

## A. Mathematical Formulation

Let  $\theta$  and  $\phi$  denote the parameters of the joint policies and the probabilistic classifier, respectively. Then, the whole learning process corresponds to the following bi-level optimization:

$$\begin{aligned} \max_{\theta \in \Theta} \quad & J(\theta, \phi^*(\theta)) \\ \text{s.t.} \quad & \phi^*(\theta) = \arg \min_{\phi' \in \Phi} \mathcal{L}(\phi', \theta), \end{aligned}$$

where  $J$  is the RL objective with intrinsic reward,  $\mathcal{L}$  is the loss function of the probabilistic classifier, and  $\phi$  is an implicit function of  $\theta$ . Therefore, to solve this optimization, we can iteratively update  $\theta$  by

$$\frac{dJ(\theta, \phi^*(\theta))}{d\theta} = \left. \frac{\partial J(\theta, \phi)}{\partial \theta} \right|_{\phi=\phi^*(\theta)} + \frac{d\phi^*(\theta)}{d\theta} \left. \frac{\partial J(\theta, \phi)}{\partial \phi} \right|_{\phi=\phi^*(\theta)}$$

where

$$\frac{d\phi^*(\theta)}{d\theta} = - \left( \frac{\partial^2 \mathcal{L}(\phi, \theta)}{\partial \phi \partial \phi^T} \right)^{-1} \left( \frac{\partial^2 \mathcal{L}(\phi, \theta)}{\partial \phi \partial \theta^T} \right) \Big|_{\phi=\phi^*(\theta)}$$

which is obtained by the implicit function theorem. In practice, the second-order term is neglected due to high computational complexity, without incurring significant performance drop, such as in meta-learning and GANs. Therefore, we can solve the bi-level optimization by the first-order approximation with iterative updates:

$$\begin{aligned} \phi_{k+1} &\approx \arg \min_{\phi} \mathcal{L}(\phi, \mathcal{B}_k) \\ \theta_{k+1} &= \theta_k + \zeta_k \nabla_{\theta} J(\theta, \phi_{k+1}). \end{aligned}$$

## B. Hyperparameters

The hyperparameters of EOI and the baselines in each scenario are summarized in Table 1. Since QMIX and MAAC are off-policy algorithms with replay buffer, we do not need to maintain the buffer  $\mathcal{B}$  but build the training data from the replay buffer  $\mathcal{D}$ . For EDTI, ROMA, and HC, we use their default settings.

Table 1. Hyperparameters

Hyperparameter	Pac-man	Windy Maze	Firefighters	Battle	10_vs_10
runs with different seeds	5	10	5	5	5
horizon ( $T$ )	30	15	20	100	100
discount ( $\gamma$ )		0.98		0.96	0.995
replay buffer size			$2 \times 10^4$		$1 \times 10^4$
actor learning rate		$1 \times 10^{-3}$		-	$3 \times 10^{-4}$
critic learning rate		$1 \times 10^{-4}$		-	$1 \times 10^{-4}$
QMIX learning rate			$1 \times 10^{-4}$		-
# MLP units			(128, 128)		
batch size			128		
MLP activation			ReLU		
optimizer			Adam		
$\phi$ learning rate		$1 \times 10^{-3}$		$1 \times 10^{-4}$	$1 \times 10^{-4}$
$\alpha$ in QMIX		0.05		0.02	-
$\alpha$ in MAAC		0.2		-	0.04
$\beta_1$		0.04		0.05	0.05
$\beta_2$		0.1		0.05	0.05
$\Delta t$			4		