

# Creating Pro-Level AI for a Real-Time Fighting Game Using Deep Reinforcement Learning

Inseok Oh<sup>ID</sup>, Seungeun Rho, Sangbin Moon, Seongho Son, Hyoil Lee<sup>ID</sup>, and Jinyun Chung<sup>ID</sup>

**Abstract**—Reinforcement learning (RL) combined with deep neural networks has performed remarkably well in many genres of games recently. It has surpassed human-level performance in fixed game environments and turn-based two-player board games. However, to the best of our knowledge, current research has yet to produce a result that has surpassed human-level performance in modern complex fighting games. This is due to the inherent difficulties with real-time fighting games, including: vast action spaces, action dependencies, and imperfect information. We overcame these challenges and made 1v1 battle AI agents for the commercial game *Blade and Soul*. The trained agents competed against five professional gamers and achieved a winning rate of 62%. This article presents a practical RL method that includes a novel self-play curriculum and data skipping techniques. Through the curriculum, three different styles of agents were created by reward shaping and were trained against each other. Additionally, this article suggests data-skipping techniques that could increase data efficiency and facilitate explorations in vast spaces. Since our method can be generally applied to all two-player competitive games with vast action spaces, we anticipate its application to game development including level design and automated balancing.

**Index Terms**—Deep learning, fighting game, imperfect information, reinforcement learning (RL), self-play curriculum learning.

## I. INTRODUCTION

**R**EINFORCEMENT learning (RL) is extending its boundaries to a variety of game genres. In player versus environment settings, such as those found in Atari 2600 games, RL agents have exceeded human level performance using various methods [5], [15], [16], [19]. Likewise, in player versus player (PVP) settings, neural networks combined with search-based methods beat the best human players in turn-based games with two or more players—such as *Go*, *Chess* [20], and *Mahjong* [28]. Recently, RL research in games has shifted focus to the PVP settings found in more complex video games such as *StarCraft2* [24], *Quake3* [10], and *Dota2* [18]. Even grand-master level RL agents have been developed for *StarCraft2* [29], which is a highly complex imperfect information game where an agent has to control multiple units at a time.

Manuscript received May 4, 2020; revised October 14, 2020 and December 17, 2020; accepted December 25, 2020. Date of publication January 6, 2021; date of current version June 16, 2022. (Inseok Oh and Seungeun Rho equally contributed to this work.) (Corresponding author: Jinyun Chung.)

The authors are with the Game AI Lab, AI Center, NCSoft, Gyeonggi-do 13494, South Korea (e-mail: ohinsuk@ncsoft.com; gloomymonday@ncsoft.com; sangbin@ncsoft.com; hingdoong@ncsoft.com; onetop21@ncsoft.com; jchung2050@ncsoft.com).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TG.2021.3049539>.

Digital Object Identifier 10.1109/TG.2021.3049539

TABLE I  
FIGHTING GAMES FROM OTHER WORKS

	Year	Commercial	Dimension	Pro-scene
FICE	2013	No	2D	No
LF2	1999	Yes	2.5D	No
SSBM	2001	Yes	2D	Yes
BAB	2013	Yes	3D	Yes

Fighting games—as one of the most representative types of complex PVP games—have been the focus of multiple studies that have made progress in this area. For instance, Monte Carlo tree search based methods [9], [11], [27] have been applied to *FightingICE* (*FICE*), a game platform made for the fighting game AI competition [13]. However, it is hard to fulfill real-time conditions when applied to heavier game engines with longer query times. Additionally, a deep RL based agent [12] was trained against a rule-based fixed opponent in “little fighter 2 (LF2).” However, since the opponent’s decision is unknown at a player’s decision time, agents trained against rule-based AIs cannot be generalized for unseen opponents. Our approach is largely similar to that of [3] in which a self-play deep RL method was applied to *Supersmash Bros. Melee* (*SSBM*). However, the complexity of state and action space is significantly limited compared to our 3-D environment with complex game rules. We created pro-level AI agents for the real-time fighting game *Blade & Soul* (*B&S*) *Arena Battle* via novel self-play based RL.

