
MatrixWorld: A Pursuit-Evasion Platform for Safe Multi-agent Coordination and Autocurricula

Lijun Sun^{*†}

11860004@mail.sustech.edu.cn

Yu-Cheng Chang[†]

Yu-Cheng.Chang@uts.edu.au

Chao Lyu[‡]

lyuchao@swu.edu.cn

Chin-Teng Lin^{†§}

Chin-Teng.Lin@uts.edu.au

Yuhui Shi^{*§}

shiyh@sustech.edu.cn

Abstract

Multi-agent reinforcement learning (MARL) achieves encouraging performance in solving complex tasks. However, the safety of MARL policies is one critical concern that impedes their real-world applications. Popular multi-agent benchmarks focus on diverse tasks yet provide limited safety support. Therefore, this work proposes a safety-constrained multi-agent environment: MatrixWorld⁵, based on the general pursuit-evasion game. Particularly, a safety-constrained multi-agent action execution model is proposed for the software implementation of safe multi-agent environments based on diverse safety definitions. It (1) extends the vertex conflict among homogeneous / cooperative agents to heterogeneous / adversarial settings, and (2) proposes three types of resolutions for each type of conflict, aiming at providing rational and unbiased feedback for safe MARL. Besides, MatrixWorld is also a lightweight co-evolution framework for the learning of pursuit tasks, evasion tasks, or both, where more pursuit-evasion variants can be designed based on different practical meanings of safety. As a brief survey, we review and analyze the co-evolution mechanism in the multi-agent setting, which clearly reveals its relationships with autocurricula, self-play, arms races, and adversarial learning. Thus, MatrixWorld can also serve as the first environment for autocurricula research, where ideas can be quickly verified and well understood.

1 Introduction

In thriving multi-agent reinforcement learning (MARL) research, increasingly complex multi-agent behaviors have emerged in increasingly complicated environments. However, in current practice, MARL is still weak and challenging in guaranteeing the safety of multi-agent coordination [36, 21], such as in autonomous driving [35]. To date, the safe MARL has not gone beyond toy domains such as grid worlds. Therefore, this paper revisits the significance of discrete pursuit-evasion game and introduces MatrixWorld, a new safety-constrained multi-agent pursuit-evasion platform.

Here, the safety concept is defined in terms of collision avoidance. For instance, collisions may stem from conflicts of interest during multi-agent cooperation, the resolution of which determines the essential contribution of a coordination algorithm [5, 6]. When collisions occur, the multi-agent software environment should deal with these conflicts correctly by presenting reasonable and practical

^{*}Department of Computer Science and Engineering, Southern University of Science and Technology, China

[†]Faculty of Engineering and Information Technology, University of Technology Sydney, Australia

[‡]College of Computer and Information Science, Southwest University, China

[§]Corresponding author.

⁵All codes are available at <https://github.com/LijunSun90/MatrixWorld>.

collision results and blaming the responsible one through e.g., rewards, based on which algorithms can learn properly. However, as indicated by Terry et al. [32], the real collision resolution mechanism implemented by MARL environments may be biased and lead to unexpected outcomes for the stochastic game (SG), where multi-agent actions should be executed simultaneously. This may result from the bugs and complexities of codes when resolving the race conditions.

Popular multi-agent grid-world benchmarks provide limited safety support, based on which general MARL studies rarely report their safety performances even though they actually matter in associated applications. For example, in the Pursuit and Battle tasks of MAgent [37], rewards are given in terms of the attacking action. However, what if two agents (of the same team) intend to go to the same position? The collision resolution mechanism is transparent without any reward feedback. Similar situations apply to the predator-prey task of multi-agent particle environments (MPE) [14], the Pursuit task of PettingZoo [32], etc. In these cases, any biased or impractical collision resolution mechanisms with transparent information to the coordination algorithm will impede its safety learning. Moreover, in practice, negative rewards for collisions are insufficient for guaranteeing the safety of MARL policies, and significant efforts are expected in open safe MARL research.

Another safety limitation of popular MARL grid-world environments is that their collision resolutions are not diverse enough to satisfy various safety requirements in applications. This is because they focus more on providing diverse multi-agent behavior scenarios than diverse safety circumstances. For example, in the same pursuit-evasion framework, if agents are drones, any collisions will lead both to crashes (death and disappearance). However, if pursuers are predators and evaders are prey, predators will survive and prey will die after their collisions since they are unequal. Furthermore, if pursuers are mobile agents while evaders are stationary targets, pursuers and evaders can collide while collisions between pursuers are not allowed, as in the multi-agent path finding (MAPF) problem. These various safety definitions may bring different inequalities to the learning processes of agents and thus their resultant behaviors.

In addition, apart from the safety issue, we also consider the topic of autotutorials [12] in adversarial multi-agent research. As a kind of automatic curriculum learning method, its related works include co-evolution, self-play, arms races, adversarial learning, etc. These studies have achieved inspiring results in many fields, such as the genetic adversarial networks (GAN) [10], the Grandmaster level agent in StarCraft II [33], agents with human-level performance in Quake III Arena [11], and emerged multi-agent strategic tool use and coordination in hide-and-seek [1]. However, many adversarial benchmarks are too heavy, making it computationally costly to test ideas and tune hyperparameters, such as the above video games. Thus far, few studies have comprehensively discussed the relationships among these similar concepts, which may confuse new researchers in this field.

Motivated by the above concerns, we propose a new safety-constrained multi-agent environment: MatrixWorld. In particular, the pursuit-evasion game is selected as the basic environmental framework since it is general enough to be cooperative, competitive, or mixed. Moreover, in the context of multi-agent co-evolution, a pursuit-evasion game may be the most intuitive framework due to its biological inspirations. The main contributions of this work are as follows.

- For safe multi-agent coordination and MARL research, we propose a safety-constrained multi-agent action execution model to overcome limited safety support of popular MARL environments. To this end, we systematically analyze the agent-environment interaction models, the types of collisions and their resolutions in the adversarial multi-agent setting (Section 3).
- For open autotutorials research, we present several pursuit-evasion variants as lightweight co-evolution frameworks. As a brief survey, we review and analyze the co-evolution mechanism in the multi-agent setting, which clearly reveals its relationships with autotutorials, self-play, arms races, and adversarial learning (Section 2, 3).

2 Brief survey on co-evolution, autotutorials, and arms races

In this section, we present a brief survey and clarify the relationships between several similar concepts: co-evolution, self-play, autotutorials, adaptation, arms races, and adversarial learning, as shown in Figure 1. We aim to better understand and highlight the role of co-evolution in multi-agent settings, more accurately use the related terminologies, and discuss some common misleading expectations and magics, especially for researchers who are new to this field. Finally, we propose that MatrixWorld

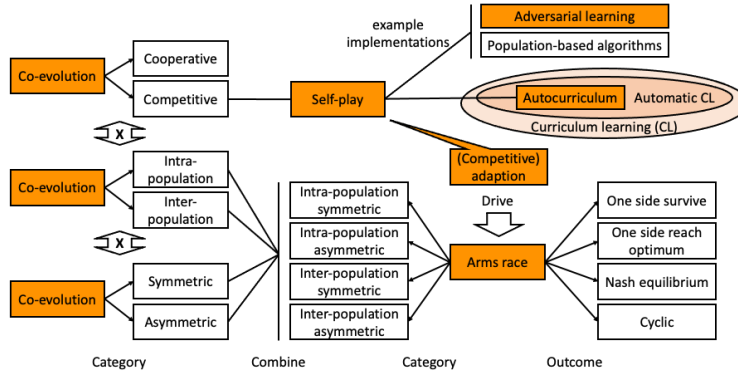


Figure 1: Relationships between co-evolution, self-play, autocurricula, arms races, and adversarial learning.

can serve as the first environment for autocurriculum research, where ideas can be quickly verified and well understood.

