

1 Checklist

2 1. For all authors...

3 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
4 contributions and scope? [Yes]

5 (b) Did you describe the limitations of your work? [Yes] See Sec. ??

6 (c) Did you discuss any potential negative societal impacts of your work? [N/A]

7 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
8 them? [Yes]

9 2. If you are including theoretical results...

10 (a) Did you state the full set of assumptions of all theoretical results? [N/A]

11 (b) Did you include complete proofs of all theoretical results? [N/A]

12 3. If you ran experiments...

13 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
14 mental results (either in the supplemental material or as a URL)? [Yes]

15 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
16 were chosen)? [Yes]

17 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
18 ments multiple times)? [Yes]

19 (d) Did you include the total amount of compute and the type of resources used (e.g., type
20 of GPUs, internal cluster, or cloud provider)? [Yes]

21 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

22 (a) If your work uses existing assets, did you cite the creators? [Yes] We use the code
23 for running the experimental environments, and we reference them in Sec. ?? and
24 appendix.

25 (b) Did you mention the license of the assets? [Yes]

26 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]

27

28 (d) Did you discuss whether and how consent was obtained from people whose data you’re
29 using/curating? [N/A]

30 (e) Did you discuss whether the data you are using/curating contains personally identifiable
31 information or offensive content? [N/A]

32 5. If you used crowdsourcing or conducted research with human subjects...

33 (a) Did you include the full text of instructions given to participants and screenshots, if
34 applicable? [N/A]

35 (b) Did you describe any potential participant risks, with links to Institutional Review
36 Board (IRB) approvals, if applicable? [N/A]

37 (c) Did you include the estimated hourly wage paid to participants and the total amount
38 spent on participant compensation? [N/A]

39 A Environment illustrations and descriptions

40 Pac-Man [2] is a mixed cooperative-competitive maze game with one pac-man player and several
41 ghost players (Figure 1). We consider three pac-man scenarios containing two scenarios (OpenClassic
42 (Figure 1 (a)) and MediumClassic (Figure 1 (b))) with two ghost players and one pac-man player and
43 the complex scenario (Figure 1 (c)) with four ghost players and one pac-man player. The pac-man
44 player’s goal is to eat as many pills (denoted as white circles in the grids) as possible and avoid
45 the pursuit of ghost players. For ghost players, they aim to capture the pac-man player as soon as
46 possible. In our settings, we aim to control ghost players and the pac-man player is the opponent

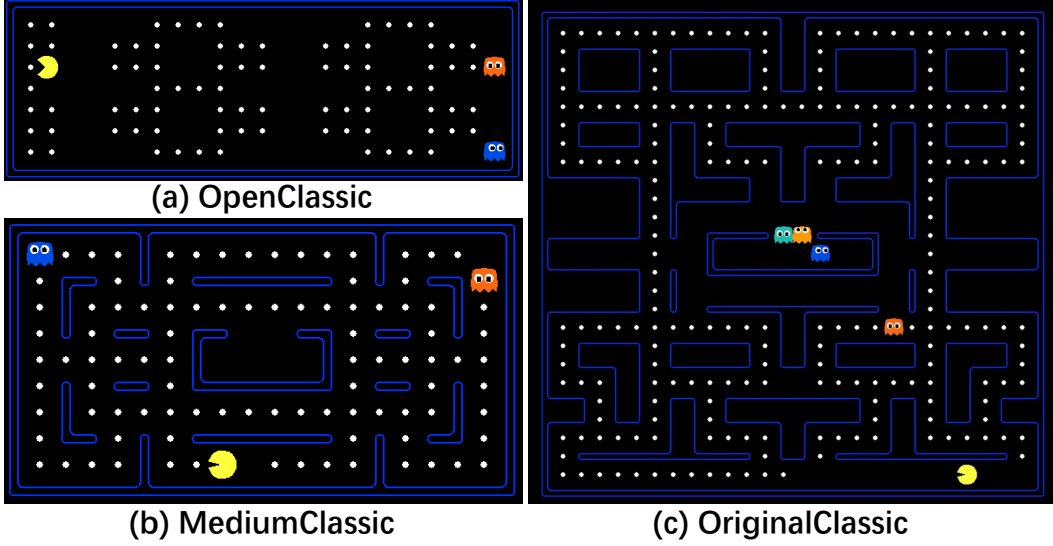


Figure 1: Pac-Man.

