

# FIGHTLADDER: A BENCHMARK FOR COMPETITIVE MULTI-AGENT REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent advances in reinforcement learning (RL) heavily rely on a variety of well-designed benchmarks, which provide environmental platforms and consistent criteria to evaluate existing and novel algorithms. Specifically, in multi-agent RL (MARL), a plethora of benchmarks based on cooperative games have spurred the development of algorithms that improve the scalability of cooperative multi-agent systems. However, for the competitive setting, a lightweight and open-sourced benchmark with challenging gaming dynamics and visual inputs has not yet been established. In this work, we present FightLadder, a real-time fighting game platform, to empower competitive MARL research. Along with the platform, we provide implementations of state-of-the-art MARL algorithms for competitive games, as well as a set of evaluation metrics to characterize the performance and exploitability of agents. We demonstrate the feasibility of this platform by training a general agent that consistently defeats 12 built-in characters in single-player mode, and expose the difficulty of training a non-exploitable agent without human knowledge and demonstrations in two-player mode. FightLadder offers meticulously crafted environments to tackle essential challenges in competitive MARL research, heralding a new era of discovery and advancement.

## 1 INTRODUCTION

As an active branch of artificial intelligence (AI), deep reinforcement learning (DRL) has achieved huge success in recent years, especially in solving various well-designed games, which serve as a natural testbed to evaluate RL algorithms. Thanks to the growing computation power, the RL community has accomplished superhuman performances in a range of games, including video games like Atari 2600 (Bellemare et al., 2013; Mnih et al., 2013), Starcraft II (Vinyals et al., 2019) and Dota 2 (Berner et al., 2019), board games like chess (Schrittwieser et al., 2020) and Go (Silver et al., 2016), as well as card games like Texas hold'em (Brown & Sandholm, 2018; 2019), etc. Besides innovation on the algorithm side, such remarkable achievements also heavily rely on many platforms built upon these games that can unify evaluation protocol, compare state-of-the-art methods, and motivate better solutions. However, most existing game platforms focus on single-player or cooperative multi-player settings. A few platforms simulate competitive games with compact representations and relatively simple dynamics (e.g., board games), while others based on complex game engines require huge computational resources and human knowledge (e.g., StarCraft and Dota). To prosper advances in competitive multi-agent reinforcement learning (MARL) and transform game-theoretical results into practical applications, a fully competitive game platform that lies on the sweet spot among complexity, efficiency, and generality is urgently needed.

Games with more than one player are known to be hard to solve, due to the additional non-stationarity caused by other players. Among different game structures, fully competitive settings can be more difficult. People have a long history of playing fighting games since the early age of computers, as well as building computer players (CPU) to make the game more challenging and hence intriguing. Previous AI research has investigated the solutions of competitive games using RL, but mostly for small-scale games like Backgammon (Tesauro et al., 1995) or other board games (Schrittwieser et al., 2020; Brown & Sandholm, 2018; 2019). Moreover, this line of work mostly uses state vectors as inputs, which is much easier than directly learning from raw pixel inputs for AIs, which is common in most popular video games. As a result, we aim to build a platform for a series of fighting games, with image inputs and complex fighting dynamics, to serve as a challenging competitive multi-player platform for the broad AI research community.



Figure 1: FightLadder currently supports various cross-platform video fighting games: *Street Fighter II* (Genesis platform), *Street Fighter III* (Arcade platform), *Fatal Fury 2* (Genesis platform), *Mortal Kombat* (Genesis platform), and *The King of Fighters '97* (Neo Geo platform). Motion and attack action spaces are shown on the right.

Apart from the game platform, the evaluation criteria and benchmark results for certain game settings are essential for boosting the field. MARL has been greatly investigated in the past few years for solving multi-player games, from both theoretical and empirical perspectives. A large number of algorithms have been proposed according to specific settings (Sunehag et al., 2017; Yu et al., 2022; Lowe et al., 2017; Silver et al., 2018; Lanctot et al., 2017; Vinyals et al., 2019; Ding et al., 2022). Nonetheless, for competitive game settings, there is a lack of unified evaluation criteria with thorough comparisons among different approaches.

In this work, we present FightLadder, a competitive two-player games benchmark. Our contributions are three-fold: We build the FightLadder platform to support five two-player fighting games, with ease to extend to other games in the future. The games support various observation spaces involving rendered images. Based on prior work, we provide implementations of the most popular algorithms for solving these competitive games, including an AlphaStar league training algorithm (Vinyals et al., 2019) and policy space response oracle (Lanctot et al., 2017). Furthermore, a unified evaluation framework with *Elo rating* and *exploitability* tests are provided alongside the game platforms and algorithm library. We report experimental results using the above toolkits to serve as the baselines for two-player competitive game settings. We empirically demonstrate that *despite existing algorithms enabling constant improvement in terms of Elo score, the non-exploitable strategies for competitive two-player games with visual observations are hard to achieve with present algorithms*, thus raising a challenge for the research community.

## 2 RELATED WORK

**RL for Games.** RL has been widely adopted in a vast variety of popular games, from single-player video games like Atari 2600 (Bellemare et al., 2013; Mnih et al., 2013), VizDoom (Kempka et al., 2016), OpenAI Gym (Brockman et al., 2016), Retro (Nichol et al., 2018a), and Minecraft (Fan et al., 2022; Ding et al., 2023; Yuan et al., 2023) to multi-player games like StarCraft 2 (Vinyals et al., 2019) and Dota 2 (Berner et al., 2019), as well as two-player board games like Backgammon (Tesauro et al., 1995), Go (Silver et al., 2016; 2017), chess (Schrittwieser et al., 2020), multi-player poker games (Brown & Sandholm, 2018; 2019) and other card games (Zha et al., 2020). A few previous works have also explored learning policies in fighting games (Palmas, 2022; Go et al., 2023), but are limited to tackling one specific game with prior human knowledge or lack explicit criteria for two-player scenarios and the interface to integrate new fighting games as a comprehensive testbed, which is not aligned with our benchmark’s aim of motivating algorithms to solve general competitive MARL tasks.

**Multi-Agent Learning Environments.** Present multi-agent learning environments could be categorized into three types according to the payoff structure of the game: *fully cooperative*, *fully competitive*, and a *mixture* of both. Existing environments for *fully cooperative* games are designed for various scenarios, including simulated games like MAMuJoCo (Peng et al., 2021), card games like Hanabi (Bard et al., 2020), video games like small-scale StarCraft SMAC (Samvelyan et al., 2019) and Google Research Football (Kurach et al., 2020), as well as practical scenarios like Traffic Junction (Sukhbaatar et al., 2016) in a grid world, Flatland (Mohanty et al., 2020) for railway networks, network load balancing (Yao & Ding, 2022) and CityFlow (Zhang et al., 2019) for city traffic. On the other hand, the *fully*

*competitive* game setting is relatively underdeveloped, with few environments presented previously, such as Pommerman (Resnick et al., 2018) focusing on scenarios with low-dimensional discrete control and homogeneous agent characters. Environments supporting a *mixture* of cooperative and competitive games include MPE (Mordatch & Abbeel, 2018), MAgent (Zheng et al., 2018), Hide-and-Seek (Baker et al., 2019), DMLab2D (Beattie et al., 2020), Arena (Song et al., 2020), Smarts (Zhou et al., 2020), Neural MMO (Suarez et al., 2021), PettingZoo (Terry et al., 2021), MATE (Pan et al., 2022), etc.

**Multi-Agent Learning Baselines.** For solving multi-agent learning tasks, the research community has proposed algorithms and built libraries for ease of usage. PyMARL (Samvelyan et al., 2019) is an initial MARL library built for solving SMAC tasks, while PyMARL2 (Hu et al., 2021) extends PyMARL with QMIX (Rashid et al., 2020). EPyMARL (Papoudakis et al., 2020) is also an extension of PyMARL, as a unified library for cooperative games supporting different learning paradigms including centralized and decentralized learning, value decomposition, etc. MARLlib (Hu et al., 2023) includes major cooperative MARL algorithms like VDN (Sunehag et al., 2017), MAPPO (Yu et al., 2022), MADDPG (Lowe et al., 2017), etc. There some recent libraries Pantheonrl (Sarkar et al., 2022), MALib (Zhou et al., 2023), etc. These libraries mainly support MARL algorithms for cooperative games, lacking support for solving competitive games.

On the other hand, there is a line of research solving competitive games based on self-play (Silver et al., 2018), fictitious play (Brown, 1951), Nash Q-learning (Hu & Wellman, 2003; Ding et al., 2022), double oracle (McMahan et al., 2003), policy space response oracle (Lanctot et al., 2017) and league training (Vinyals et al., 2019). However, there is a lack of a unified framework to evaluate these algorithms efficiently on the same tasks, especially when combining these algorithms with deep RL.

To address the shortage of fully competitive environments, we focus on two-player zero-sum games and propose a platform for fighting-style fully competitive games, as well as the baseline implementation of those popular algorithms.

### 3 MULTI-AGENT REINFORCEMENT LEARNING

FightLadder is designed to motivate novel algorithms for fully competitive two-player games in the domains of MARL and game theory. Markov Games (MGs) (Shapley, 1953) generalize single-player Markov Decision Processes (MDPs) into multi-player settings. Each player has its own utility and optimizes its policy to maximize the utility. The two-player zero-sum setting in MG represents a competitive relationship between the two players. With a shaped dense reward, the games can be generalized to general-sum.