*B&S* is a commercial massively multiplayer online role-playing game. It supports duels between two players called *B&S Arena Battles* (*BABs*). As given in Table I, *BAB* is a more modern fighting game compared to the games considered in other works; hence, it has much more complex game dynamics and heavier game engines. Additionally, a large number of people play *BAB* and it has more active professional scenes<sup>1</sup> than other fighting games. *BAB*’s larger number of active professional scenes stands out more significantly when compared to *FightingICE*, which was designed solely for research purposes.

Fig. 1 displays a scene from *BAB*. In *BAB*, two players fight against each other to reduce their opponent’s health point (HP)

<sup>1</sup>Nine regional league winners from all over the world (including KOR, NA, EU, RUS, and CHN) participated in the 2018 *B&S* world championship (fourth annual event). The winning prize was approx. \$50k (compared to *Tekken7*: \$30k)



Fig. 1. Scene from the *B&S* arena battle.

to zero within three minutes. To master *BAB*, an agent must be able to deal with multiple challenges.

First, an agent must manage vast action and state spaces compared to other fighting games [3], [11], [12]. An agent must make skill, move, and targeting decisions simultaneously, which yields many possible combinations. As a rough estimate, there are 144 potential actions for each time step: 8 (avg. # of avail. skills) \* 9 (8 directional + no move) \* 2 (facing opponent or moving direction). Since the average game length is 1200 time steps (120 s), numerous scenarios are possible—not considering the opponent’s actions.

Moreover, an agent must consider the dependencies between skills: e.g., a skill may become available only for a short period of time following the use of another skill. As a result, out of the 45 skills in total (including “no-op”), the set of skills available at a given time constantly changes. The agent must also consider the properties of each skill because they have different cooldown times (required interval for reusing a skill) and skill point (SP) consumptions, and serve one or more of five different functions: damage dealing, crowd control (a set of skills that reduces the number of possible actions the opponent can utilize; abbreviated CC) [32], resistance (which functions to make the player immune or resistant to CC skills), escape, and dash. In *BAB*, crowd control refers to a set of skills that reduces the number of possible actions the opponent can utilize. For example, when you “stun” a Destroyer (one of the classes in *BAB*), it can only use the skill “escape.” When a destroyer is “groggy,” it becomes limited to the skills “escape” and “retreat.” When it is “down,” the *Destroyer* is limited to 5-6 possible skills.

Finally, an agent must deal with imperfect information settings. Because *BAB* is a real-time game, two players make their decisions simultaneously. This indicates that an agent is required to make decisions without knowing the opponent’s decision or strategy. Hence, *BAB* can be considered an imperfect information game [30], [33]. For example, when a player uses a resistance skill and the opponent uses a crowd control skill at the same time, the player gains advantage over the opponent. As a result, the essence of the problem is to approximate a Nash equilibrium strategy so that the agent can respond appropriately to any opposing strategy.

To tackle these challenges, we have made improvements to vanilla self-play algorithm by diversifying opponent pools and skipping data to facilitate exploration. The main contributions of this article are as follows:

- 1) We devised a novel self-play curriculum [35] with agents of different styles. The curriculum made these agents compete against each other and reinforced the agents simultaneously, rendering the agents capable of handling a variety of opponents. We empirically demonstrate that our curriculum outperforms vanilla self-play method.
- 2) We diversified the fighting style of the game-playing AIs by reward shaping [17]. We created three types of agents with different fighting styles: aggressive; defensive; and balanced. We anticipate its application to game development including level design and automated balancing.
- 3) We introduced data skipping techniques to enhance exploration in vast space. These can be generally applied to any two-player real-time fighting games.
- 4) We evaluate our agents by pitting them against professional players in the 2018 B&S World Championship blind match. Our AI agents won three out seven matches, while the aggressive one beating all professional players both in the live event and pretest.

