**Algorithm 1:** Proximal Policy Optimization (PPO)

**Input:** Initial policy parameters $\theta$, value function parameters $\phi$,
clipping threshold $\epsilon$, total batch size $B$, minibatch size $M$,
number of update epochs $K$, trajectory horizon $T$, learning
rates $\lambda_\pi$, $\lambda_V$

1 **for** *each* ***epoch*** **do**
2   Initialize buffer $\mathcal{D} \leftarrow \emptyset$;
3   Reset environment and observe initial state $s_0$;
4   **while** *buffer $\mathcal{D}$ not full* **do**
5     Sample action $a_t \sim \pi_\theta(\cdot|s_t)$ and store $\log \pi_\theta(a_t|s_t)$;
6     Execute $a_t$ in environment;
7     Observe $r_t$, $s_{t+1}$, done signal $d_t$;
8     Store $(s_t, a_t, r_t, \log \pi_\theta(a_t|s_t), d_t)$ in $\mathcal{D}$;
9     **if** $d_t$ *is True* **then**
10       Reset environment and observe new $s_{t+1}$;
11     **end**
12     $s_t \leftarrow s_{t+1}$;
13   **end**
14   Compute advantage estimates $\hat{A}_t$ and returns $\hat{R}_t$ using GAE or TD;
15   **for** $k = 1$ **to** $K$ **do**
16     Shuffle $\mathcal{D}$ and split into minibatches of size $M$;
17     **for** *each minibatch* **do**
18       Compute importance ratio:

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$$

19       Compute clipped objective:

$$L^{\text{CLIP}}(\theta) = \frac{1}{M} \sum \min\left(r_t(\theta)\hat{A}_t,\ \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)\hat{A}_t\right)$$

20       Update policy:
$$\theta \leftarrow \theta + \lambda_\pi \nabla_\theta L^{\text{CLIP}}(\theta)$$

21       Compute value loss:

$$L^V(\phi) = \frac{1}{M} \sum \left(V_\phi(s_t) - \hat{R}_t\right)^2$$

22       Update value function:

$$\phi \leftarrow \phi - \lambda_V \nabla_\phi L^V(\phi)$$

23     **end**
24   **end**
25 **end**

1