Specifically, we consider a finite-horizon two-player general-sum partially observable MG, denoted as POMG( $\mathcal{S}, \mathcal{O}, \mathcal{A}, \mathcal{B}, \mathbb{P}, \mathbb{O}, \{r_i\}_{i=1}^2, H$ ).  $\mathcal{S}$  is the state space, which can be partially observable for players and transformed through an observation emission function  $\mathbb{O}: \mathcal{S} \rightarrow \mathcal{O}$  to the observation space  $\mathcal{O}$  (e.g., video frames).  $\mathcal{A}$  and  $\mathcal{B}$  are action spaces for two players, respectively.  $\mathbb{P}(\cdot|s, a, b)$  is the state transition distribution,  $r_i: \mathcal{S} \times \mathcal{A} \times \mathcal{B} \rightarrow \mathbb{R}$  is the reward function for the  $i$ -th player. In the zero-sum setting, two reward functions satisfy the zero-sum payoff structure  $r_1 + r_2 = 0$ .  $H$  is the horizon length. We denote the policies of two players as  $\mu$  and  $\nu$ , respectively.  $V_i^{\mu, \nu}: \mathcal{S} \rightarrow \mathbb{R}$  represents the value function for player  $i$  evaluated with policies  $\mu$  and  $\nu$ , which can be expanded as the expected cumulative reward starting from the state  $s$ ,

$$V_i^{\mu, \nu}(s) := \mathbb{E}_{\mu, \nu} \left[ \sum_{h=1}^{\infty} r_i(s_h, a_h, b_h) \mid s_1 = s \right].$$

In zero-sum games, we have  $V_1^{\mu, \nu}(s) = -V_2^{\mu, \nu}(s), \forall s \in \mathcal{S}$  and define  $V^{\mu, \nu}(s) = V_1^{\mu, \nu}(s)$  for simplicity.

**Definition 3.1 (Best Response).** For any policy of the first player  $\mu$ , there exists a *best response* (BR) against it from the second player, which is a policy  $\nu^\dagger(\mu)$  satisfying  $V_{2,h}^{\mu, \nu^\dagger(\mu)}(s) = \max_{\nu} V_{2,h}^{\mu, \nu}(s)$  for any  $(s, h) \in \mathcal{S} \times [H]$ . We denote  $V_{2,h}^{\mu, \dagger} := V_{2,h}^{\mu, \nu^\dagger(\mu)}$  for simplification.  $V_{2,h}^{\mu, \nu}(s)$  is the value function of the second player. BR against the second player can be defined similarly.

**Definition 3.2 (Nash Equilibrium).** The *Nash equilibrium* (NE) in zero-sum setting is defined as a pair of policies  $(\mu^*, \nu^*)$  satisfying the following minimax equation:

$$\max_{\mu} \min_{\nu} V^{\mu, \nu}(s) = V^{\mu^*, \nu^*}(s) = \min_{\nu} \max_{\mu} V^{\mu, \nu}(s).$$

**Definition 3.3** (Exploitability). The exploitability for a policy  $\mu$  of the first player is defined as  $V_2^{\mu, \dagger}(s_1) - V_2^{\mu^*, \nu^*}(s_1)$ , i.e., the value of its BR policy  $\nu^\dagger(\mu)$  or the suboptimality gap from the NE value. The exploitability of the other side policy  $\nu$  can be defined accordingly.

Note that NE strategies will always lead to zero exploitability, thus approaching the non-exploitable strategies is a reasonable pursuit for the game.

## 4 FIGHTLADDER

In this section, we present technical details of FightLadder. In the following part, we first introduce different game settings of FightLadder, followed by elaborating elements of MGs corresponding to the environment, and conclude with highlighting features of our benchmark.

### 4.1 SCENARIOS

FightLadder provides a flexible interface between modern game emulators (Murphy, 2013; Nichol et al., 2018b) and algorithm developers. Thanks to its flexibility, FightLadder can support a wide range of classical fighting games over the past decades<sup>1</sup>, including Street Fighter, Mortal Kombat, Fatal Fury, and The King of Fighters, some of which are still very popular nowadays. Figure 1 shows screenshots of several fighting games provided by FightLadder. With this diverse set of supported games, we can benchmark algorithms on various fighting scenarios differing in backgrounds, characters, and moving dynamics, which can further motivate novel algorithms that are general rather than overfitting to one specific game. For better readability and clarity, we would use Street Fighter as an example for illustration and evaluation in the rest of the paper. The other fighting games are very similar, and readers could refer to Appendix A.2 for more details. We name each scenario in the form  $[game\ alias]_{[character\ left]}_{vs}_{[character\ right]}$ , for example  $sf\_ryu\_vs\_ryu$  in Street Fighter.

While FightLadder mainly focuses on the competitive two-player setting, the nature of fighting games allows it to be seamlessly deployed to the single-player scenario where the agent’s task is to compete against a built-in game AI (e.g.,  $sf\_ryu\_vs\_ryu(cpu)$ ). Under this single-player setting, users have the freedom to choose characters and set up the difficulty of the scripted AI opponent. Moreover, our benchmark also supports training in a much more challenging full-game scenario (e.g.,  $sf\_ryu\_full\_game$ ), where the agent needs to defeat all 12 characters controlled by computers with the difficulty progressively increasing. As we shall see in later experiments, this scenario could also serve as a sanity check for our baseline algorithms to see whether they could learn effective behaviors from the environment.

### 4.2 STATE AND OBSERVATIONS

We define the state space  $\mathcal{S}$  as the complete set of attributes stored in the game emulator after each step of action. Same as human players, the agent is not allowed to access the underlying full state but can only access the observation space  $\mathcal{O}$  of pixels, which forms a  $128 \times 100$  RGB image corresponding to the rendered screen. This image includes the position and movement of both sides of the players, as well as the hit-point bar and the round timer on the top of the screen. At every step, a configurable number of images are stacked as the input of the agent.

While we use pixels as default observations, we also provide an interface for users to access additional information about the game status, including position, hit-point, and exact countdown number for agents on both sides. Users can leverage these attributes to better understand the agent’s behavior or augment feature representations. More details are provided in Appendix A.2.

### 4.3 ACTION SPACE

In fighting games, two players share the same action space  $\mathcal{A}$ . The native *human action space*  $\mathcal{A}_{\text{human}}$  is designed to mimic the joystick control of arcade games, which is a 12-dimensional binary space ( $\{ 'B', 'A', 'MODE', 'START', 'UP', 'DOWN', 'LEFT', 'RIGHT', 'C', 'Y', 'X', 'Z' \}$ ) with each dimension representing a button being pressed or not. Note that due to the nature of fighting game

<sup>1</sup>FightLadder is compatible with Gym Retro and MAMEToolkit, hence it should support most of the fighting games running on Arcade or emulators that support the Libretro API.

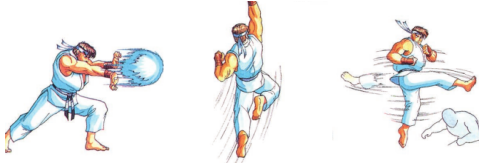


Figure 2: Example of special moves for character Ryu in StreetFighter II (left to right): Fireball, Dragon Punch, Hurricane Kick.

Table 1: FPS and memory usage of several open-sourced platforms.

Environment	Speed (FPS)	Memory (MB)
FightLadder (Ours)	1935.76	195.46
SMACv2	146.72	876.96
PettingZoo Atari	6268.18	32.13
DMLab2D	1144.27	47.41

engines, this space contains many redundant actions that are invalid, for instance, moving in opposite directions or moving and attacking at the same moment. To filter out these redundant actions and to construct a more structured space, we develop a categorical *transformed action space*  $\mathcal{A}_{\text{trans}}$  through an encoding function  $F: \mathcal{A}_{\text{human}} \rightarrow \mathcal{A}_{\text{trans}}$ . Specifically,  $\mathcal{A}_{\text{trans}}$  is the joint set of a direction move set  $\mathcal{A}_{\text{motion}} = \{\text{defense, forward, jump, crouch, back flip, front flip, offensive crouch, defensive crouch}\}$  and an attack move set  $\mathcal{A}_{\text{attack}} = \{\text{light punch, medium punch, hard punch, light kick, medium kick, hard kick}\}$ , as shown in Figure 1. Each action will remain a number of frames according to users’ configuration. The games also have special techniques called *close attack*, i.e., Throws and Holds, which can be applied in certain regions near the opponent.

In addition to the standard move set, one signature element of fighting games is *special moves*, a kind of powerful attack or maneuver that requires the player to follow a specific action sequence (i.e., sequential keys combination, or combination of key holding and key pressing), as an example depicted in Figure 2. These moves usually have special properties (e.g., invincibility frames, larger coverage, etc.) and play a critical role in the strategy and depth of the game. They are especially useful for higher levels of play, from which players could create complex combos and outperform opponents. However, we observe that learning to perform special moves from scratch can be challenging to baseline algorithms, as it requires the agent to memorize frames and actions in previous steps and accurately perform the next action in the action sequence of special moves. Moreover, the special moves can be different from character to character, which increases the difficulty of the game. Therefore, to alleviate this challenge, we also include hard-coded special move lists as one part of the action space so that the agent can directly access special moves with one single action.