## II. BACKGROUND

### A. Reinforcement Learning

In RL [23], agent and environment can be formalized as a Markov decision process (MDP) [8]. For every discrete time step  $t$ , an agent receives a state  $s_t \in S$  and sends an action  $a_t \in A$  to the environment. Then, the environment makes a state transition from  $s_t$  to  $s_{t+1}$  with the state transition probability  $P_{ss'}^a = P[s', a]$  and gives a reward function  $R : S \times A \rightarrow \mathbb{R}$ , with the reward signal  $r_t = R(s_t, a_t) \in \mathbb{R}$  given to the agent. Therefore, this process can be expressed with  $\{S, A, P, R, \gamma\}$ , where  $\gamma \in [0, 1]$  is a discount factor, which represents the preference for immediate reward over long-term reward. Here, the agent samples an action from a policy  $\pi : S \rightarrow P(A)$ , where  $P(A)$  represents the set of probability distributions. The learning process modifies the policy to encourage good actions and suppress bad actions. The objective of the learning is to find the optimal policy  $\pi^*$  that maximizes the expected discounted cumulative reward

$$\pi^* = \operatorname{argmax}_{\pi} E_{\pi} [\sum_t \gamma^t r_t].$$

### B. Real-Time Two Player Game

In a real-time two player game, there are two players, namely, the agent and the opponent. Both of them send an action to the environment at the same time. Let us denote the policy of the agent as  $\pi^{\text{ag}}$ , and the policy of the opponent as  $\pi^{\text{op}}$ . Each samples an action from its own policy for every time step

$$a_t^{\text{ag}} \sim \pi^{\text{ag}}(a_t^{\text{ag}} | s_t), \quad a_t^{\text{op}} \sim \pi^{\text{op}}(a_t^{\text{op}} | s_t).$$

Then, the environment makes a state transition by considering those two actions jointly

$$s_{t+1} \sim P(s_{t+1} | s_t, a_t^{\text{ag}}, a_t^{\text{op}}), \quad r_{t+1} = R(s_t, a_t^{\text{ag}}, a_t^{\text{op}}).$$

Here, the MDP can be expressed as  $\{S, A^{\text{ag}}, A^{\text{op}}, P, R, \gamma\}$ . If  $\pi^{\text{op}}$  is fixed, then we can regard the opponent as a part of the

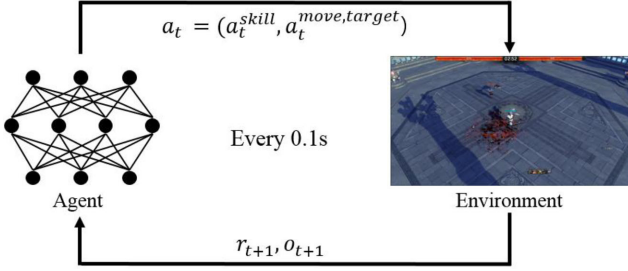


Fig. 2. Agent-environment plot in BAB.

environment by marginalizing the policy of the opponent. This way, we can obtain  $P'$  and  $R'$

$$\begin{aligned}
 P' (s_{t+1} | s_t, a_t^{ag}) &= \sum_{a_t^{op}} \pi^{op} (a_t^{op} | s_t) * P (s_{t+1} | s_t, a_t^{ag}, a_t^{op}) \\
 R' (s_t, a_t^{ag}) &= \sum_{a_t^{op}} \pi^{op} (a_t^{op} | s_t) * R (s_t, a_t^{ag}, a_t^{op}).
 \end{aligned}$$

Then, the MDP expression turns into a simpler form with  $P'$  and  $R'$ :  $\{S, A^{ag}, P', R', \gamma\}$ . This expression is coherent with the one player MDP. Therefore, any methods for the original MDP work in this form as well.

### C. B&S Arena Battles as Markov Decision Process

If we assume  $\pi^{op}$  or the pool of  $\pi^{op}$  is fixed, *BAB* can be expressed as an MDP [34]. Fig. 2 illustrates the agent-environment framework in *BAB*. Long short-term memory (LSTM) [6] based agents interact with the *BAB* simulator, which acts as the environment. For every time step with 0.1 sec intervals, state  $s_t$  is constructed from the history of observations  $H_t = \{o_1, o_2, \dots, o_t\}$ . To be specific,  $s_t$  is composed of any information that a human can access during a game, such as HP, SP, distance from opponent, distance from the arena wall, current position, remaining game time, remaining cooldown times for all 44 skills, an agent's status info (midair, stun, down, kneel, etc.), and so on. Then, the agent decides on an action  $a_t = (a_t^{skill}, a_t^{move, target})$  for every time step. Note that the targeting action (i.e., orientation) space was originally continuous. We discretized the space into two actions—facing the opponent and facing away from the opponent—and jointly considered it along with the quantized move decision. Following this, the action is then sent to the environment and a state transition occurs accordingly.