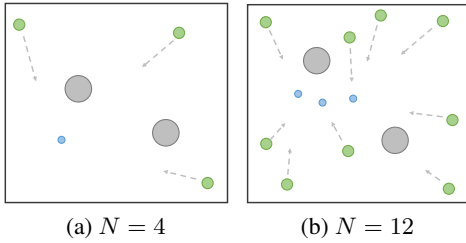


Figure 2: Predator-prey.

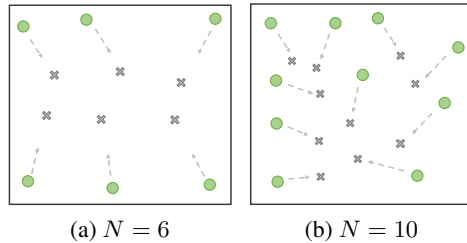


Figure 3: Cooperative Navigation.

47 controlled by a well pre-trained PPO policy. The game ends when one ghost catches the pac-man
 48 player or the episode exceeds 100 steps. Each ghost player receives -0.01 penalty for each step and
 49 $+5$ reward for catching the pac-man player.

50 MPE [1] is a multiagent particle world with continuous observation and discrete action space. We
 51 choose two scenarios of MPE: predator-prey (Figure 2), and cooperative navigation (Figure 3). The
 52 predator-prey contains three (nine) agents (green) which are slower and want to catch one (three)
 53 adversaries (blue) (rewarded $+10$ by each hit). Adversaries are faster and want to avoid being hit by
 54 the other three (nine) agents. Obstacles (grey) block the way. The cooperative navigation contains
 55 six (ten) agents (green), and six (ten) corresponding landmarks (cross). Agents are penalized with a
 56 reward of -1 if they collide with other agents. Thus, agents have to learn to cover all the landmarks
 57 while avoiding collisions. At each step, each agent receives a reward of the negative value of the
 58 distance between the nearest landmark and itself. Both games end when exceeding 100 steps.

59 State Description

60 **Pac-Man** The layout size of two scenarios are 25×9 (OpenClassic), 20×11 (MediumClassic) and
 61 28×27 (OriginalClassic) respectively. The observation of each ghost player contains its position,
 62 the position of its teammate, walls, pills, and the pac-man, which is encoded as a one-hot vector.
 63 The input of the network is a 68-dimension in OpenClassic, 62-dimension in MediumClassic and
 64 111-dimension in OriginalClassic.

65 **MPE** The observation of each agent contains its velocity, position, and the relative distance between
 66 landmarks, blocks, and other agents, which is composed of 18-dimension in predator-prey with four

67 agents (36-dimension with twelve agents), 36-dimension with six agents (60-dimension with ten
 68 agents) as the network input.

69 B Network structure and parameter settings

70 The experiments are conducted on a device with CPU of 64 cores, GPU of RTX2080TI and 256G
 71 Memory.

72 **Network Structure** Here we provide the network structure for PPO, MADDPG, QMIX and MAPTF
 73 respectively. 1) PPO: for each agent i , the actor network has two fully-connected hidden layers both
 74 with 64 hidden units, the output layer is a fully-connected layer that outputs the action probabilities
 75 for all actions; the critic network contains two fully-connected hidden layers both with 64 hidden
 76 units and a fully-connected output layer with a single output: the state value;

77 2) MADDPG: the actor network has two fully-connected hidden layers, one with 128 hidden units,
 78 the second layer with 64 hidden units; the output layer is a fully-connected layer that outputs one
 79 single action; the critic network contains two fully-connected hidden layers, one with 128 hidden
 80 units, the second layer with 64 hidden units; and a fully-connected output layer with a single output:
 81 the state-action value;

82 3) QMIX: for each agent i , the Q network has two fully-connected hidden layers, both with 128
 83 hidden units; the output layer is a fully-connected layer that outputs the Q-values for all actions; the
 84 mixing network contains two hypernetworks with 128 hidden units a mixing layer with 32 hidden
 85 units; and a fully-connected output layer with a single output: the joint state-action value;

4) SRO network structure is provided in Figure 4.