#### 4.4 REWARDS

**Sparse Reward.** Both sides of the agents are to maximize their win rate for each round of the game. The *sparse reward*  $r_{\text{sparse}}$  assigns +1 for the winner and -1 for the loser at the end of each episode.

**Win Rate.** For two players  $A$  and  $B$ , policy  $\pi_A$  winning against policy  $\pi_B$  can be defined as a reward relationship  $r_{\text{sparse}}^A(\pi_A, \pi_B) > r_{\text{sparse}}^B(\pi_A, \pi_B)$  in a single match, with  $r_{\text{sparse}}^A$  and  $r_{\text{sparse}}^B$  as the sparse reward for players  $A$  and  $B$  in the zero-sum setting. The win rate is defined as the probability of the winning as:  $p(\pi_A \succ \pi_B)$ .

**Shaped Dense Reward.** While sparse reward is straightforward for evaluation, we discover that baseline algorithms could not effectively learn to behave well from such a sparse signal. To address this issue, we introduce a *shaped dense reward*  $r_{\text{dense}}$  for training, which is a weighted sum of the hit-point damage inflicted by the agent on the opponent and the damage it receives, together with a bonus (penalty) for winning (losing) the game. Specific format of this reward refers to Appendix A.1. The dense reward  $r_{\text{dense}}$  is chosen to coincide with the win rate of the policy, such that  $\pi_A \succ \pi_B$  will always lead to  $r_{\text{dense}}^A(\pi_A, \pi_B) > r_{\text{dense}}^B(\pi_A, \pi_B)$  in expectation. The dense reward also offers some flexibility, that the user can control the agent’s aggressiveness by configuring the weighing scales in the reward function.

#### 4.5 FEATURES

We remark on the following features of the proposed benchmark that could benefit MARL research.

**Rich Strategy Space.** One key feature of our benchmark is the rich strategy space as the nature of fighting games, which is particularly beneficial to the development of game-theoretical algorithms. To name a few, fighting games require players to consider **(a) character diversity**: each character

has a unique skill set with different strengths and weaknesses, so one needs to master the strategy and counter-strategy of all possible opponents, and even reason how to select and order characters when they have the freedom to do so; **(b) complexity of mechanics:** fighting games are designed with sophisticated mechanics such as invincibility frame, hitboxes, and combo systems, which are challenging for micro-management of characters; and **(c) adversarial opponents:** opponents may progressively adapt their policies to players’ policies, thus finding non-exploitable policies is crucial in mastering fighting games.

**Computational Efficiency.** FightLadder also enjoys efficient computation for its usage, and the comparison with several other popular game environments is shown in Table 1. The frame rate is 13 times faster than SMACv2, with one-fourth usage of the memory. While it is less efficient than FightLadder is the PettingZoo Atari, it provides more game complexity. The balance of complexity and low computational cost is important for evaluating algorithms at scale.

**Fidelity and Popularity.** FightLadder allows testing agents in full-length fighting games with an interface similar to human perception, thus providing a high-fidelity evaluation of competitive RL algorithms. Moreover, fighting games have been gaining popularity since they were released, making it easier to test the learned RL agents against human expert players.

**Open-Source, Compatibility, and Flexibility.** FightLadder is designed for the broad RL research community, so we make efforts to improve the ease of usage and make it accessible<sup>2</sup> to all potential users. It is compatible with the Gym (Brockman et al., 2016) interface so that users can leverage off-the-shelf RL algorithms implementation. Furthermore, it is extremely flexible as users can configure environments’ specifications, add new fighting games, and save or load the environment state at any moment during the game.

## 5 EVALUATION METRICS

**Versus Built-In Game AIs.** Directly competing with the built-in AIs of the games provides a straightforward way of measuring policy performance. Typically, fighting games offer a hierarchical structure of levels, enabling players to adjust the difficulty setting (for example, Street Fighter features eight distinct levels). This structure allows for the empirical evaluation of the policy against the game’s scripted AI at varying levels of challenge. It is important to acknowledge, however, that the limitations associated with hard-coded adversaries restrict the extent to which this metric can accurately reflect the policy’s real capability. For brevity, we shall refer to such agents as CPU.

**Elo Ratings.** The skills of agents can be ranked through the FIDE rating system (Elo & Sloan, 1978), which is an incremental learning system that increases the Elo of winners and decreases the Elo of losers. The larger the difference in Elo between players  $A$  and  $B$ , the higher the probability that the player with the higher Elo,  $A$ , beats the player with the lower Elo,  $B$ . The Elo score calculation takes the following procedures:

First, the probability of player  $A$  winning is estimated with,

$$p_A := p(\pi_A \succ \pi_B) = (1.0 + 10^{\frac{\text{Elo}_B - \text{Elo}_A}{400}})^{-1}.$$

Then the Elo rating for player  $A$  as  $\text{Elo}_A$  will be updated with following formula:

$$\text{Elo}_A = \text{Elo}_A + k \cdot (\mathbb{1}[\text{winner} = A] - p_A),$$

where  $k$  is a constant of update rate. The update is symmetric for player  $B$ , as well as any other player in the ranking system.

**Versus AI Exploiters.** As discussed in Section 3, exploitability (as Definition 3.3) measures the distance of a policy to the Nash equilibrium of the game. Specifically, the exploitability of a policy  $\mu$  is measured by the win rate of its BR policy  $\nu^\dagger(\mu)$  against  $\mu$ , since  $V^{\mu^*, \nu^*}(s_1) = 0$  for symmetric zero-sum game and  $V_2^{\mu, \dagger}(s_1) = 1 \cdot p(\nu \succ \mu) + 0 \cdot p(\nu \preceq \mu) = p(\nu \succ \mu)$  for sparse reward setting. In practice, we can use any single-agent deep RL algorithm as an exploiter to approximately learn the BR policy  $\nu^\dagger(\mu)$ . For fair comparisons, we should use one consistent exploiter (same RL algorithm with same configurations) to evaluate the exploitability of different baselines.

<sup>2</sup>We will open-source FightLadder once the paper is published.

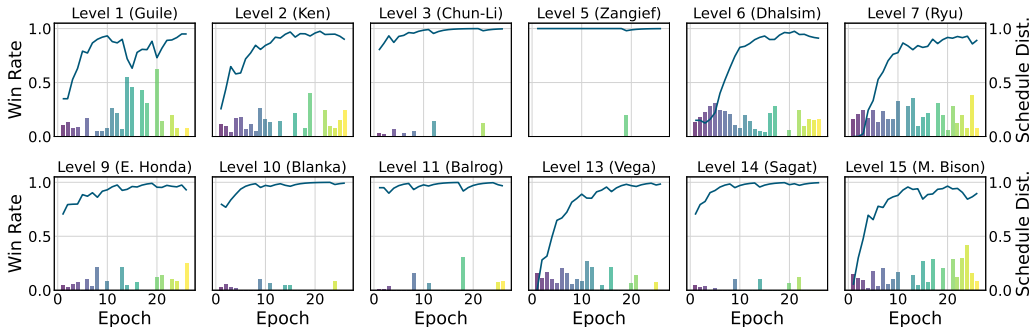


Figure 3: The win rate curves and the scheduling distribution bar plot in *sf\_ryu\_full\_game* via the proposed PPO with curriculum learning. Opponents of different characters are marked with different levels. Levels 4, 8, and 12 are omitted as they are bonus levels without fighting.

**Versus Human Players.** While Definition 3.3 is a general metric to measure exploitability, it may be limited to the capability of deep RL algorithms in usage. Therefore, we also provide an interface for human players such that they can play with any learned model with convenience. This feature will show the strengths and weaknesses of agents directly and visibly, and motivate developers to improve their algorithms to be more non-exploitable in general. Given the remarkable success of modern RL algorithms outperforming expert human players in various video games (Mnih et al., 2013; Vinyals et al., 2019; Berner et al., 2019), we believe that FightLadder will emerge as a promising platform for the broad competitive MARL community and researchers will eventually build AI agents that could beat world champions in a much richer set of strategic games with significantly less engineering efforts.

## 6 FIGHTLADDER-BASELINES

For the convenience of the community to evaluate existing methods and new algorithms on FightLadder platform, we open-source the implementation of several state-of-the-art (SOTA) competitive MARL algorithms, including independent learning (de Witt et al., 2020), two-timescale learning (Daskalakis et al., 2020), fictitious self-play (Heinrich et al., 2015), policy-space response oracle (Lanctot et al., 2017) and league training (Vinyals et al., 2019). Our codebase supports decentralized learning across multiple GPUs, and it is built upon Stable-Baselines3 (Raffin et al., 2021) so that users can leverage off-the-shelf implementations of RL algorithms. We choose proximal policy optimization (PPO) (Schulman et al., 2017) as the backbone policy optimization algorithm in our experiments. More details of baseline algorithms refer to Appendix B.

## 7 RESULTS

In this section, we provide benchmark results on a selected game in FightLadder—the Street Fighter. We aim to answer the following questions through our benchmark: **(a)** Can existing RL algorithms solve the full video game in the single-player scenario? **(b)** How does the performance of state-of-the-art baseline algorithms in the two-player competitive setting compare? and **(c)** Does multi-agent training help to improve the non-exploitability?