Here, exact rewards should also be determined. Rewards are closely related to high performance in *BAB*. We provided  $r_t^{WIN}$ , which is the reward for winning a game, and  $r_t^{HP}$ , the reward for the changes in HP margin. These rewards are designed based on the assumption that the more a player wins, and with more remaining HP, the better that player's performance is.  $r_t^{WIN}$  is given at the terminal step of each episode with +10 for a win and -10 for a loss.  $r_t^{HP}$  may occur at every time step when the agent deals damage to the opponent and vice versa. Since HP is

normalized to  $[0, 10]$ ,  $r_t^{WIN}$  and  $r_t^{HP}$  have the same scale.

$$\begin{aligned}
 r_t &= r_t^{WIN} + r_t^{HP} + r_t^{EXTRA} \\
 r_t^{HP} &= (HP_t^{ag} - HP_{t-1}^{ag}) - (HP_t^{op} - HP_{t-1}^{op}) \\
 r_t^{EXTRA} &= -(r_t^{time} + r_t^{distance})
 \end{aligned}$$

where  $r_t^{EXTRA}$  is an additional reward for guiding battle styles. It is the sum of the time penalty and the distance penalty.  $r_t^{time}$  is a reward based on the game length, and  $r_t^{distance}$  is a reward based on the distance between agents. These additional rewards are described in further detail in the next section. The value of  $\gamma$  is set to 0.995, which is close to 1.0, since all episodes in *BAB* are forced to terminate after 1800 time steps (=3 min).

### III. SELF-PLAY CURRICULUM WITH DIVERSE STYLES

$\pi^{op}$  needs to be fixed to formalize *BAB* as an MDP. However,  $\pi^{op}$  is not fixed in general and our agent does not know which  $\pi^{op}$  it is going to face. We propose a self-play curriculum with a diversified pool of  $\pi^{op}$  to solve this issue. Existing self-play methods ([21], [22]) generally use opponent pools for training. Parameters of a network are stored at regular intervals during training to create a pool of past selves. Opponents are then sampled from this pool.

Although the self-play method of RL offers a way to learn the Nash equilibrium strategy [4], high coverage of strategy space is essential to efficiently find one. Vanilla self-play alone does not guarantee enough coverage for games with large problem spaces. To tackle this problem, AlphaStar [25] diversified the opponent pool by imitating different human strategies and introducing three types of agents with different match making scheme. The Poker AI, Pluribus [1], hand-tuned three different strategies on top of basic blueprint strategy. The three strategies are biased toward raising, calling, and folding, respectively.

Concurrently, we devised a novel self-play curriculum. We enforced diversity of agents' strategies by introducing a range of different battle styles, and agents of different styles were made to compete against each other.

#### A. Guiding Battle Styles Through Reward Shaping

One of the most noticeable fighting styles to invest with is the degree of aggressiveness. We used three dimensions of rewards to control the degree of aggressiveness. The first dimension is the "time penalty." The aggressive agent receives larger penalties per time step, and this motivates it to finish the match in a shorter period of time. The second dimension is the relative importance of the agent's HP to the opponent's HP. Aggressive players will try to reduce the opponent's HP rather than preserving their own HP, while defensive players tend to act the opposite way. The final dimension is the "distance penalty." Defensive players tend to ensure a certain distance from their opponents to respond appropriately against attacks, while aggressive players tend to approach their opponents and attack relentlessly. To realize these properties, the aggressive agent received larger penalties in proportion to the distance between itself and its opponent. The specific reward weights used for each style are given in Table II.



