  
    Clear gradients  $d\theta \leftarrow 0$   
    Synchronize thread-specific parameters  $\theta' = \theta$   
     $t_{start} = t$   
    Get state  $s_t$   
    **repeat**  
        Take action  $a_t$  according to the  $\epsilon$ -greedy policy based on  $Q(s_t, a; \theta')$   
        Receive reward  $r_t$  and new state  $s_{t+1}$   
         $t \leftarrow t + 1$   
         $T \leftarrow T + 1$   
    **until** terminal  $s_t$  **or**  $t - t_{start} == t_{max}$   
     $R = \begin{cases} 0 & \text{for terminal } s_t \\ \max_a Q(s_t, a; \theta^-) & \text{for non-terminal } s_t \end{cases}$   
    **for**  $i \in \{t - 1, \dots, t_{start}\}$  **do**  
         $R \leftarrow r_i + \gamma R$   
        Accumulate gradients wrt  $\theta'$ :  $d\theta \leftarrow d\theta + \frac{\partial (R - Q(s_i, a_i; \theta'))^2}{\partial \theta'}$   
    **end for**  
    Perform asynchronous update of  $\theta$  using  $d\theta$ .  
    **if**  $T \bmod I_{target} == 0$  **then**  
         $\theta^- \leftarrow \theta$   
    **end if**  
**until**  $T > T_{max}$ 

---