### 7.1 SINGLE-PLAYER FULL VIDEO GAME

To answer question **(a)**, we evaluate PPO’s performance in the scenario *sf\_ryu\_full\_game* as a feasibility check. As mentioned in Section 4, this scenario requires the agent to learn a generalizable policy to compete against all different characters with increasing difficulty levels. *Curriculum learning* is applied to train the policy from easy to hard cases. Furthermore, to improve learning efficiency we develop a curriculum scheduler for opponent sampling to match with the learner after each epoch. More specifically, for the current learner  $L$  with policy  $\pi_L$ , we sample its opponent  $C$  from the entire character set  $\mathcal{C}$ , with the following inverse-weight scheduling distribution:

$$C \sim \Delta(\mathcal{C}) \propto 1 - p(\pi_L \succ \pi_C),$$

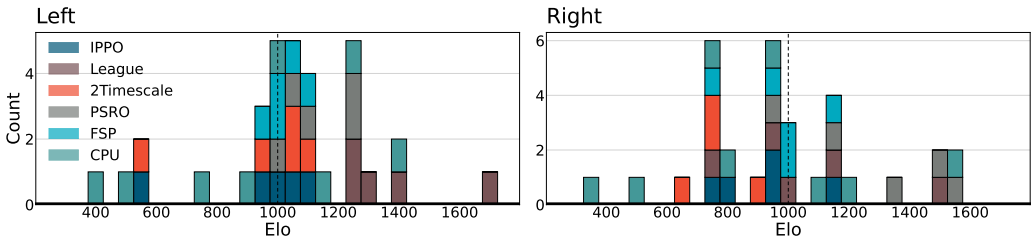


Figure 4: The distribution of Elo ratings for top ten agents from each baseline.

where  $p(\pi_L \succ \pi_C)$  is the win rate of the learner against the opponent and  $\Delta(\cdot)$  is the simplex. Intuitively, such a curriculum will encourage the agent to focus on the hardest opponents, similarly to prioritized experience play (Schaul et al., 2015). We defer other implementation details to Appendix C.

Figure 3 shows the performance of our proposed method during training. With 20 epochs of training, the agent is capable of defeating characters at each level with a win rate close to 1. In addition to beating each character with a high probability, the trained policy can complete the full video game with over 0.6 win rate, outperforming human players with hours of playing experience. This result shows that existing RL algorithms can already learn a well-behaved policy to solve the full single-player video game, which provides a good starting point for exploring the multi-agent setting.

## 7.2 PERFORMANCE OF TWO-PLAYER BASELINE ALGORITHMS

To answer question **(b)**, we evaluate five SOTA algorithms mentioned in Section 6: independent PPO (IPPO), two-timescale IPPO (2Timescale), fictitious self-play (FSP), policy-space response oracles (PSRO), and league training (League) in the scenario *sf\_ryu\_vs\_ryu*. IPPO and 2Timescale can be categorized into the independent learning paradigm, while FSP, PSRO, and League can be categorized into the population-based learning paradigm. For each algorithm, we initialize the population of agents with a pretrained policy in *sf\_ryu\_vs\_ryu(cpu)* against the most difficult CPU<sup>3</sup>. We use the transformed actions  $\mathcal{A}_{\text{trans}}$  with hard-coded special moves to unleash the full potential for agents. As a fair comparison, we use the same codebase (FightLadder-Baselines) and fix the hyperparameters of the backbone PPO algorithm. We train IPPO and 2Timescale for approximately 50M steps until the Elos saturate across all three seeds, FSP and PSRO for approximately 250M steps, and League for approximately 700M steps due to a larger population. Please refer to Appendix C for more implementation details.

For each algorithm, we report the training Elos of agents in the population during the course of training, respectively. The results are shown in Appendix D, which reveal that all baseline algorithms are improving their policies at the onset of training. Subsequently, IPPO and 2Timescale gradually converge and oscillate around the peak Elos, where FSP, PSRO, and League continue to increase their scores. This suggests that IPPO and 2Timescale may suffer from optimization issues during training and population-based methods may be more suitable for policy learning in fighting games.

To compare different baseline algorithms, we select the top ten agents (five on each left or right side) from each algorithm to form a new population, and compute the test Elos for this group of agents and CPU policies. We report the highest Elos for each algorithm in Table 2 and the distribution of these agents’ Elos in Figure 4, where we find that League and PSRO significantly outperform other baselines, and population-based methods deliver better results than independent learning counterparts, which is aligned with our previous observation inspecting Elos of baselines individually. On the other hand, we notice that CPU policies may defeat most of the agents in this group except for a few best-performing agents, suggesting that it is still very challenging for existing SOTA algorithms to reach an advanced or superhuman level of performance in these fighting games. We also noticed that two sides of agents reveal asymmetric strengths in terms of Elos in both individual evaluation for each algorithm (Appendix D Figure 8-12) and overall evaluations across algorithms (Table 2). Such an imbalance may result from various factors, for instance, optimizing instability, variance from the population or Elos computation, etc, and can be an interesting research question for future work.

<sup>3</sup>We do not pre-train in *sf\_ryu\_full\_game* as *sf\_ryu\_vs\_ryu* does not require skills to compete with other characters rather than Ryu.



Table 2: Comparison of training steps and the best Elo ratings among baselines, with CPU’s Elos as references.

Method	Training Steps (Left/Right)	Elo (Left/Right)
IPPO	46M / 46M	1082 / 1164
League	647M / 630M	<b>1682 / 1503</b>
2Timescale	51M / 46M	1080 / 919
PSRO	176M / 161M	1262 / <b>1517</b>
FSP	262M / 244M	1079 / 1150
CPU	N/A	1395 / 1541

Table 3: Comparison of methods’ exploitability. A lower number indicates the evaluated policy is more robust to exploitation.

Method	Exploitability (Left/Right)
IPPO	$0.96 \pm 0.03 / 0.91 \pm 0.03$
League	<b><math>0.94 \pm 0.05</math></b> / $0.94 \pm 0.00$
2Timescale	$0.96 \pm 0.02 / 0.90 \pm 0.05$
PSRO	$0.97 \pm 0.02 / $ <b><math>0.88 \pm 0.05</math></b>
FSP	$1.00 \pm 0.00 / 0.95 \pm 0.01$
PPO	$0.99 \pm 0.02 / 0.99 \pm 0.01$

### 7.3 NON-EXPLOITABILITY OF TRAINED AGENTS

To answer question (c), we measure the non-exploitability of baseline algorithms according to the evaluation approaches proposed in Section 5. More specifically, we choose models with the highest Elos from each two-player baseline algorithm respectively, and compare their exploitability with the single-player pretrained model used for initializing the population-based methods in Section 7.2.

**Single-agent RL Exploiters.** We use PPO as the algorithm for training exploiters, given its decent performance in both single-player and two-player scenarios shown in previous experiments. Table 3 shows the exploitability of comparing methods evaluated across three seeds, from which we observe that the single-player pretrained policy via PPO is easier to exploit and suffers from higher exploitability than almost all selected policies from two-player baselines. Therefore, this result indicates that two-player learning algorithms such as League and PSRO can help to improve the robustness of learned policies. On the other hand, the PPO exploiter eventually learns to beat policies from all baselines (with a win rate greater than 0.5), which means that none of these algorithms can result in the exact Nash equilibrium policies, or even close to it. Therefore, closing this gap is a challenging direction for future research.

**Human Players as Exploiters.** In addition to exploiting the learned models with RL algorithms, we also attempt to exploit their policies with human effort. During human evaluations, the evaluated models reveal some robustness to human players (e.g., defend when a human player attacks), but some simple strategies (e.g., defensive posture combined with low kicks at proper timing) could still defeat them rather consistently. Visualizations are provided in Appendix E.

Therefore, based on two exploiting experiments, we observe that *existing competitive MARL algorithms are found hard to learn non-exploitable strategies in competitive fighting games like Street Fighter*, thus raising a new challenge for the research community.

## 8 CONCLUSION

In this paper, we present the FightLadder platform and evaluation benchmarks as a novel testbed for competitive MARL research. The platform supports various video action games including the popular Street Fighter series, with flexible support for new game integration.

We further provide experimental evaluations of present RL and MARL algorithms in both single-player and two-player modes of one specific game Street Fighter. In the single-player setting, we proposed a learning scheme based on curriculum learning. It trains a general RL agent that can consistently beat CPUs across different characters. In the two-player setting, the Elo rating and exploitability test are conducted as part of the proposed evaluation criteria. Our implementation of league training and PSRO provides stronger agents than FSP and IPPO in terms of Elo ratings. However, both single-agent RL and human players are capable of exploiting all agents learned by current widely adopted algorithms.

This motivates further research in developing more efficient and effective self-play algorithms finding non-exploitable strategies. We hope that our platform prompts general interest and more extensive research in competitive MARL and serves as a standard benchmark for developing practically useful self-play training paradigms.

