**Algorithm 1:** ABR Model Training

---

1: Initialize critic networks $Q_{\theta_1}$, $Q_{\theta_2}$, actor network $\pi_\phi$
2: Initialize target networks $Q_{\theta_1^{tar}}$, $Q_{\theta_2^{tar}}$, $\pi_{\phi^{tar}}$
3: Initialize replay buffer $\mathcal{D}$, mini-batch $\Omega$
4: Initialize soft update factor $\tau$, delay update para $\delta$
5: **for** *each video session* **do**
6:     **for** segment $t = 1$ to $T$ **do**
7:         Select action (bitrate): $a_t = \pi_\phi(s_t)$
8:         Observe reward (QoE) $r_t$, and done $d_t$
9:         Store tuple $(s_t, a_t, r_t, d_t)$ in $\mathcal{D}$
10:        Sample a mini-batch of $j$ sequences, each
          with $n$ consecutive steps from $\mathcal{D}$:

11:
$$(s_{t'_0+i}, a_{t'_0+i}, r_{t'_0+i}, d_{t'_0+i})_{i=0,1,2,\dots,n-1}$$
$$(s_{t'_1+i}, a_{t'_1+i}, r_{t'_1+i}, d_{t'_1+i})_{i=0,1,2,\dots,n-1}$$
$$\vdots$$
$$(s_{t'_{j-1}+i}, a_{t'_{j-1}+i}, r_{t'_{j-1}+i}, d_{t'_{j-1}+i})_{i=0,1,2,\dots,n-1}$$

12:        $R, m = \text{ComputeCumulativeReward}(n)$
13:        $Q_{tar} = \text{ComputeTargetQValue}(R, m, n)$
14:        Update $Q_{\theta_1}$, $Q_{\theta_2}$ by minimizing:
15:        $E_\Omega[(Q_{\theta_{i=1,2}}(s_{t'}, a_{t'}) - Q_{tar})^2]$
16:        **if** $t \bmod \delta == 0$ **then**
17:            Update $\pi_\phi$ by maximizing:
18:            $E_\Omega[Q_{\theta_1}(s_{t'}, \pi_\phi(s_{t'}))]$
19:            Softly update target networks:
20:            $\phi^{tar} = \tau\phi^{tar} + (1-\tau)\phi$
21:            $\theta_1^{tar} = \tau\theta_1^{tar} + (1-\tau)\theta_1$
22:            $\theta_2^{tar} = \tau\theta_2^{tar} + (1-\tau)\theta_2$

23: **Function** ComputeCumulativeReward($n$)**:**
24:     Initialize cumulative reward: $R = 0$
25:     Initialize video terminal tag: $m = n$
26:     **for** $k = 0$ to $n$-1 **do**
27:         Update cumulative reward:
28:         $R = R + \gamma^k r_{t'+k}$
29:         **if** $d_{t'+k} == 1$ **then**
30:            Mark the termination step: $m = k$
31:            **return** $R, m$
32:     **return** $R, m$
33: **Function** ComputeTargetQValue($R, m, n$)**:**
34:     **if** $m < n$ **then**
35:         **return** $R$
36:     **else**
37:         Calculate target bitrate with Gaussian noize:
38:         $a_{tar} = \pi_{\phi^{tar}}(s_{t'+n}) + \mathcal{N}(\mu, \sigma^2)$
39:         Calculate minimum target Q value:
40:         $Q_{tar_1} = Q_{\theta_1^{tar}}(s_{t'+n}, a_{tar})$
41:         $Q_{tar_2} = Q_{\theta_2^{tar}}(s_{t'+n}, a_{tar})$
42:         $min\_value = \min(Q_{tar_1}, Q_{tar_2})$
43:         **return** $R + \gamma^n \times min\_value$