2.1 Co-evolution

Co-evolution is a natural phenomenon that witnesses the evolutionary progress and co-adaptations of, such as biological species. As an evolutionary and learning mechanism, its appealing properties, including autonomous co-adaptation and arms races, the natural modeling of coupled relationships, and performance improvement, contribute to its sustained success in diverse research fields and applications.

The most distinct feature of co-evolution is that the fitness evaluation of one individual depends on other individuals, where each individual is a separately evolving or learning agent. Co-evolution happens when adaptations occur between different agents, such as pursuers and evaders, environmental configurations and the agents therein. Its effectiveness is promoted by iterative coordinated learning among diverse individuals, which may induce a steadier and more guaranteed improvement. Therefore, co-evolution is more akin to a framework that can integrate different evolutionary computations (ECs), reinforcement learning (RL) algorithms, etc.

The available co-evolution methods can be broadly classified into three categories. First, there are competitive and cooperative co-evolution, which describe the relationships between individuals. However, they do not determine the number of populations that co-evolve. Individuals in one population can be competitive, while individuals from multiple populations may be cooperative. The number of co-evolving populations defines the second category: intra-population and inter-population co-evolution methods. Third, co-evolution can be classified into symmetric and asymmetric contests. Symmetric co-evolution is for homogeneous agents in the sense that they learn similar skills for the same tasks. In contrast, asymmetric co-evolution is for heterogeneous agents that target different skills for different, perhaps competitive, tasks.

2.2 Self-play

Self-play is a kind of competitive co-evolution process. It generally refers to the competitive interactions of multiple agents but can also occur between different "minds" (policy models) of the same physical agent, such as Alice and Bob in the research of Sukhbaatar et al. [29]. Its idea is that learning progress is only achieved by the interplay or interaction between agents without external interference. For example, Jaderberg et al. [11] trained symmetric adversarial agents in a competitive game (capture the flag) by letting the agents in a population play with each other. Baker et al. [1] trained asymmetric adversarial agents in another competitive game: hide-and-seek through multi-agent competitions.

2.3 Adaptation and arms races

Arms races between two rival sides are common in nature, such as those between predators and prey, and between parasites and hosts. Biologically, Dawkins et al. [7] discussed an arm race as an evolutionary process of reciprocal counter-adaptations and the resultant challenges faced by two sides. Based on this definition, the most crucial feature for identifying an arms race in a co-evolution process is the driving force derived from mutual counter-adaptations. However, note that, adaptations and challenges do not necessarily mean more complex, stronger, or more general skills, which will be further explored in Section 2.4.

In addition, Dawkins et al. [7] proposed a two-way classification paradigms for arms races, i.e., symmetric or asymmetric, and interspecific or intraspecific, which form four possible combinations with biological examples. This is also consistent with the categories of co-evolution. Moreover, they summarized four classes of arms race endpoints: one side wins by driving the other to extinction, one side wins by first reaching its optimum, equilibrium endpoints, and cyclic endings. One of the interesting discussions provided by Dawkins et al. [7] indicates that the average success of one side, such as the capture rate of predators, does not necessarily increase with the evolutionary time in an arms race. This is because both rival sides are improving. Then, the question of how to evaluate an arms race is progressing is actually the same as asking how the arms race will terminate, as described by the four possible outcomes of arms races.

2.4 Curriculum learning (CL), automatic CL, and autotutorials

Curriculum learning (CL) [3] refers to a training strategy that learns from a sequence of related tasks organized with increasing difficulty, which generally achieves better results faster than normal learning methods. Rather than manually designing a set of tasks and the time to switch between them, automatic CL automates part or all of the training process [34]. In particular, the autotutorial concept proposed by Leibo et al. [12] specifically refers to automatic CL in multi-agent settings, where the sequence of tasks or challenges is self-generated from mutual counter-adaptations in multi-agent interactions, i.e., an arms race. In this sense, autotutorial can be seen as equivalent to self-play.

However, based on both biological observations [7] and artificial intelligence (AI) investigations [12], it is noted that challenges do not definitely mean more difficult or complex tasks, and adaptations to these challenges do not necessarily involve increasing strong abilities. Even the commonly adopted difficulty ranking of tasks, i.e., from easy to hard, in curriculum learning is not guaranteed to be optimal in all cases [34]. The conditions for obtaining different arms race outcomes have attracted researchers' attention. Particularly, people are more interested in generating stronger models, such as agents with more complex behaviors. Compared with various optimal rankings of task difficulties, the conclusions on how to avoid policy cycles are more consistent. It has been shown that the policy cycling issue may occur in both symmetric and asymmetric arms races, for example, in a symmetric intra-population adversarial game of rock-paper-scissors [12] where the same agent strategy can play two different roles, and an asymmetric game [19]. Nolfi et al. [19] investigated the policy cycling problem in the co-evolution and proposed that it could be reduced by evaluating the current policy with all past policies, which successfully drives the evolutionary process to an arms race with increasing complexity. The same conclusion was also identified by Leibo et al. [12] and Vinyals et al. [33]: self-play algorithms that do forgetting past policies provide a way to break out the policy cycling problem. In addition, Leibo et al. [12] pointed out that the ceiling of the intelligence that can be achieved by an arms race is determined by the "environment's carrying capacity", or the task definition that the autotutorial is trying to solve. Their research actually states that optima are defined by the problem itself. For example, in a pursuit-evasion game without tools, novel tool usages will never emerge.

2.5 Adversarial learning

Adversarial learning (AL) is a training framework that optimizes adversarial models with conflicting objectives and is typically implemented in an alternative learning mode. Different from minimax search algorithms [26] in which an agent optimizes itself against an optimal adversary, AL can be more general. Subsequently, AL is not the only way to co-evolve adversarial agents and is more of an example implementation of self-play. For example, in the symmetric adversarial game "capture the

flag" [11], population-based RL was used to train a population of agents with similar skills that could play two adversarial game roles.

In addition, one main application of AL is to target an ultimate protagonist model rather than an adversary model, such as the generator model in generative adversarial networks (GANs) [10] and primary agents in the robust adversarial reinforcement learning (RARL) [23, 15, 22]. Last, note that not all adversarial work refers to co-evolutionary adversarial learning, and they may be merely normal learning approaches in adversarial settings. In these works, only one side in the adversarial setting is intentionally trained to be an opponent of the other, such as the adversarial policies in RL [9] and adversarial examples.

2.6 MatrixWorld: A lightweight co-evolution environment

Biologically, predators pursue their prey for dinner, while prey evade predators for life, and their co-evolutionary arms race and co-adaptations never stop. In AI, the co-evolutionary pursuit-evasion has been studied since 1994 [25, 18, 17, 16], which is slightly later than the research on pursuit or evasion alone, which has taken place since 1953 [13]. Compared with complex 3D first-person video games such as Quake III Arena [11], real-time strategy games such as StarCraft II [27], or adversarial games that need millions of episodes for arms races [1], MatrixWorld is a lightweight and simple co-evolution framework based on the pursuit-evasion game.