## REFERENCES

- Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor Mordatch. Emergent tool use from multi-agent autotutorials. *arXiv preprint arXiv:1909.07528*, 2019.
- Nolan Bard, Jakob N Foerster, Sarath Chandar, Neil Burch, Marc Lanctot, H Francis Song, Emilio Parisotto, Vincent Dumoulin, Subhodeep Moitra, Edward Hughes, et al. The hanabi challenge: A new frontier for ai research. *Artificial Intelligence*, 280:103216, 2020.
- Charles Beattie, Thomas Köppe, Edgar A Duéñez-Guzmán, and Joel Z Leibo. Deepmind lab2d. *arXiv preprint arXiv:2011.07027*, 2020.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- George W Brown. Iterative solution of games by fictitious play. *Act. Anal. Prod Allocation*, 13(1): 374, 1951.
- Noam Brown and Tuomas Sandholm. Superhuman ai for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424, 2018.
- Noam Brown and Tuomas Sandholm. Superhuman ai for multiplayer poker. *Science*, 365(6456): 885–890, 2019.
- Constantinos Daskalakis, Dylan J Foster, and Noah Golowich. Independent policy gradient methods for competitive reinforcement learning. *Advances in neural information processing systems*, 33: 5527–5540, 2020.
- Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- Zihan Ding, Dijia Su, Qinghua Liu, and Chi Jin. A deep reinforcement learning approach for finding non-exploitable strategies in two-player atari games. *arXiv preprint arXiv:2207.08894*, 2022.
- Ziluo Ding, Hao Luo, Ke Li, Junpeng Yue, Tiejun Huang, and Zongqing Lu. Clip4mc: An rl-friendly vision-language model for minecraft. *arXiv preprint arXiv:2303.10571*, 2023.
- A. Domahidi, E. Chu, and S. Boyd. ECOS: An SOCP solver for embedded systems. In *European Control Conference (ECC)*, pp. 3071–3076, 2013.
- Melvin Dresher, Lloyd S Shapley, and Albert William Tucker. *Advances in Game Theory.(AM-52), Volume 52*, volume 52. Princeton University Press, 2016.
- Arpad E Elo and Sam Sloan. The rating of chessplayers: Past and present. (*No Title*), 1978.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35: 18343–18362, 2022.
- Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.

- Shao-Xiang Go, Yu Jiang, and Desmond K Loke. A phase-change memristive reinforcement learning for rapidly outperforming champion street-fighter players. *Advanced Intelligent Systems*, 5(11): 2300335, 2023.
- Johannes Heinrich, Marc Lanctot, and David Silver. Fictitious self-play in extensive-form games. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 805–813, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/heinrich15.html>.
- Jian Hu, Siyang Jiang, Seth Austin Harding, Haibin Wu, and Shih-wei Liao. Rethinking the implementation tricks and monotonicity constraint in cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2102.03479*, 2021.
- Junling Hu and Michael P Wellman. Nash q-learning for general-sum stochastic games. *Journal of machine learning research*, 4(Nov):1039–1069, 2003.
- Siyi Hu, Yifan Zhong, Minquan Gao, Weixun Wang, Hao Dong, Xiaodan Liang, Zhihui Li, Xiaojun Chang, and Yaodong Yang. Marllib: A scalable and efficient multi-agent reinforcement learning library. *Journal of Machine Learning Research*, 24(315):1–23, 2023.
- Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaśkowski. Vizdoom: A doom-based ai research platform for visual reinforcement learning. In *2016 IEEE conference on computational intelligence and games (CIG)*, pp. 1–8. IEEE, 2016.
- Karol Kurach, Anton Raichuk, Piotr Stańczyk, Michał Zajac, Olivier Bachem, Lasse Espeholt, Carlos Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research football: A novel reinforcement learning environment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 4501–4510, 2020.
- Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.
- H Brendan McMahan, Geoffrey J Gordon, and Avrim Blum. Planning in the presence of cost functions controlled by an adversary. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 536–543, 2003.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Sharada Mohanty, Erik Nygren, Florian Laurent, Manuel Schneider, Christian Scheller, Nilabha Bhattacharya, Jeremy Watson, Adrian Egli, Christian Eichenberger, Christian Baumberger, et al. Flatland-rl: Multi-agent reinforcement learning on trains. *arXiv preprint arXiv:2012.05893*, 2020.
- Igor Mordatch and Pieter Abbeel. Emergence of grounded compositional language in multi-agent populations. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- David Murphy. Hacking public memory: Understanding the multiple arcade machine emulator. *Games and Culture*, 8(1):43–53, 2013.
- Alex Nichol, Vicki Pfau, Christopher Hesse, Oleg Klimov, and John Schulman. Gotta learn fast: A new benchmark for generalization in rl. *arXiv preprint arXiv:1804.03720*, 2018a.

- Alex Nichol, Vicki Pfau, Christopher Hesse, Oleg Klimov, and John Schulman. Gotta learn fast: A new benchmark for generalization in rl. *arXiv preprint arXiv:1804.03720*, 2018b.
- Alessandro Palmas. Diambra arena: a new reinforcement learning platform for research and experimentation. *arXiv preprint arXiv:2210.10595*, 2022.
- Xuehai Pan, Mickel Liu, Fangwei Zhong, Yaodong Yang, Song-Chun Zhu, and Yizhou Wang. Mate: Benchmarking multi-agent reinforcement learning in distributed target coverage control. *Advances in Neural Information Processing Systems*, 35:27862–27879, 2022.
- Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V Albrecht. Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks. *arXiv preprint arXiv:2006.07869*, 2020.
- Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Böhrer, and Shimon Whiteson. Facmac: Factored multi-agent centralised policy gradients. *Advances in Neural Information Processing Systems*, 34:12208–12221, 2021.
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *The Journal of Machine Learning Research*, 21(1):7234–7284, 2020.
- Cinjon Resnick, Wes Eldridge, David Ha, Denny Britz, Jakob Foerster, Julian Togelius, Kyunghyun Cho, and Joan Bruna. Pommerman: A multi-agent playground. *arXiv preprint arXiv:1809.07124*, 2018.
- 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. *arXiv preprint arXiv:1902.04043*, 2019.
- Bidipta Sarkar, Aditi Talati, Andy Shih, and Dorsa Sadigh. Pantheonrl: A marl library for dynamic training interactions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 13221–13223, 2022.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10): 1095–1100, 1953.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419): 1140–1144, 2018.

- Yuhang Song, Andrzej Wojcicki, Thomas Lukasiewicz, Jianyi Wang, Abi Aryan, Zhenghua Xu, Mai Xu, Zihan Ding, and Lianlong Wu. Arena: A general evaluation platform and building toolkit for multi-agent intelligence. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7253–7260, 2020.
- Joseph Suarez, Yilun Du, Clare Zhu, Igor Mordatch, and Phillip Isola. The neural mmo platform for massively multiagent research. *arXiv preprint arXiv:2110.07594*, 2021.
- Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropagation. *Advances in neural information processing systems*, 29, 2016.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the tenth international conference on machine learning*, pp. 330–337, 1993.
- 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.
- Gerald Tesauro et al. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- 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.
- Zhiyuan Yao and Zihan Ding. Learning distributed and fair policies for network load balancing as markov potential game. *Advances in Neural Information Processing Systems*, 35:28815–28828, 2022.
- Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35:24611–24624, 2022.
- Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. Plan4mc: Skill reinforcement learning and planning for open-world minecraft tasks. *arXiv preprint arXiv:2303.16563*, 2023.
- Daochen Zha, Kwei-Herng Lai, Songyi Huang, Yuanpu Cao, Keerthana Reddy, Juan Vargas, Alex Nguyen, Ruzhe Wei, Junyu Guo, and Xia Hu. Rlcard: A platform for reinforcement learning in card games. In *IJCAI*, 2020.
- Huichu Zhang, Siyuan Feng, Chang Liu, Yaoyao Ding, Yichen Zhu, Zihan Zhou, Weinan Zhang, Yong Yu, Haiming Jin, and Zhenhui Li. Cityflow: A multi-agent reinforcement learning environment for large scale city traffic scenario. In *The world wide web conference*, pp. 3620–3624, 2019.
- 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. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Ming Zhou, Jun Luo, Julian Villeda, Yaodong Yang, David Rusu, Jiayu Miao, Weinan Zhang, Montgomery Alban, Iman Fadarar, Zheng Chen, et al. Smarts: Scalable multi-agent reinforcement learning training school for autonomous driving. *arXiv preprint arXiv:2010.09776*, 2020.
- Ming Zhou, Ziyu Wan, Hanjing Wang, Muning Wen, Runzhe Wu, Ying Wen, Yaodong Yang, Yong Yu, Jun Wang, and Weinan Zhang. Malib: A parallel framework for population-based multi-agent reinforcement learning. *Journal of Machine Learning Research*, 24(150):1–12, 2023.

## A DETAILS OF FIGHTLADDER

### A.1 DENSE REWARD

The shaped dense reward for the  $i$ -th agent at step  $t$  is defined as follows:

$$r_{i,t} = \alpha[\lambda(\text{HP}_{-i,t-1} - \text{HP}_{-i,t}) - (\text{HP}_{i,t-1} - \text{HP}_{i,t}) + r_{i,\text{bonus}}], \quad (1)$$

where  $\alpha$  is a scaling factor,  $\text{HP}_{i,t}$  denotes agent  $i$ 's hit-point at step  $t$  and  $\lambda$  control the aggressiveness of learned agents, and  $-i$  denotes the opponent agent. At the end of the game, the agent  $i$  will receive a bonus reward  $r_{i,\text{bonus}}$ , which is positively correlated to  $\text{HP}_i$  if it wins and negatively correlated to  $\text{HP}_{-i}$  if it loses. By default, we choose  $\lambda = 3$  in SF2, FF2, and MK, and  $\lambda = 1$  in SF3 and KOF97, for the consideration of practical performances.