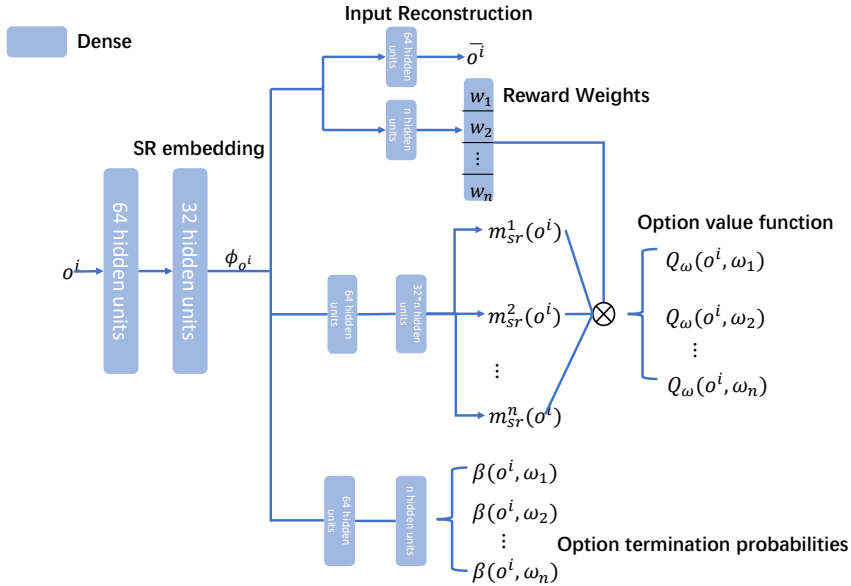


Figure 4: Network structures.

86

87 Parameter Settings

88 Here we provide the hyperparameters for MAPTF, DVM as well as three baselines, PPO, MADDPG
 89 and QMIX shown in Table 1, 2 and 3 respectively.

Table 1: Hyperparameters for all methods based on PPO.

Hyperparameter	Value
Learning rate	$3e - 4$
Length of trajectory segment T	32
Gradient norm clip λ	0.2
Optimizer	Adam
Discount factor γ	0.99
Batch size B of the option module	32
Replay memory size	$1e5$
Learning rate	$1e - 5$
μ	$5e - 4$
ξ	$5e - 3$
Action-selector	ϵ -greedy
ϵ -start	1.0
ϵ -finish	0.05
ϵ anneal time	$5e4$ step
target-update-interval τ	1000
distillation-interval for DVM	$2e5$ step
distillation-iteration for DVM	2048 step

Table 2: Hyperparameters for all methods based on MADDPG.

Hyperparameter	Value
Learning rate	$1e - 2$
Batch size	1024
Optimizer	Adam
Discount factor γ	0.99
Batch size B of the option module	32
Replay memory size	$1e5$
Learning rate	$1e - 5$
μ	$5e - 4$
ξ	$5e - 3$
Action-selector	ϵ -greedy
ϵ -start	1.0
ϵ -finish	0.05
ϵ anneal time	$5e4$ step
target-update-interval τ	1000

90 References

- 91 [1] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent
 92 actor-critic for mixed cooperative-competitive environments. In *Proceedings of NeurIPS*, pages
 93 6379–6390, 2017.
- 94 [2] Tycho van der Ouderaa. Deep reinforcement learning in pac-man. 2016.

Table 3: Hyperparameters for all methods based on QMIX.

Hyperparameter	Value
Learning rate	$3e - 4$
Batch size	64
Optimizer	Adam
Discount factor γ	0.99
ϵ -start	1.0
ϵ -finish	0.05
ϵ anneal time	$5e3$ step
Batch size B of the option module	32
Replay memory size	$1e5$
Learning rate	$1e - 5$
μ	$5e - 4$
ξ	$5e - 3$
Action-selector	ϵ -greedy
ϵ -start	1.0
ϵ -finish	0.05
ϵ anneal time	$5e4$ step
target-update-interval τ	1000