As discussed above, research on co-evolution is still open and it is supposed to be promising for automatically generating more complex agent behaviors and intelligence from the arms races in multi-agent interactions. We propose that MatrixWorld can serve as the first environment for quickly verifying and effectively understanding research ideas in autotutorials.

Table 1: Multi-agent-environment interaction models for distributed adversarial multi-agent systems. Task environment properties [26]: **Single-agent vs. multi-agent**: whether the other agent(s) optimizes some objectives that depend on the current agent. **Cooperative vs. non-cooperative**: whether all agents share a common goal. **Static vs. dynamic**: whether the environmental state or the agent’s observation changes if the agent does nothing regardless of the flying time. **Deterministic vs. nondeterministic (stochastic)**: whether “the next state s' is completely determined by the current state s and the action a executed by the agent(s)”, i.e., whether the transition function $P(s'|s, a) = 1$ or not.

Game of	Two-swarm turn taking	Agent-agent turn taking	Agent-environment interaction model	Environment model & properties
Single-agent system				Markov decision process
Adversarial multi-agent system (e.g., pursuit-evasion)	X	X		Stochastic / Markov game - multi-agent - non-cooperative - dynamic - stochastic
	✓	X		Two-swarm turn-taking game - multi-agent - non-cooperative - dynamic - stochastic
	✓	✓		Agent environment cycle game - multi-agent - non-cooperative - dynamic - stochastic

3 MatrixWorld: Safety-constrained multi-agent pursuit-evasion games

In this section, we introduce the propose MatrixWorld, a safety-constrained adversarial multi-agent environment based on co-evolutionary pursuit-evasion games. In particular, we propose a safety-constrained multi-agent action execution model for the software implementation of safe multi-agent environments. It determines the execution of agents’ actions (the agent-environment interaction model in Section 3.1), their results (the collision resolution mechanism in Section 3.2), and who should be responsible, i.e., the evidence to give right rewards, etc. In addition, we provide safety related information in the API for safe MARL research.

3.1 Multi-agent-environment interaction model

For a single-agent system, the typical agent-environment interaction (AEI) loop is that: at each time step, (1) an agent receives an observation from the environment, (2) based on the observation the agent (policy) makes its decision and outputs an action, (3) the action is executed in the environment, and (4) the environment changes by responding to the action and other factors [31].

For adversarial multi-agent settings, e.g., pursuit-evasion games, their agent-environment interaction models can be categorized based on two dimensions: whether agents take turns in a swarm-by-swarm manner, and whether agents take turns in an agent-by-agent manner, as shown in Table 1. Therefore, the third row in the table is a strict stochastic game where all agents concurrently observe, make decisions, and execute actions. The fourth row in the table is a two-swarm turn-taking game and a one-swarm stochastic game for each swarm. Similar to the two-player turn-taking game, a two-swarm turn-taking game can be defined as a game in which two swarms of agents take turns, i.e., swarm-by-swarm, observing, making decisions, and executing actions until the end of the game. Finally, the last row in Table 1 is a multi-agent turn-taking game, rather than a single-agent system, because the other agents also optimize some objectives in terms of the current agent in an AEI loop [26].

3.2 Safety-constrained collision resolution mechanism

There are four types of conflicts: (1) vertex conflict, (2) edge conflict, (3) following conflict, and (4) cycle conflict in MAPF [8, 28], as shown in the first row of Figure 2. Since the following and cycle conflicts do not actually cause collisions if all agents move simultaneously in a stochastic game, we allow these two “conflicts” as the literature works. Since the vertex and edge conflicts will physically cause the collision of agents, we forbid these two conflicts as the convention of MAPF. However, the first difference is that the vertex conflict type in MAPF is for homogeneous / cooperative agents while our work extends it to heterogeneous / adversarial settings by further classifying it to three types of collisions based on the agent type (see the second row of Figure 2). The second difference is that rare works systematically summarize the environment execution outcome for each forbidden conflict type while our work proposes three types of resolutions based on the safety definition (see the third row of Figure 2).

In detail, we implement the safety constraints in MAS by designing a collision resolution mechanism based on meaningful and practical safety definitions. As shown in Figure 2 and 3, collisions in the general multi-agent setting are classified into three types: (1) agent-obstacle collisions, (2) cooperative agent-agent collisions, and (3) competitive agent-adversary collisions. Agent-obstacle collisions are common in single and multiple agent systems, where only one of the two colliding objects is the decision-maker that should be responsible. Collisions between cooperative agents are not hostile and arise from conflicts of interest during cooperation, where both sides are accountable. In contrast, competitive agent-adversary collisions may be intentional from one side and unexpected from the other, where the reward strategy depends on the given case. Furthermore, we consider three outcomes for each collision: (1) reach (target) and disappear, (2) reach (target) and alive, and (3) bounce back (stay put) and alive. These collision types, outcomes, and resultant reward strategies form the collision resolution mechanism that satisfies the safety requirements of different scenarios in Table 2.

3.3 Pursuit-evasion game variants

Pursuit-evasion game variants can be defined based on three dimensions: (1) the agent-environment interaction model (Section 3.1), (2) the collision resolution mechanism (Section 3.2), and (3) the capture behavior. In view of the literature conventions and real-world applications, we propose

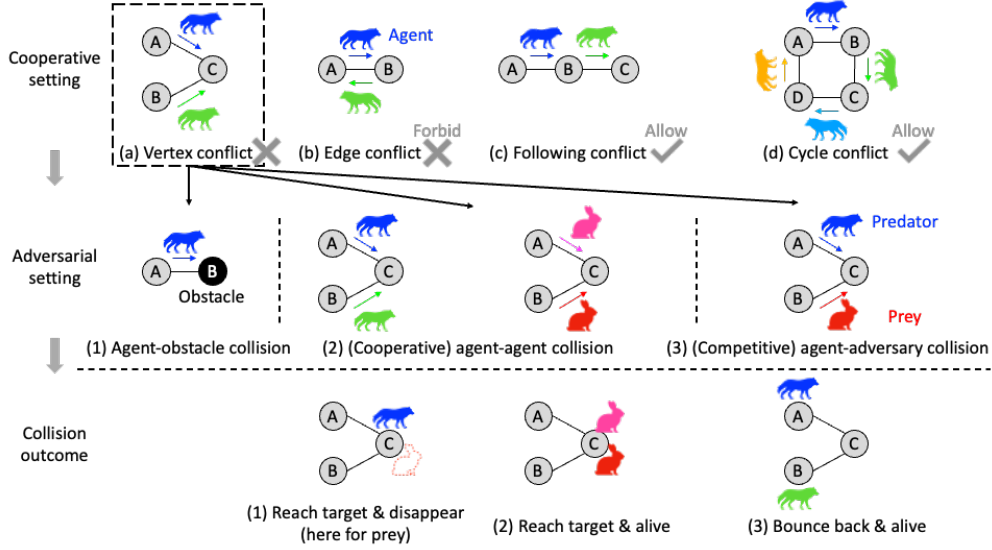


Figure 2: Illustration of the types of collisions and outcomes in general multi-agent interactions. First row: four types of conflicts for homogeneous / cooperative agents from [8, 28]. Second row: the extension of vertex conflict to heterogeneous / adversarial settings, i.e., the collision types. Third row: the types of outcomes for each collision.