### A.2 GAME SETTINGS

Table 4 illustrates the observation, action, and rewards as well as other elements in the environment for all supported games — Street Fighter II (SF2), Fatal Fury 2 (FF2), Mortal Kombat (MK), Street Fighter III (SF3), and The King of Fighters '97 (KOF97).

Table 4: Specification of supported games in FightLadder.

	SF2	FF2	MK	SF3	KOF97
Observation (Pixels)	100×128×3	112×128×3	112×160×3	112×192×3	112×192×3
Human Action Supported	Yes	Yes	Yes	Yes	Yes
Transformed Action Supported	Yes	Yes	Yes	No	No
Shaped Dense Reward	Yes	Yes	Yes	Yes	Yes
Default Frames Per Step	8	8	8	3	3
Default Frames Stacked <sup>4</sup>	12	12	12	9	9
Additional Available Info	HPs, Countdown, Scoreboard, Positions	HPs, Countdown	HPs, Countdown, Scoreboard	HPs	HPs, Countdown, Positions, Power Status

## B BASELINE ALGORITHMS OF FIGHTLADDER-BASELINES

**Independent Learning (IPPO).** Independent learning is a straightforward extension of single-agent RL into MARL. It decomposes the joint optimization into individual ones for each agent while regarding all other agents as part of the environment. It can be implemented easily by simultaneously running single-agent RL algorithms for each player. Theoretically, this independent learning paradigm suffers from suboptimality (Tan, 1993; Foerster et al., 2018), because the environment becomes non-stationary while other agents are updating their policies. However, recent work (de Witt et al., 2020; Yu et al., 2022) finds that with modest hyperparameter tuning, IPPO can serve as a strong baseline compared to other state-of-the-art algorithms in some cooperative MARL tasks.

**Two-timescale Learning (2Timescale).** Two-timescale learning follows the independent learning paradigm, but requires two players to update gradients according to the two-timescale rule, i.e., one player uses a much smaller step size than the other one. As a result of this modification, two-timescale learning enjoys some nice theoretical properties — it is proven that under some mild assumptions, independent policy gradient algorithms satisfying two-timescale converge to a Nash equilibrium in two-player zero-sum stochastic games (Daskalakis et al., 2020).

**Population-Based Methods.** The independent learning framework is only training agents against the current version of their opponents, which may fail or converge slowly due to the lack of diversity (Dresher et al., 2016). Population-based methods are proposed to increase policy diversity by maintaining a pool of policies in previous iterations, and using them as a curriculum to update the current policy. More specifically, for  $t$ -th update, the agent  $\mu^t$  plays with previous versions of its opponent  $\tilde{\nu}$  sampled from the meta-strategy  $\rho_\nu$ , which is a distribution over  $\nu^0, \nu^1, \dots, \nu^{t-1}$ . Algorithm 1 presents the pseudo-code for general population-based methods. With different choices of sampling distribution, we can recover several state-of-the-art baselines:

<sup>4</sup>We uniformly sample the stacked frames as observations to improve the computational efficiency.

- **Fictitious Self-Play (FSP)**, where  $\rho_\nu$  is the uniform distribution  $U(\nu^0, \nu^1, \dots, \nu^{t-1})$  (Heinrich et al., 2015).
- **Policy-Space Response Oracles (PSRO)**, where  $\tilde{o}$  is sampled from a meta-strategy by solving Nash equilibrium of the payoff matrix game between  $\mu^0, \mu^1, \dots, \mu^{t-1}$  and  $\nu^0, \nu^1, \dots, \nu^{t-1}$  (Lanctot et al., 2017).
- **League Training (League)**, where three types of agents — main agents, league exploiters, and main exploiters, are introduced into the population. Main agents train against themselves as well as all previous versions of agents in the population; league exploiters train against all previous agents; and main exploiters optimize the best response of main agents. Each type of agent adopts a different sampling distribution which is a mixture of self-play and prioritized fictitious self-play. We refer readers to (Vinyals et al., 2019) for more implementation details.

---

**Algorithm 1** Population-Based Methods for MGs
 

---

```

1: Initialize policies  $\mu^0 = \{\mu_h\}, \nu^0 = \{\nu_h\}, h \in [H]$ 
2: Initialize policy sets:  $\mu = \{\mu^0\}, \nu = \{\nu^0\}$ 
3: Initialize meta-strategies:  $\rho_\mu = [1.], \rho_\nu = [1.]$ 
4: for  $t = 1, \dots, T$  do
5:   if  $t \% 2 == 0$  then
6:      $\nu^t = \text{BEST\_RESPONSE}(\rho_\mu, \mu)$ 
7:      $\nu = \nu \cup \{\nu^t\}$ 
8:     Update  $\rho_\nu$  according to specific algorithms
9:   else
10:     $\mu^t = \text{BEST\_RESPONSE}(\rho_\nu, \nu)$ 
11:     $\mu = \mu \cup \{\mu^t\}$ 
12:    Update  $\rho_\mu$  according to specific algorithms
13:   end if
14: end for
15: Return  $\mu, \rho_\mu, \nu, \rho_\nu$ 

```

---

## C EXPERIMENT DETAILS

### C.1 HYPERPARAMETERS (TABLE 5 AND 6)

Hyperparameters	Value
feature extractor	CNN (Mnih et al., 2015)
rollout steps for each environment	512
batch size	1024
epochs per update	4
$\gamma$	0.94
GAE $\lambda$	0.95
learning rate	linear schedule from 2.5e-4 to 2.5e-6
clipping range	linear schedule from 0.15 to 0.025
advantage normalization	True
entropy coefficient	0.0
gradient clipping	0.5
value function coefficient	0.5

Table 5: Training hyperparameters for PPO, which is the backbone for both single-player and two-player algorithms in the experiment.

### C.2 TRAINING DETAILS

Figure 5, 6, and 7 report the payoff matrix of policies within the population for FSP, PSRO, and League, respectively, with the value representing the win rate of the left player against the right player.

<b>FSP</b>		<b>PSRO</b>		<b>League</b>	
# envs per learner	24	# envs per learner	24	# envs per learner	24
steps for BR	10M	steps for BR	10M	steps for BR	10M
total steps	50M	total steps	250M	total steps	700M
# main agent	1	# main agent	1	# main agent	1
		Nash solver	ECOS (Domahidi et al., 2013)	# main exploiter	1
				# league exploiter	2

Table 6: Training hyperparameters for FSP, PSRO, and League. We omit the details of League’s opponent scheduling here as it strictly follows the pseudocode provided in (Vinyals et al., 2019).

## D INDIVIDUAL ELO RESULTS

D.1 IPPO (FIGURE 8)

D.2 2TIMESCALE (FIGURE 9)

D.3 FSP (FIGURE 10)

D.4 PSRO (FIGURE 11)

D.5 LEAGUE (FIGURE 12)

## E VISUALIZATION OF HUMAN EXPLOITERS

Figure 13 visualizes how human players can exploit learned models with a simple strategy.



