

Algorithm 1 PPO, Actor-Critic Style

for iteration=1, 2, ... **do**

for actor=1, 2, ..., N **do**

 Run policy $\pi_{\theta_{\text{old}}}$ in environment for T timesteps

 Compute advantage estimates $\hat{A}_1, \dots, \hat{A}_T$

end for

 Optimize surrogate L wrt θ , with K epochs and minibatch size $M \leq NT$

$\theta_{\text{old}} \leftarrow \theta$

end for