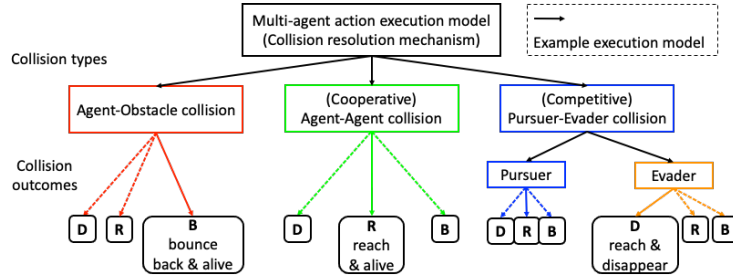


Figure 3: Safety-constrained multi-agent collision resolution mechanism for the multi-agent environment modeled by stochastic game.

Table 2: Rationality and example applications of various collision outcomes in the multi-agent setting. **D**: reach (target) and disappear. **R**: reach (target) and alive. **B**: bounce back (stay put) and alive.

Collision type	Agent type	Outcome			Rationality (example application)
		D	R	B	
Agent-obstacle	Agent	✓	-	-	Vehicle crash.
		-	✓	-	Allowed in some MARL environments.
		-	-	✓	Practical in many real-world applications.
Agent-agent	Cooperative	✓	-	-	Vehicle crash.
		-	✓	-	Allowed in some MARL environments.
		-	-	✓	Bumper cars.
Agent-adversary	Evader	✓	-	-	Vehicle crash.
		✓	-	-	
	Pursuer	✓	-	-	Predator eats prey.
		-	✓	-	
	Pursuer	-	-	✓	Pursuer is stronger.
		-	✓	-	
Pursuer	-	-	✓	Bumper cars.	
	-	-	✓		

nine pursuit-evasion tasks, i.e., the -D, -R, -B, -O, -S, -SB, -SD, -SDB, and -TO variants of the pursuit-evasion task in Table 3.

We mainly consider the first multi-agent-environment interaction model in Table 1 for the stochastic game due to its generality in modeling multi-agent games, based on which most of the pursuit-evasion variants are designed. In addition, the last Pursuit-Evasion-TO task is based on the two-swarm turn-taking and agent-agent turn-taking model due to its advantages in theory analysis and the past sustained research [4].

For cooperative agent-agent collisions, typical MARL tasks allow agents to collide with each other without fatal consequences, i.e., “reach and alive”. Therefore, we include this in -R, -O, -S, and -SD variants of the Pursuit-Evasion task, which may ease the learning of complex strategies. However, in some real-world applications, such collisions may be illegal, where the “bounce back and alive” setting is actually applied. Hence, the -B, -SB, and -SDB variants are designed accordingly to upgrade their safety.

For capture behaviors in the grid world, two main categories are considered. One is the occupation-based capture in which an evader is captured if one or more pursuers occupy its position. The other is the surrounding-based capture in which an evader is captured if enough pursuers and environment boundaries or obstacles encircle it such that it cannot move.

Table 3: **MatrixWorld**: Task definition. Collision type: **A-O**: agent-obstacle; **A-A**: cooperative agent-agent; **P-E**: pursuer-evader. Collision outcome: **D**: reach (target) and disappear; **R**: reach (target) and alive; **B**: bounce back (stay put) and alive.

Task environment (Pursuit-Evasion)	Agent-environment interaction		Collision resolution mechanism			Capture behavior	Related work and applications (examples)	
	Two-swarm turn-taking	Agent-agent turn-taking	A-O	A-A	P-E			
-D	Evader Pursuer	✗	✗	D	D	D	Occupation-based	Real-world drone and vehicle swarm
-R	Evader Pursuer	✗	✗	B	R	R	Occupation-based	Multi-agent path finding (MAPF) MPE predator-prey [14]
-B	Evader Pursuer	✗	✗	B	B	R	Occupation-based	Multi-agent path finding (MAPF)
-O	Evader Pursuer	✗	✗	B	R	D R	Occupation-based	Predator-prey; Spatial search
-S	Evader Pursuer	✗	✗	B	R	B R	Surrounding-based	Benda et al. [2]; Sun et al. [30]
-SB	Evader Pursuer	✗	✗	B	B	B R	Surrounding-based	Predator-prey
-SD	Evader Pursuer	✗	✗	B	R	B R	Surrounding-based disappear if captured	Predator-prey
-SDB	Evader Pursuer	✗	✗	B	B	B R	Surrounding-based disappear if captured	Predator-prey
-TO	Evader Pursuer	✓	✓	B	R	D R	Occupation-based	Cops-and-robbers or vertex-pursuit [24, 20, 4]; PettingZoo Pursuit [32]

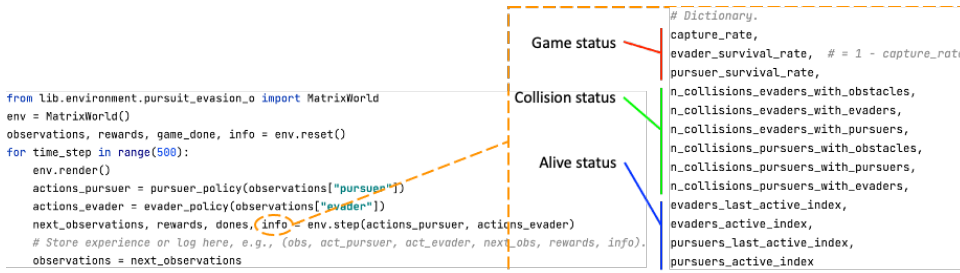


Figure 4: Illustration of basic usage of MatrixWorld.

3.4 API

The application programming interface (API) of MatrixWorld is designed based on the convention of RL community, as shown in Figure 4. The interface is the same while keeping the background multi-agent-environment interaction model in Table 1 transparent for users' convenience. In the dictionary information returned from the environmental step function, we provide (1) the game status data to monitor and compare the game progress regardless of the specific values in a reward structure; (2) the collision status data to give safety related feedback for performance evaluation, reward shaping, safety constraint construction, etc.; and (3) the "alive" status data of the agents to track their remains and population decreases due to death.

4 Conclusion

In this paper, we propose a safety-constrained multi-agent benchmark: MatrixWorld. Compared with popular MARL environments, it covers diverse safety definitions in various applications, by considering reasonable collision types, collision outcomes, and reward strategies in the collision resolution mechanism of multi-agent action execution model. In addition, safety-related environmental information, such as different kinds of collisions in each time step, is provided in the API for conducting open safe MARL research, e.g., reward shaping.

Furthermore, MatrixWorld is also a lightweight co-evolution framework for autotutor research based on the general pursuit-evasion game. By revisiting the concepts of co-evolution, autotutor, self-play, arms race, and adversarial learning, we clarify their relationships to achieve a better understanding and more accurate terminology usage.

Finally, this paper's work is only a start toward the safe MARL research. For the rewarding mechanism, this paper's work only determines who should be responsible, but limited work is done on how to optimally utilize the safety related feedback and shape the rewards, which are crucial for the success of safe MARL algorithms. Therefore, more work can be done based on the safety support of MatrixWorld and more multi-agent environments can be implemented based on the proposed general safety-constrained multi-agent action execution model. It is also suggested to explore autotutor in the lightweight MatrixWorld.

5 Acknowledgements

This work is partially supported by the Shenzhen Fundamental Research Program under Grant No. JCYJ20200109141235597, the National Science Foundation of China under Grant No. 61761136008, the Shenzhen Peacock Plan under Grant No. KQTD2016112514355531, the Program for Guangdong Introducing Innovative and Entrepreneurial Teams under Grant No. 2017ZT07X386, and the Australian Research Council (ARC) under Discovery Grant DP210101093 and DP220100803.