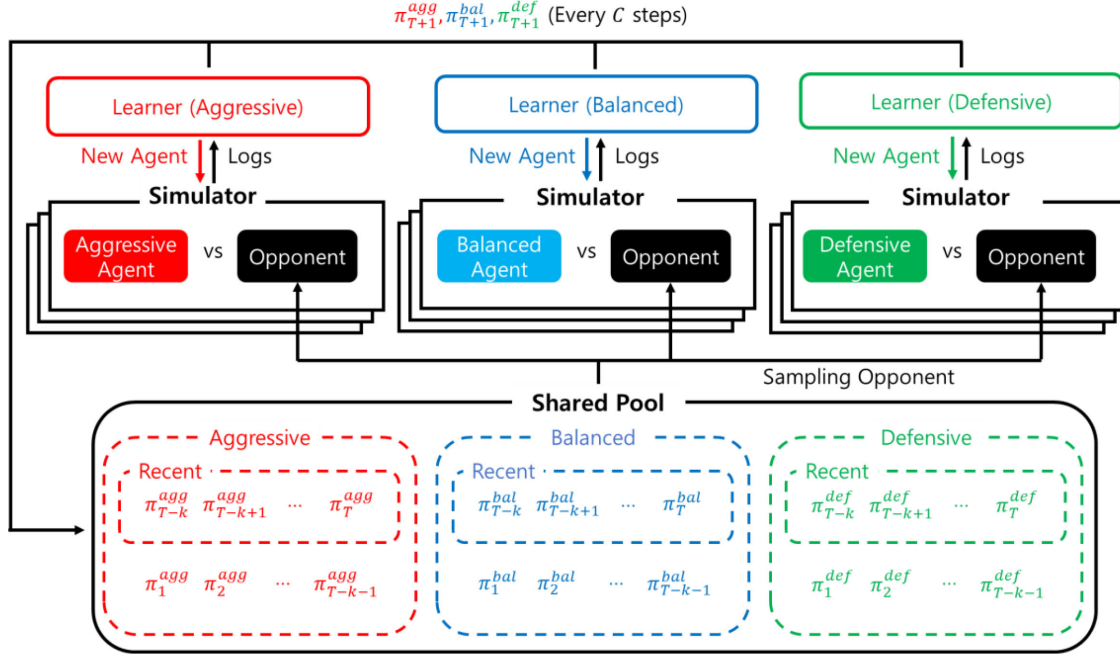


Fig. 3. Overview of self-play curriculum with three different styles.

TABLE II  
REWARD DETAILS OF EACH STYLE

	Aggressive	Balanced	Defensive
Time penalty	0.008	0.004	0.0
HP ratio	5:5	5:5	6:4
Distance penalty	0.002	0.0002	0.0

Note that each of these three dimensions can take continuous values. This means that it is possible to create a spectrum of different fighting styles with varying degrees of aggressiveness. However, to effectively demonstrate the viability of this method, we limited the number of fighting styles to three. By using any type of additional reward signals along with  $r_t^{WIN}$  and  $r_t^{HP}$ , this method could be applied to other fighting games in general to create agents with various fighting styles.

### B. Our Self-Play Curriculum

Fig. 3 shows an overview of the proposed self-play curriculum with three different types of agents. Agents of each style have their own learning process, and all three agent types were trained in a concurrent manner.

Each learning process consists of a learner and multiple simulators. The learner and the simulators work asynchronously. In the simulators, an agent constantly plays matches against randomly sampled opponents from the shared pool. The most recent  $k$  models (see Table VI in the Appendix) of each style are uniformly selected with total probability mass of  $p$ , while other models are chosen uniformly with probability  $1-p$ . As training goes on,  $p$  is linearly annealed from 0.8 to 0.1. A higher  $p$  assists in swift adaptation to the latest opponents,

while a lower  $p$  stabilizes the learning process by alleviating catastrophic forgetting. Each simulator sends a match log to the learner at the end of every match and updates its agent with the latest parameters received from the learner. The same procedure continues to be used through subsequent games.

The learner trains its agents in an off-policy manner using logs gathered from multiple simulators and sends the latest network parameters to the simulators on request. In addition, the learner sends its network parameters to the shared pool every  $C$  steps (e.g.,  $C = 10000$ ) of update. Thus, the pool has varying policies that come from the different learning processes of the different styles. These sets of model parameters are provided as opponents to each learning process. By sharing a pool, every learning agent encounters opponents of every style during training and learns how to deal with them. Therefore, agents trained via our self-play curriculum can ultimately learn how to face opponents with varying fighting styles while maintaining their own battle styles.

## IV. DATA SKIPPING TECHNIQUES

In this section, we detail two data skipping techniques: “no-op” and “maintain move decision.” Data skipping techniques refer to the process of dropping certain data during training and evaluation procedures.