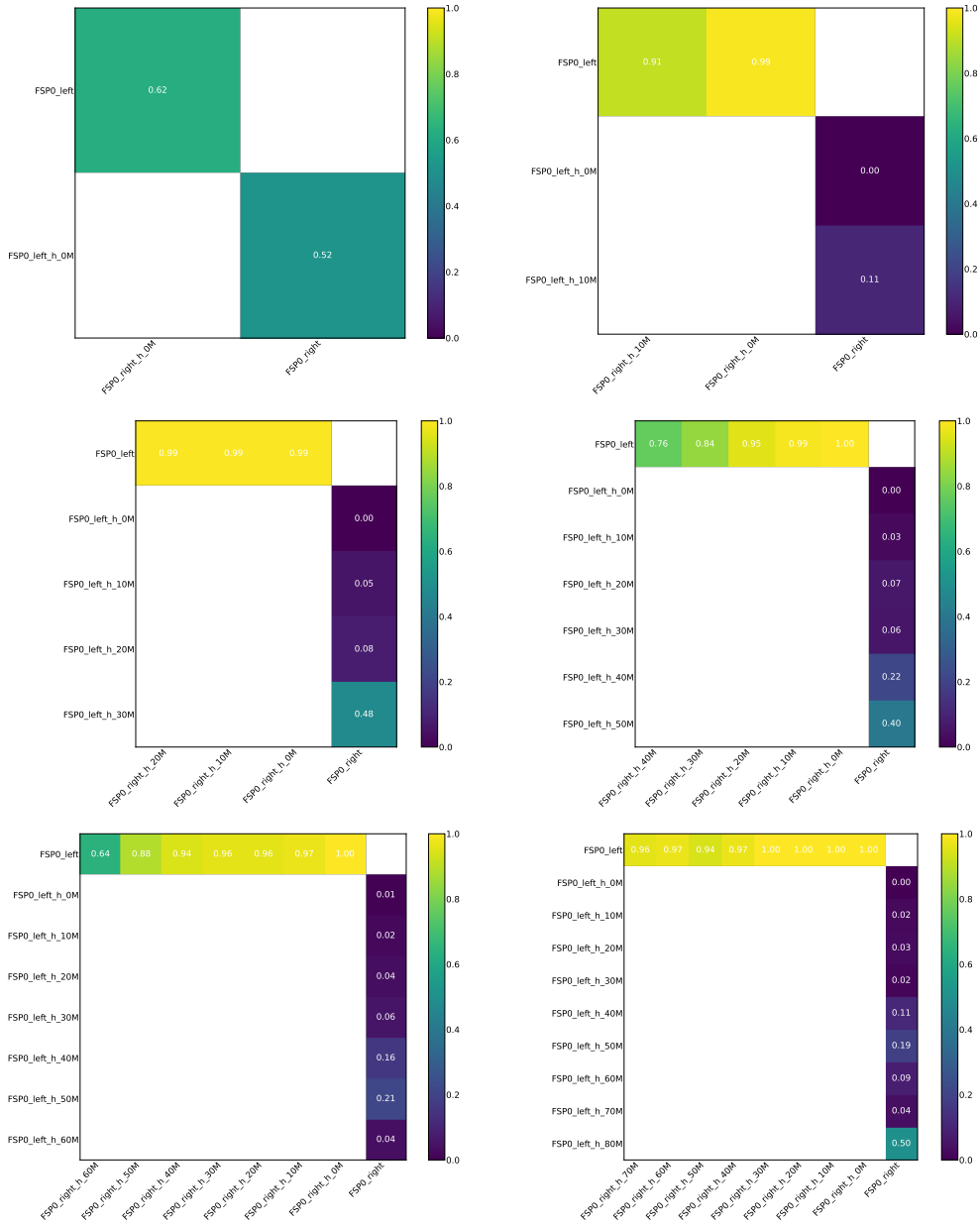


Figure 5: FSP details (training order from top left to bottom right): For FSP, there is one agent for each side (left or right). The name of each row indicates the agent information as Character\_Side\_Checkpoint. Checkpoint=h\_xM represents a previous version of agent saved at x million steps. The value indicates the win rate of the left (row) player against the right (column) player.

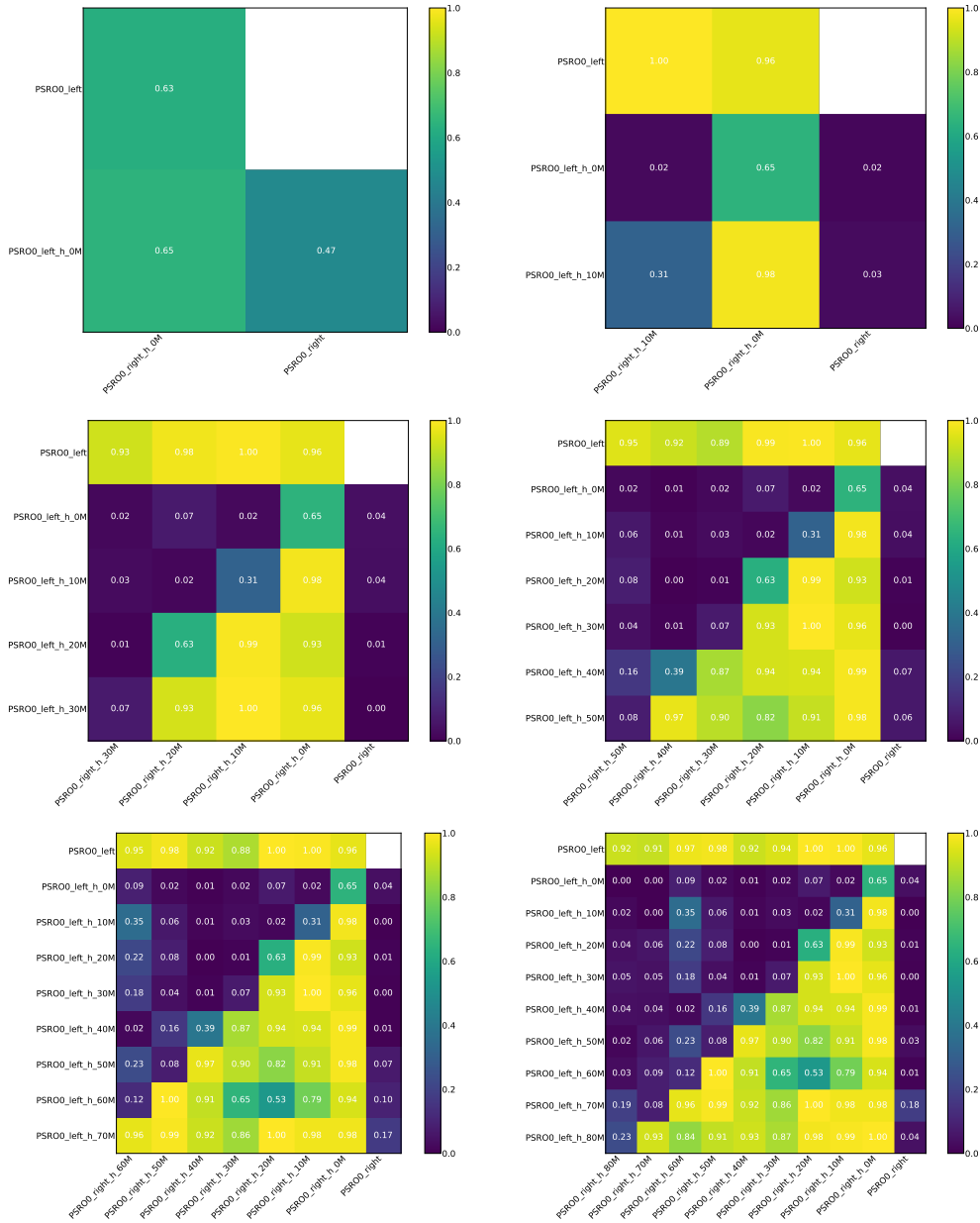


Figure 6: PSRO details (training order from top left to bottom right): For PSRO, there is one agent for each side (left or right). The name of each row indicates the agent information as Character\_Side\_Checkpoint. Checkpoint=h\_xM represents a previous version of agent saved at x million steps. The value indicates the win rate of the left (row) player against the right (column) player.

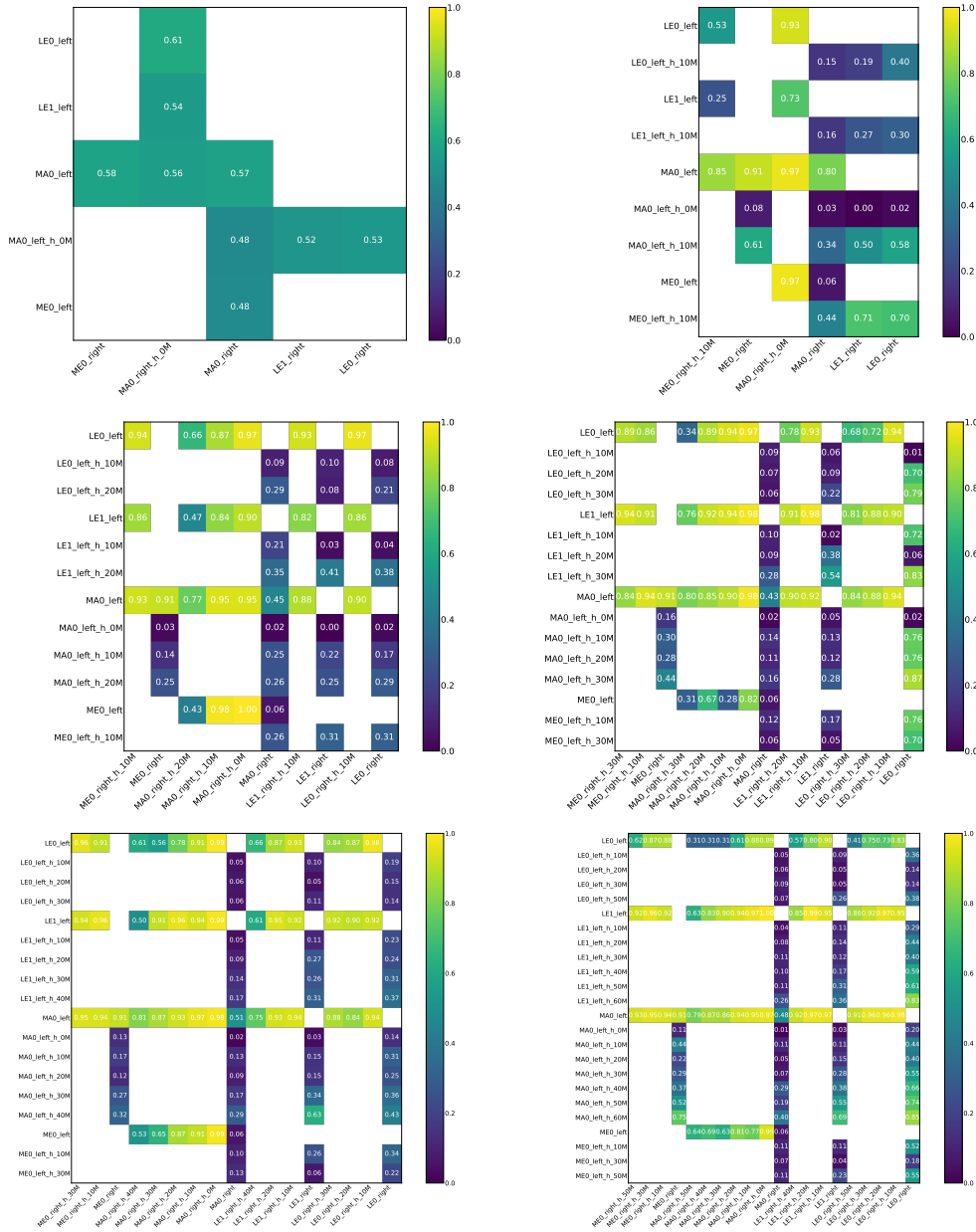


Figure 7: League training details (training order from top left to bottom right): For league training, there is one main agent (MA), two league exploiters (LE0, LE1), and one main exploiter (ME) for each side (left or right). The name of each row indicates the agent information as Character\_Side\_Checkpoint. Checkpoint=h\_xM represents a previous version of agent saved at x million steps. The value indicates the win rate of the left (row) player against the right (column) player.