References

- [1] Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor Mordatch. Emergent tool use from multi-agent autotutor. In *International Conference on Learning Representations*, 2019.
- [2] M Benda, V Jagannathan, and R Dodhiawala. On optimal cooperation of knowledge sources-an empirical investigation. Technical report, BCS-G2010-28, Boeing Advanced Technology Center, Boeing Computing Services, Seattle, Washington, 1986.
- [3] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- [4] Anthony Bonato. *The game of cops and robbers on graphs*. American Mathematical Soc., 2011.
- [5] Craig Boutilier. Planning, learning and coordination in multiagent decision processes. In *TARK*, volume 96, pages 195–210. Citeseer, 1996.

- [6] Craig Boutilier. Sequential optimality and coordination in multiagent systems. In *IJCAI*, volume 99, pages 478–485, 1999.
- [7] Richard Dawkins and John Richard Krebs. Arms races between and within species. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 205(1161):489–511, 1979.
- [8] Jianqi Gao, Yanjie Li, Xinyi Li, Kejian Yan, Ke Lin, and Xinyu Wu. A review of graph-based multi-agent pathfinding solvers: From classical to beyond classical. *Knowledge-Based Systems*, page 111121, 2023.
- [9] Adam Gleave, Michael Dennis, Cody Wild, Neel Kant, Sergey Levine, and Stuart Russell. Adversarial policies: Attacking deep reinforcement learning. *arXiv preprint arXiv:1905.10615*, 2019.
- [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- [11] Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al. Human-level performance in 3d multiplayer games with population-based reinforcement learning. *Science*, 364(6443):859–865, 2019. doi: 10.1126/science.aau6249.
- [12] Joel Z Leibo, Edward Hughes, Marc Lanctot, and Thore Graepel. Autocurricula and the emergence of innovation from social interaction: A manifesto for multi-agent intelligence research. *arXiv preprint arXiv:1903.00742*, 2019.
- [13] John Edensor Littlewood. *A mathematician’s miscellany*. Methuen & Co. Ltd., London, 1953.
- [14] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Neural Information Processing Systems (NIPS)*, 2017.
- [15] Xiaobai Ma, Katherine Driggs-Campbell, and Mykel J Kochenderfer. Improved robustness and safety for autonomous vehicle control with adversarial reinforcement learning. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 1665–1671. IEEE, 2018.
- [16] Pattie Maes, Maja J Mataric, Jean-Arcady Meyer, Jordan Pollack, and Stewart W Wilson. Co-evolution of pursuit and evasion ii: Simulation methods and results. 1996.
- [17] Geoffrey F Miller and Dave Cliff. *Co-evolution of pursuit and evasion I: Biological and game-theoretic foundations*. School of Cognitive and Computing Sciences, University of Sussex Brighton, 1994.
- [18] Geoffrey F Miller and Dave Cliff. Protean behavior in dynamic games: Arguments for the co-evolution of pursuit-evasion tactics. *From animals to animats*, 3:411–420, 1994.
- [19] Stefano Nolfi and Dario Floreano. How co-evolution can enhance the adaptive power of artificial evolution: Implications for evolutionary robotics. In *Evolutionary Robotics: First European Workshop, EvoRobot98 Paris, France, April 16–17, 1998 Proceedings 1*, pages 22–38. Springer, 1998.
- [20] Richard Nowakowski and Peter Winkler. Vertex-to-vertex pursuit in a graph. *Discrete Mathematics*, 43(2):235 – 239, 1983. ISSN 0012-365X. doi: [https://doi.org/10.1016/0012-365X\(83\)90160-7](https://doi.org/10.1016/0012-365X(83)90160-7).
- [21] Afshin Oroojlooy and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement learning. *Applied Intelligence*, pages 1–46, 2022.
- [22] Xinlei Pan, Daniel Seita, Yang Gao, and John Canny. Risk averse robust adversarial reinforcement learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8522–8528. IEEE, 2019.

- [23] Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial reinforcement learning. In *International Conference on Machine Learning*, pages 2817–2826. PMLR, 2017.
- [24] Alain Quilliot. *Jeux et pointes fixes sur les graphes*. PhD thesis, Université de Paris VI, 1978.
- [25] Craig W Reynolds. Competition, coevolution and the game of tag. In *Proceedings of the Fourth International Workshop on the Synthesis and Simulation of Living Systems*, pages 59–69, 1994.
- [26] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach, 4th Edition*. Pearson Education, 2021.
- [27] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pages 2186–2188, 2019.
- [28] Roni Stern, Nathan Sturtevant, Ariel Felner, Sven Koenig, Hang Ma, Thayne Walker, Jiaoyang Li, Dor Atzmon, Liron Cohen, TK Kumar, et al. Multi-agent pathfinding: Definitions, variants, and benchmarks. In *Proceedings of the International Symposium on Combinatorial Search*, volume 10, pages 151–158, 2019.
- [29] Sainbayar Sukhbaatar, Zeming Lin, Ilya Kostrikov, Gabriel Synnaeve, Arthur Szlam, and Rob Fergus. Intrinsic motivation and automatic curricula via asymmetric self-play. *arXiv preprint arXiv:1703.05407*, 2017.
- [30] Lijun Sun, Yu-Cheng Chang, Chao Lyu, Ye Shi, Yuhui Shi, and Chin-Teng Lin. Toward multi-target self-organizing pursuit in a partially observable markov game. *Information Sciences*, 648: 119475, 2023.
- [31] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [32] J Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Pettingzoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34:15032–15043, 2021.
- [33] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- [34] Xin Wang, Yudong Chen, and Wenwu Zhu. A survey on curriculum learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):4555–4576, 2021.
- [35] Cathy Wu, Abdul Rahman Kreidieh, Kanaad Parvate, Eugene Vinitsky, and Alexandre M Bayen. Flow: A modular learning framework for mixed autonomy traffic. *IEEE Transactions on Robotics*, 38(2):1270–1286, 2021.
- [36] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, pages 321–384, 2021.
- [37] Lianmin Zheng, Jiacheng Yang, Han Cai, Ming Zhou, Weinan Zhang, Jun Wang, and Yong Yu. Magent: A many-agent reinforcement learning platform for artificial collective intelligence. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), Apr. 2018.

A Experiments

In this section, through adversarial learning, we achieve various arms race outcomes of different co-evolution mechanisms in MatrixWorld. Based on experiments, arms races with steady and converging improvement are more practical for increasingly complex behaviors, while policy cycles between two rival sides are useful for producing diverse policies. In particular, we find that the passive (evasive) policy learning benefits more from co-evolution than active (pursuing) policy learning in an asymmetric adversarial game. An arms race can drive the passive policy to a higher level than that in normal RL. Finally, based on experiments in Pursuit-Evasion-S, we show the demanding safety assurance regarding the reward shaping or other techniques in safe MARL.

A.1 Adversarial learning algorithms for pursuit-evasion tasks

We particularly explore three representative co-evolution frameworks and their influences on the arms race process and outcomes.