### A. Discarding Passive “No-Op”

In fighting games, using skills generally consumes resources, such as SP and cooldown time. Therefore, if a player overuses a certain skill, it will not be available for use during actual times of need. Thus, players should strategically use and retain their skills to ensure their availability when needed. To take this aspect into account, we concatenated a “no-op” action to the output of

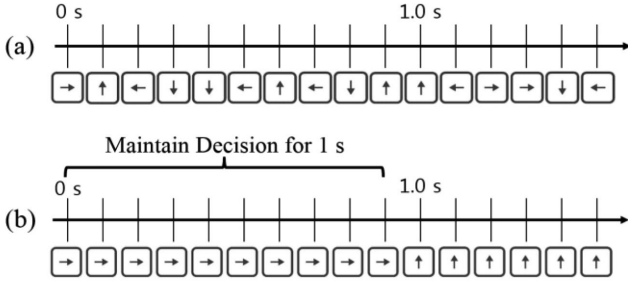


Fig. 4. Examples of (a) regular move decisions and (b) maintaining decisions for 1 s.

the policy network, allowing the agent to choose “no-op” and do nothing for a certain period if necessary. This means that our action space has 44 skills, plus an additional “no-op” action. This is significant because human play logs of *BAB* show that “no-op” actions take up the largest portion of skill usage among human players.

“No-op” decisions can be categorized as passive and active use cases. The passive use of “no-op” implies that an agent chooses “no-op” because there is no skill available for use. For example, when an agent is out of resources or is hit by an opponent’s CC skill, an agent has no option but to choose “no-op.” The active use of “no-op” means that an agent selects “no-op” strategically, even though other skills are available for use.

We discarded passive “no-op” data from both the training and evaluation phases because passive “no-ops” are not used deliberately by an agent. In addition, the method enables LSTM to reflect representations of longer time horizons because the data is not provided to the network. We show in the experiment section that skipping passive “no-ops” greatly improves learning efficiency. Note that this methodology is generally applicable to other domains where a set of available skills changes constantly and the “no-op” action is a valid option to choose.

### B. Maintain Move Decision

Although a single skill decision can have a substantial influence on the subsequent states, the effect of a single move decision is relatively limited. The reason is that the distance a character moves in a single time step (0.1 s) is very short considering its speed. In order for any moving decision to have a meaningful effect, the agent should make the same moving decision consecutively for several ticks in a row. This allows the agent to literally “move” and leads to changes in subsequent states and rewards. Therefore, it is difficult to train a move policy from the initial policy with random move decisions. Since the chance of a random policy making the same decision consecutively is very low, exploration is extremely limited. We therefore propose maintaining the move decision for a fixed number of time steps.

Fig. 4 shows how “maintain move decision” works with an example. If the agent selects a move action, it skips the move decision for the following  $n - 1$  time steps. This means that the agent maintains the same move decision for  $n$  steps in total. Note

that our method has different purpose from frame skip technique [15] in Atari domain. Frame skip technique was introduced for simulator’s efficiency. However, we cannot just skip the frames because skill decisions must still be made. Although we could not enjoy advantage in the simulator’s efficiency, “maintain move decision” still facilitates training and this is solely because maintaining the move decision increases the influence of a single move decision, as we will confirm with experiments. In this sense, “maintain move decision” rather can be viewed as “amplifying advantage” from [14].

## V. EXPERIMENTS

### A. Implementation Details

1) *Network*: The network is composed of LSTM-based architecture which has four heads with a shared state representation layer. Each head consists of  $\pi_{\text{skill}}$ ,  $Q_{\text{skill}}$ ,  $\pi_{\text{move,target}}$ , and  $Q_{\text{move,target}}$ .  $Q_{\text{skill}}$  and  $Q_{\text{move,target}}$  are used for the gradient update of  $\pi_{\text{skill}}$  and  $\pi_{\text{move,target}}$ , respectively. Before the network output goes into the softmax layer, a Boolean vector indicating the availability of each skill operates to make the output of unavailable skill to negative infinity.

