

**Algorithm 1** DRL-based Decentralized CW Optimization

---

```

  ▷ ### Initialization ###
1: Initialize the observation buffer,  $O$ , with zeroes
2: Initialize the weights,  $\theta$ , of the agent
3: Get the action function,  $A_\theta$ , which the agent uses to choose the action according to the
   current stage
4: Initialize the algorithm's interaction period with the environment,  $envStepTime$ 
5: Initialize the training stage period,  $trainingPeriod$ 
6: Set  $trainingFlag \leftarrow True$  to tell the algorithm is in the training stage
7: Initialize the experience replay buffer,  $E$ , with zeroes.
8:  $trainingStartTime \leftarrow currentTime$ 
9:  $lastUpdate \leftarrow currentTime$ 
10:  $\mu_{prev}(i) \leftarrow 0$  (previous mean value)
11:  $\sigma_{prev}^2(i) \leftarrow 0$  (previous variance value)
12: Set  $useQueueLevelFlag \leftarrow True$  to use the averaged normalized transmission queues'
   level as observation.
13:  $CW \leftarrow 15$ 

14: for  $t = 1, \dots, \infty$  do
15:   for  $t = 1, \dots, \infty$  do

     ▷ ### Pre-learning stage ###
16:      $N_t(i) \leftarrow$  get number of transmitted frames
17:      $N_r(i) \leftarrow$  get number of received frames
18:      $observation(i) \leftarrow \frac{N_t(i) - N_r(i)}{N_t(i)}$ 
19:      $O(i).append(observation(i))$ 

20:     if  $currentTime \geq lastUpdate + envStepTime$  then

       ▷ ### Learning and operational stages ###
21:        $\mu(i), \sigma^2(i) \leftarrow preprocess(O(i))$ 
22:        $a(i) \leftarrow A_{\theta(i)}(\mu(i), \sigma^2(i), trainingFlag)$ 
23:        $CW(i) \leftarrow 2^{a(i)+4} - 1$ 

24:       if  $trainingFlag == True$  then
25:          $N_{RP}(i) \leftarrow$  get the number of received packets.
26:          $tput(i) \leftarrow \frac{N_{RP}}{envStepTime}$ 
27:         Send the throughput of each station to the access point.
28:          $r \leftarrow normalize(tput(i))$ 
29:         Broadcast the new  $r$  reward value to all associated stations
30:          $E(i).append((\mu(i), \sigma^2(i), a(i), r, \mu_{prev(i)}, \sigma_{prev}^2(i)))$ 
31:          $\mu_{prev(i)} \leftarrow \mu(i)$ 
32:          $\sigma_{prev}^2(i) \leftarrow \sigma^2(i)$ 
33:          $mb(i) \leftarrow$  get random mini-batch from  $E(i)$ 
34:         Update  $\theta(i)$  based on  $mb(i)$ 
35:       end if

36:        $lastUpdate \leftarrow currentTime$ 
37:     end if

     ▷ ### Makes the transition between learning and operational stages ###
38:     if  $currentTime \geq trainingStartTime + trainingPeriod$  then
39:        $trainingFlag \leftarrow False$ 
40:     end if
41:   end for

```

---