- Pursuer specialist vs. evader specialist (Algorithm 1): Pursuers are trained to be specialized for the current evaders, and so are the evaders.
- Pursuer generalist vs. evader specialist (Algorithm 2): Pursuers are trained to be general, while the evaders are trained to be specialized to their opponents, which is an unfair training system but may be more practical in the real world. For example, police can learn from all past criminal data, while criminals may only access the current police data.
- Pursuer generalist vs. evader generalist (Algorithm 3): Pursuers are trained to be general and robust for all past evaders, and so are the evaders.

The alternative learning scheme of adversarial learning is applied that when one side learns the other side is fixed.

Algorithm 1: Adversarial learning: Pursuer specialist vs. evader specialist

```

1 Initialize pursuer model  $\theta_0^p$  and evader model  $\theta_0^e$ .
2 for generation  $k = 1 : N$  do
3    $\theta_{k,k^p=0}^p = \theta_{k-1}^p$ .
4   for  $k^p = 1 : N^p$  do
5     Collect experiences by  $\theta_{k,k^p-1}^p$  and  $\theta_{k-1}^e$ .
6     Pursuer model update  $\theta_{k,k^p-1}^p \rightarrow \theta_{k,k^p}^p$ .
7    $\theta_k^p = \theta_{k,N^p}^p$ .
8    $\theta_{k,k^e=0}^e = \theta_{k-1}^e$ .
9   for  $k^e = 1 : N^e$  do
10    Collect experiences by  $\theta_k^p$  and  $\theta_{k,k^e-1}^e$ .
11    Evader model update  $\theta_{k,k^e-1}^e \rightarrow \theta_{k,k^e}^e$ .
12   $\theta_k^e = \theta_{k,N^e}^e$ .
13 return  $\theta_N^p, \theta_N^e$ .
```

A.2 Experimental setup

The policy models for the pursuers and evaders are both two-layer ReLU multi-layer perceptrons (MLPs) with hidden sizes of [400, 300]. The actor-critic [31] algorithm is adopted in the centralized training and decentralized learning (CTDE) scheme, with the policy and value learning rates set to 3×10^{-4} and 10^{-3} , respectively. To stabilize the learning process, the actor model is trained 5 times, while value function is trained 1 time per epoch. $N = 30$ generations of co-evolution are conducted for Algorithms 1, 2, and 3, with $N_p = N_e = 400$ epochs in each generation. For the Pursuit-Evasion-O task, 8 pursuers compete with 30 evaders in 40×40 grid worlds, while for Pursuit-Evasion-S, 20 pursuers compete with 5 evaders since 4 pursuers are required for every evader in the surrounding-based capture paradigm. The reward functions are shown in Table 4.

Algorithm 2: Adversarial learning: Pursuer generalist vs. evader specialist

```
1 Initialize pursuer model  $\theta_0^p$  and evader model  $\theta_0^e$ .
2 for generation  $k = 1 : N$  do
3    $\theta_{k,k^p=0}^p = \theta_{k-1}^p$ .
4   for  $k^p = 1 : N^p$  do
5      $k^e = k^p \bmod k$ .
6     Collect experiences by  $\theta_{k,k^p-1}^p$  and  $\theta_{k^e}^e$ .
7     Pursuer model update  $\theta_{k,k^p-1}^p \rightarrow \theta_{k,k^p}^p$ .
8    $\theta_k^p = \theta_{k,N^p}^p$ .
9    $\theta_{k,k^e=0}^e = \theta_{k-1}^e$ .
10  for  $k^e = 1 : N^e$  do
11    Collect experiences by  $\theta_k^p$  and  $\theta_{k,k^e-1}^e$ .
12    Evader model update  $\theta_{k,k^e-1}^e \rightarrow \theta_{k,k^e}^e$ .
13   $\theta_k^e = \theta_{k,N^e}^e$ .
14  return  $\theta_N^p, \theta_N^e$ .
```

Algorithm 3: Adversarial learning: Pursuer generalist vs. evader generalist

```
1 Initialize pursuer model  $\theta_0^p$  and evader model  $\theta_0^e$ .
2 for generation  $k = 1 : N$  do
3    $\theta_{k,k^p=0}^p = \theta_{k-1}^p$ .
4   for  $k^p = 1 : N^p$  do
5      $k^e = k^p \bmod k$ .
6     Collect experiences by  $\theta_{k,k^p-1}^p$  and  $\theta_{k^e}^e$ .
7     Pursuer model update  $\theta_{k,k^p-1}^p \rightarrow \theta_{k,k^p}^p$ .
8    $\theta_k^p = \theta_{k,N^p}^p$ .
9    $\theta_{k,k^e=0}^e = \theta_{k-1}^e$ .
10  for  $k^e = 1 : N^e$  do
11     $k^p = k^e \bmod (k + 1)$ .
12    Collect experiences by  $\theta_{k^p}^p$  and  $\theta_{k,k^e-1}^e$ .
13    Evader model update  $\theta_{k,k^e-1}^e \rightarrow \theta_{k,k^e}^e$ .
14   $\theta_k^e = \theta_{k,N^e}^e$ .
15  return  $\theta_N^p, \theta_N^e$ .
```

Table 4: Reward function for all pursuit-evasion tasks. “-”: the same.

Pursuer		Evader	
Action	Reward	Action	Reward
Capture an evader	10	Being captured	-10
Neighbor an evader	0.1	Being neighbored	-0.1
Collide	-12	-	-12
Move before termination	-0.05	-	-0.05

Suggestion: Set the hyperparameters N_p and N_e to allow the (moving average) learning performance improve to some extent. In this way, an arms race can be observed in the adversarial learning; otherwise, neither rival side can learn effectively.

A.3 Autocurricula in co-evolutionary pursuit-evasion tasks

Adaptation and arms race occur between pursuers and evaders. As shown in the training performance presented in Figure 5, the learning of pursuers brings challenges to the evaders’ policies and the capture rate increases, and vice versa. The performance of the agents (pursuers and evaders) continues to improve with iterative generations. Their performance variance brought by policy adaptations tends to converge as the number of generations increases, but this is not the end of the arms race. When we observe many more generations than $N = 30$, say 100 generations, we find that the performance variance diverges again, which demonstrates the sustained adaptations during the arms race.

Policy cycles occur in the specialist-vs.-specialist framework, but they are avoided in the generalist-vs.-generalist scheme. As shown in Figure 6, we test the evolutionary performance (capture rate) of the evaders against the pursuers in each generation. The policy cycles are observed in Figure 6a, where both the pursuers and evaders learn only against their contemporary opponents. The evaders’ performances against the pursuers in a specific generation, i.e., one row in the figure, fluctuate periodically from generation to generation, or even decrease sometimes. This indicates that due to the mutual adaptations of pursuers and evaders, their skills periodically appear and disappear, which is called a policy cycle [19]. On the other hand, such policy cycles from counter-adaptations are useful for producing diverse policies with similar complexity levels. In contrast, when the pursuers are trained to be general in Algorithm 2, the policy cycle problem is less severe in Figure 6b. Furthermore, when both the pursuers and evaders learn against all past opponents, the skill of competing with a specific pursuer is preserved in later generations, which can be seen from the stable performance in Figure 6c after the evaders first learn against that pursuer. This indicates that learning against past opponents helps avoid forgetting already learned skills. This result is consistent with the idea of past literature works on eliminating policy cycles, such as [19, 33, 12].

An arms race is more useful for learning passive policies. In Figure 7, we test the generalization performance of the agent in each generation by averaging its performance against the opponents across all generations. At the same time, we compare it with three baselines.

