Algorithm: IPPO

1: Initialize: $\theta$ the parameters of individual actors $pi_i(.;\theta)$, $\phi$ the parameters of individual critics $V_i(;\phi)$
2: **while** $t < T$ **do**
3:     Initialize a rollout buffer $\mathcal{D}$ // $(episode_1, episode_2, \ldots)$
4:     **for** a number of episodes **do**:
5:         $current\_episode = \{\}$
6:         **while** $\mathbf{o}_t$ is not *done* **do**
7:             Collect observations $\{o_1^t, \ldots, o_n^t\}$
8:             Sample an action $a_i^t \sim \pi_i(.|o_i^t)$ for each agent $i$
9:             Execute joint action $\mathbf{a}^t = (a_1^t, \ldots, a_n^t)$
10:            Collect $r^t$, $done^t$
11:            Store $(\mathbf{o}^t, \mathbf{a}^t, r^t, done^t)$ in $current\_episode$
12:         **end while**
13:         Store $current\_episode$ in the rollout buffer $\mathcal{D}$
14:     **end for**
15:     Process the rollout buffer for batch training // Episodes with different lengths
16:     Compute the advantages $A_i$ and TD targets $y$
17:     Compute the actor losses:

$$\mathcal{L}(\theta) = \tfrac{1}{|\mathcal{B}|} \sum_b \tfrac{1}{L^b} \sum_t \tfrac{1}{n} \sum_i \min\Bigg( \frac{\pi(a_i^t \mid o_i^t; \theta)}{\pi(a_i^t \mid o_i^t; \theta_{\text{old}})} A_i^t, $$
$$\text{clip}\bigg( \frac{\pi(a_i^t \mid o_i^t; \theta)}{\pi(a_i^t \mid o_i^t; \theta_{\text{old}})}, 1 - \varepsilon,\, 1 + \varepsilon \bigg) A_i^t \Bigg)$$

18:     Compute the entropy bonus

$$\mathcal{H}(\theta) = \tfrac{1}{|\mathcal{B}|} \sum_b \tfrac{1}{L^b} \sum_t \tfrac{1}{n} \sum_i \mathcal{H}_i(\theta)$$

19:     Compute the critic losses:

$$\mathcal{L}(\phi) = \tfrac{1}{|\mathcal{B}|} \sum_b \tfrac{1}{L^b} \sum_t \tfrac{1}{n} \sum_i \left( y_i^{t,b} - V_i(\mathbf{o}_i^{t,b}; \phi) \right)^2$$

20:     Update $\theta$ and $\phi$ using:

$$\mathcal{L}(\theta, \phi) = -\mathcal{L}(\theta) + \alpha^{critic}\mathcal{L}(\phi) - \alpha^{entropy}\mathcal{H}(\theta)$$

21: **end while**