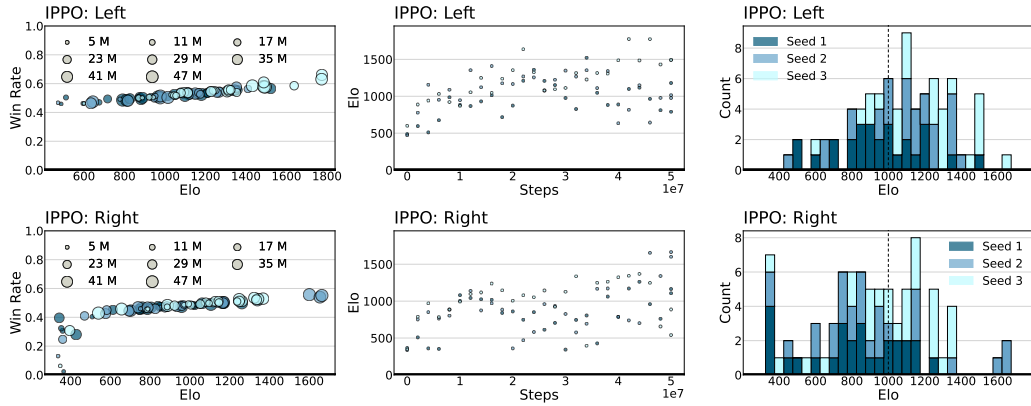


Figure 8: The Elo rating for the population of agents trained with IPPO algorithm. The upper three plots are for left-side player and the bottom three are for the right-side player. The Elo rating is plotted against the winning rate over matched policies (left figures), training steps (middle figures) and the number of policies (right figures).

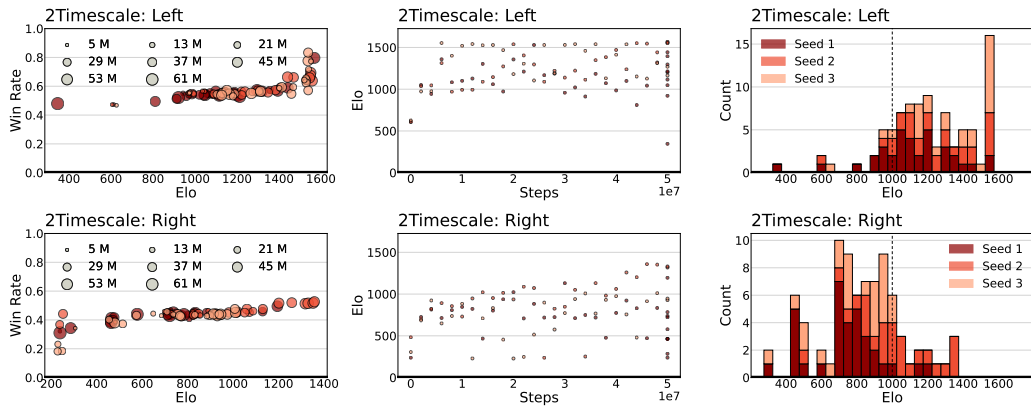


Figure 9: The Elo rating for the population of agents trained with 2Timescale algorithm. The upper three plots are for left-side player and the bottom three are for the right-side player. The Elo rating is plotted against the winning rate over matched policies (left figures), training steps (middle figures) and the number of policies (right figures).

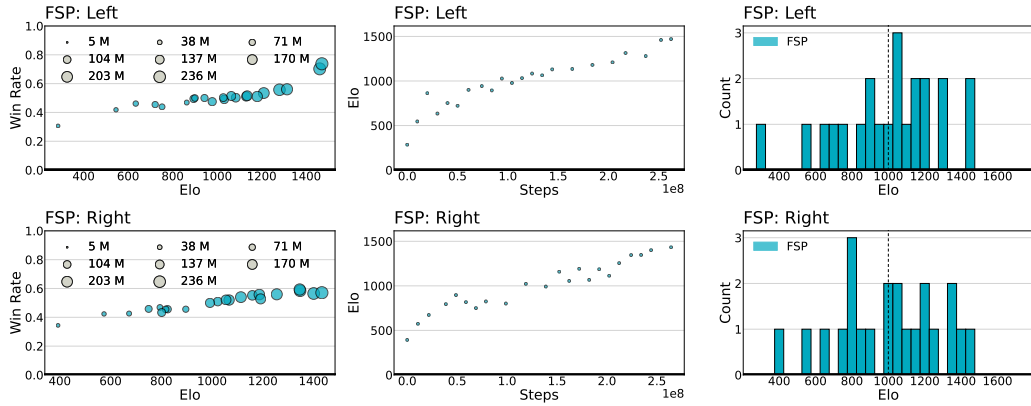


Figure 10: The Elo rating for the population of agents trained with FSP algorithm. The upper three plots are for left-side player and the bottom three are for the right-side player. The Elo rating is plotted against the winning rate over matched policies (left figures), training steps (middle figures) and the number of policies (right figures).

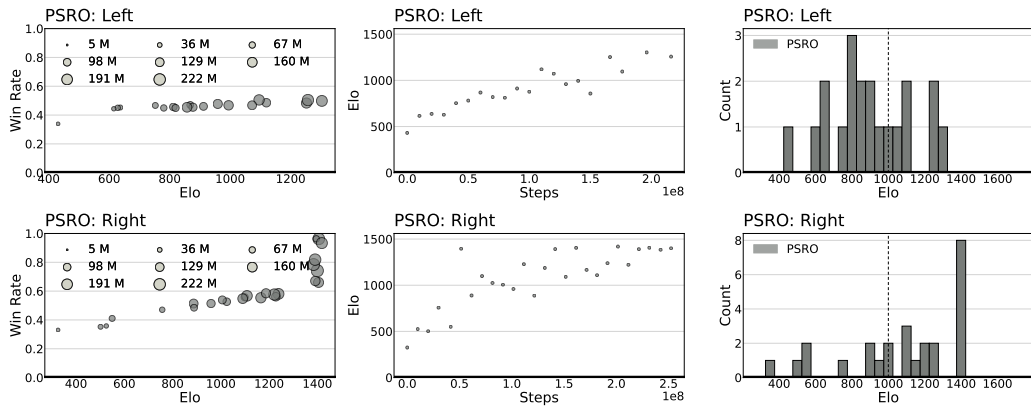


Figure 11: The Elo rating for the population of agents trained with PSRO algorithm. The upper three plots are for left-side player and the bottom three are for the right-side player. The Elo rating is plotted against the winning rate over matched policies (left figures), training steps (middle figures) and the number of policies (right figures).

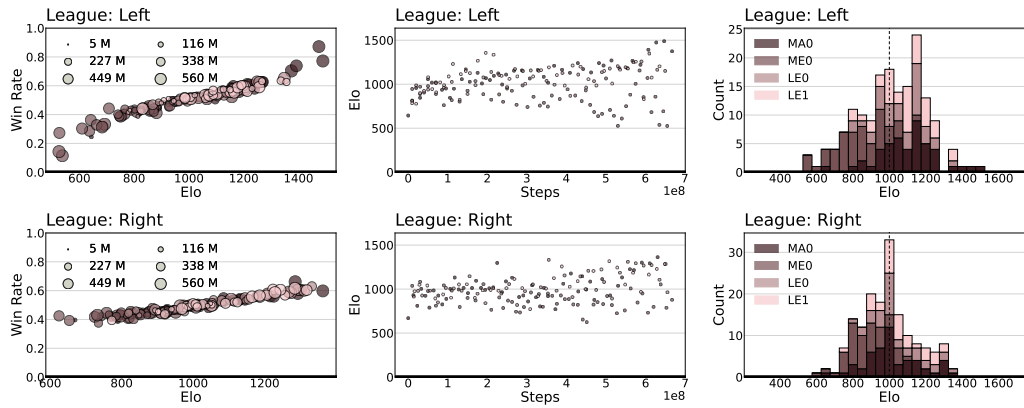


Figure 12: The Elo rating for the population of agents trained with League training. The upper three plots are for left-side player and the bottom three are for the right-side player. The Elo rating is plotted against the winning rate over matched policies (left figures), training steps (middle figures) and the number of policies (right figures).



Figure 13: Demonstration of the exploiting strategy of one human player. The human player (Ryu on the right in white) defends when the AI opponent (Ryu on the left in gray) attacks, and inflicts damage with low kicks.