- Baseline 1: the pursuer that is trained by competing with randomly walking evaders.
- Baseline 2: the evader that is trained by competing with randomly walking pursuers.
- Baseline 3: the evader that is trained by competing with well-learned pursuers. The well-learned pursuer is trained by competing with randomly walking evaders. That is, the baseline 3 is trained via manual curriculum learning.

The conclusion that an arms race is more useful for learning passive (evasive) policies can be drawn based on two consistent observations from Figure 7a to 7c. First, adversarial learned evaders have better generalization performance than the baseline evaders, while adversarial learned pursuers do not achieve significantly better generalization performance than the baseline pursuers. Second, it can be seen that the complexity of a passive evasive policy highly depends on opponents’ ability. When trained only against random pursuers, the evaders’ performance is hard to improve once it reaches a certain level, which is the worst in Figure 7. If trained against stronger pursuers, i.e., well-learned pursuer, their performance is significantly enhanced. The best evader policies are obtained by adversarial learning. This indicates that autocurriculum learning is achieved and increasingly strong pursuers drive the continuous learning of evaders.

A.4 General MARL in safe multi-agent coordination scenarios

Compared with Pursuit-Evasion-O, more multi-agent conflicts of interest may occur in Pursuit-Evasion-S since more explicit multi-agent coordination is required to capture each evader with four pursuers. Therefore, we demonstrate the difficulties faced by general MARL in guaranteeing safe multi-agent coordination with only negative rewards for collisions in Pursuit-Evasion-S. As shown in Figure 8, after training for 1000+ epochs, the reward and capture rates improve significantly. However, more conflicts of interest occur with increasingly efficient capture behavior, as shown in

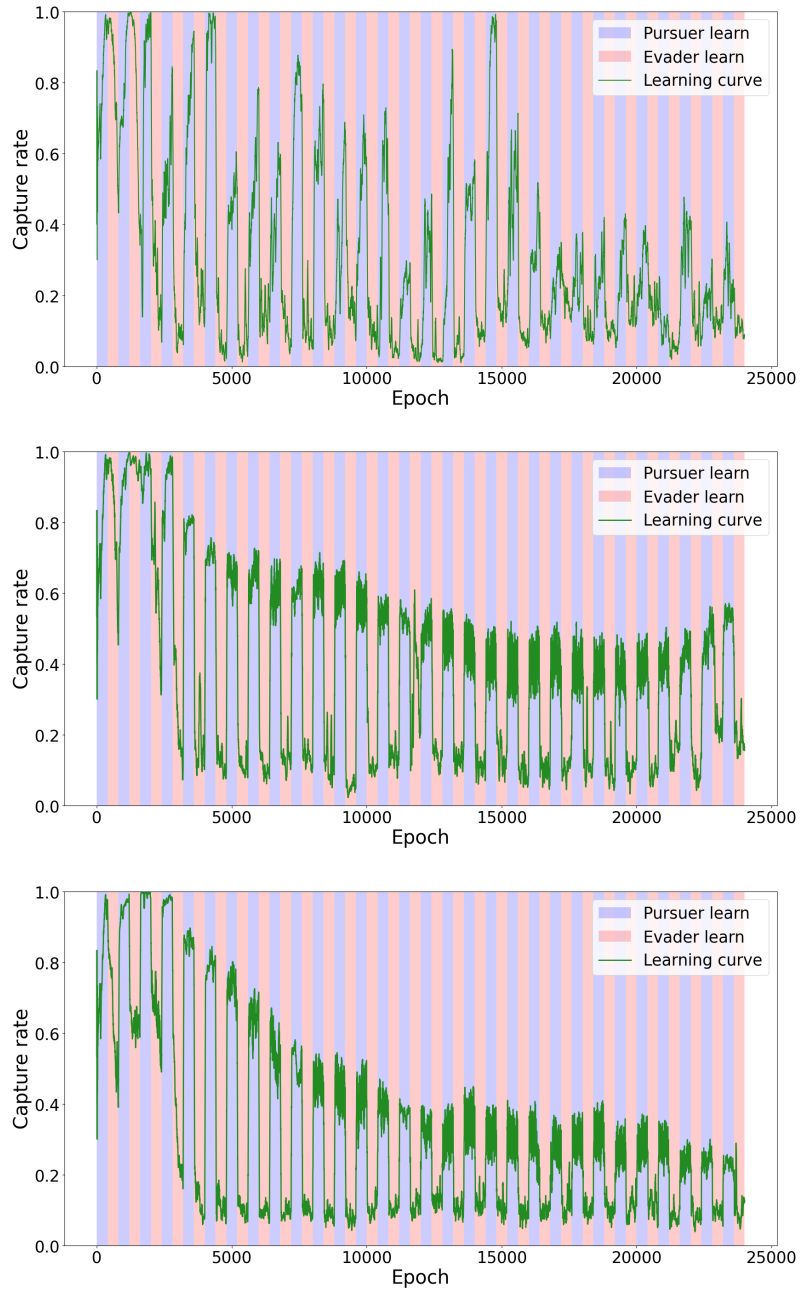
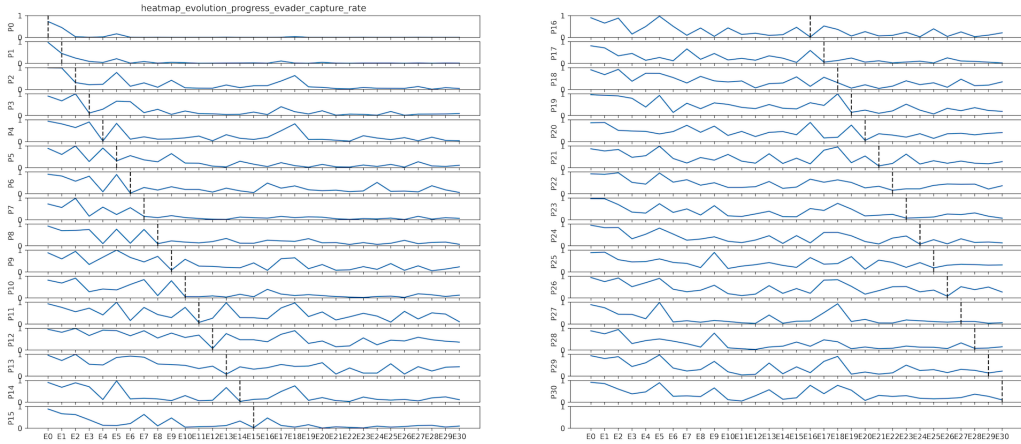
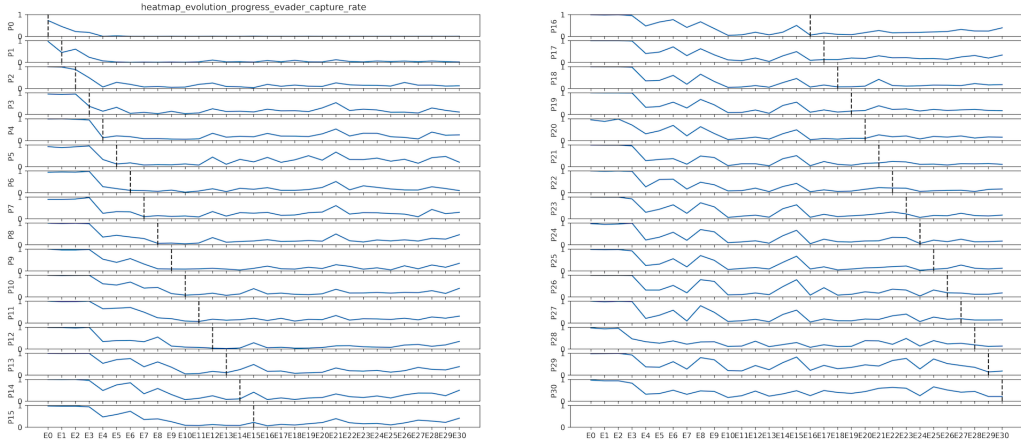


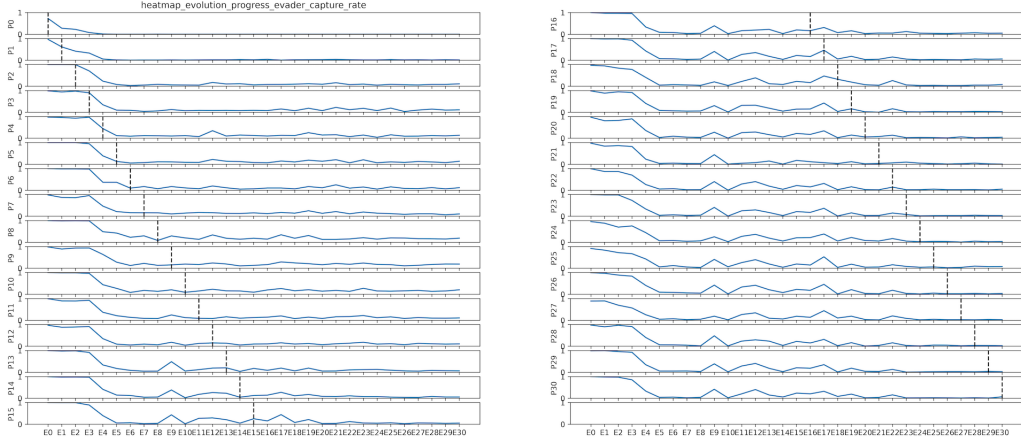
Figure 5: Training performance achieved for Pursuit-Evasion-O by Algorithms 1, 2, and 3 (from top to bottom). The curves are smoothed over 30 points.



(a) Pursuer specialist vs. evader specialist (Algorithm 1)

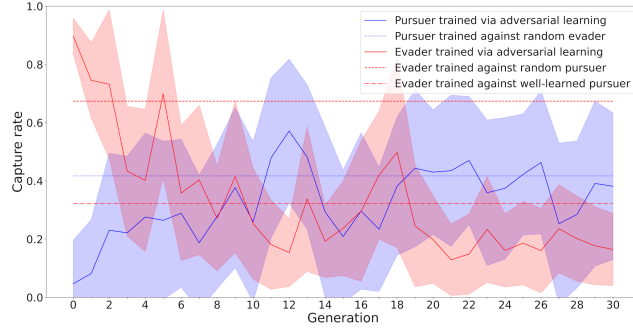


(b) Pursuer generalist vs. evader specialist (Algorithm 2)

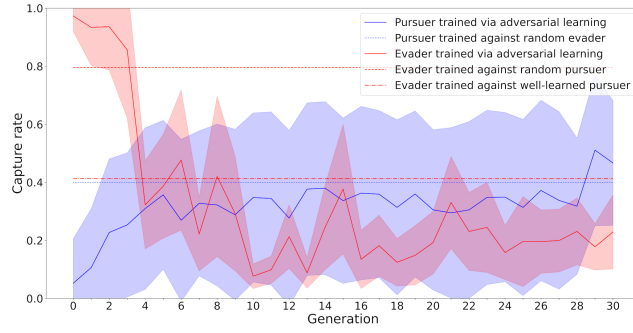


(c) Pursuer generalist vs. evader generalist (Algorithm 3)

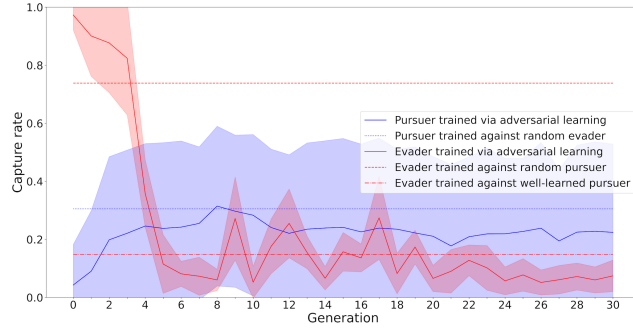
Figure 6: Evolutionary performance (capture rate) of evaders based on the testing performance achieved on Pursuit-Evasion-O, which is averaged over 10 independent runs. Horizontal axis: E0 - E30, the 30 generations of evaders. Vertical axis: P0 - P15 (left), P16 - P30 (right), the 30 generations of pursuers. E0 and P0 are the initial evaders and pursuers. A black dashed line indicates the associated pursuer and evader are in the same generation.



(a) Pursuer specialist vs. evader specialist (Algorithm 1)



(b) Pursuer generalist vs. evader specialist (Algorithm 2)



(c) Pursuer generalist vs. evader generalist (Algorithm 3)

Figure 7: Generalization performance achieved for Pursuit-Evasion-O.

the collision statuses. The multi-agent collisions do not absolutely vanish with a longer training time, especially in test scenarios, although the costs of collisions are larger than the benefits of capturing an evader in the reward structure (Table 4).

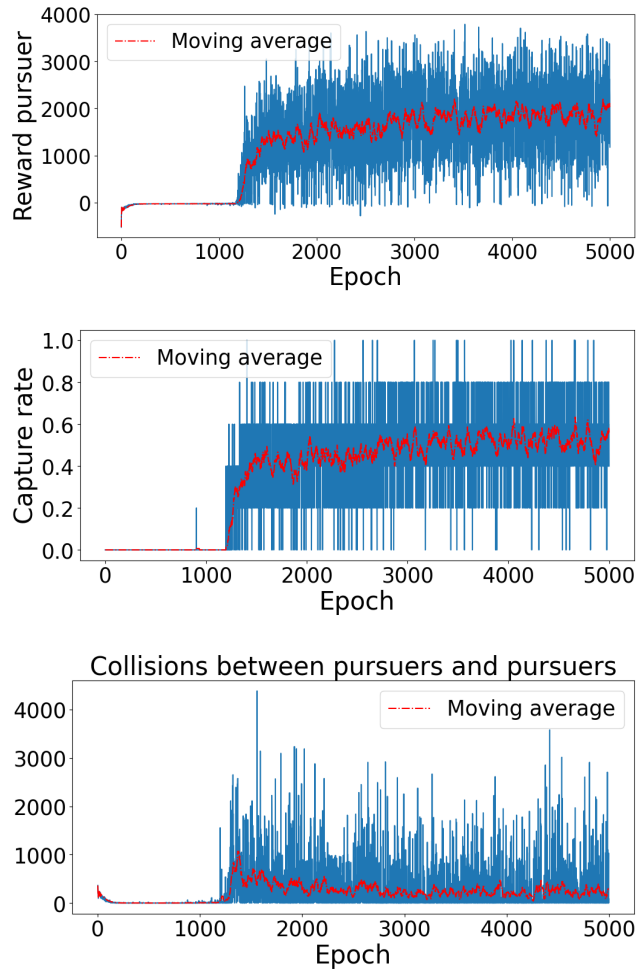


Figure 8: Training performance achieved for Pursuit-Evasion-S.