2) *Algorithm*: We used actor-critic off-policy learning algorithm [26]. It enables us to deal with policy lag between the simulators and learner through truncated importance sampling. Moreover, we could also use the advantages of stochastic policy, which responds more stably to changes in the environment due to smooth policy updates and works well in the domain of games like rock-paper-scissors where deterministic policy is vulnerable to exploitation. For this specific algorithm, both  $\pi_{\text{skill}}$  and  $\pi_{\text{move,target}}$  are updated in an alternating manner with following gradient:

$$g_t^{\text{acer}} = \bar{\rho}_t \nabla_{\theta} \log \pi_{\theta}(a_t | x_t) [Q^{\text{ret}}(x_t, a_t) - V_{\theta_v}(x_t)] + E_{a \sim \pi} \left( \left[ \frac{\rho_t(a) - c}{\rho_t(a)} \right]_+ \nabla_{\theta} \log \pi_{\theta}(a | x_t) [Q_{\theta_v}(x_t, a) - V_{\theta_v}(x_t)] \right)$$

where  $\bar{\rho}_t = \min\{c, \rho_t\}$  with behavior policy  $\mu$  and importance sampling ratio  $\rho_t = (\pi(a_t | x_t) / \mu(a_t | x_t)) \cdot [x]_+ = x$  if  $x > 0$  and zero otherwise.

3) *Learning System*: In total, there are three learning processes with each learning process consisting of a learner and 100 simulators. Each learning process is largely similar to that proposed by [7]. The final agent is trained for two weeks, which is equivalent to four years of game play.

### B. Effect of Self-Play Curriculum With Three Styles

To demonstrate the effects of the proposed self-play curriculum, we trained agents with and without the proposed curriculum. A baseline agent was trained with the vanilla self-play curriculum without any style-related rewards (only win reward and HP reward were included) and a pool of past selves was used. Meanwhile, three agents with different styles were trained with the self-play curriculum using the shared pool that we proposed.

TABLE III  
WINNING RATE OF THREE STYLE OF AGENTS AGAINST  
BASELINE (1000 GAMES EACH)

	Aggressive	Balanced	Defensive	Average
Vs. Baseline	59.5%	63.8%	63.2%	62.2%

TABLE IV  
GENERALIZATION PERFORMANCE OF THREE STYLES OF AGENTS FOR BOTH  
WITH AND WITHOUT SHARED POOL (7000 GAMES EACH)

	Aggressive	Balanced	Defensive	Average
Shared	64.8%	79.6%	75.3%	73.6%
Ind.	64.7%	72.1%	56.5%	64.4%

Our aggressive, balanced and defensive agents<sup>2</sup> then played 1000 matches each against the baseline agent to measure the performance. As given in Table III, the agents that followed the learnings from our curriculum outperformed the baseline agent.

Next, we conducted an ablation study to observe how the shared pool helps generalization. We wanted to confirm whether an agent would be able to deal with opponents of unseen style, when it experienced only a limited range of opponents during training. Thus, we created three styles of agents trained in exactly the same manner, except that they had their own independent opponent pools. We denote the three types of agents using shared pools as  $\pi_{sh}^{agg}$ ,  $\pi_{sh}^{bal}$ , and  $\pi_{sh}^{def}$ , and three type of agents using independent pools as  $\pi_{ind}^{agg}$ ,  $\pi_{ind}^{bal}$ , and  $\pi_{ind}^{def}$ . All of six agents were trained for 5M steps (equivalent to six days) each.

Our assumption is that the agent trained with the shared pool is more robust when it faces opponents it has never encountered. Thus, we compared the winning rate of  $\pi_{sh}^{agg}$  vs.  $\{\pi_{ind}^{bal}, \pi_{ind}^{def}\}$  and  $\pi_{ind}^{agg}$  versus  $\{\pi_{sh}^{bal}, \pi_{sh}^{def}\}$ . This experimental setting is based on three key ideas. First,  $\pi_{sh}^{agg}$  and  $\pi_{ind}^{agg}$  have the same training settings except for sharing the pool. Second,  $\pi_{sh}^{agg}$  and  $\pi_{ind}^{agg}$  are evaluated against the same opponents. Finally, although  $\pi_{sh}^{agg}$  has encountered other styles from its pool, it has not confronted  $\{\pi_{ind}^{bal}, \pi_{ind}^{def}\}$ , for they were trained using independent opponent pools. If our assumption is correct,  $\pi_{sh}^{agg}$  should have a higher winning rate. It is to be noted that  $\pi_{sh}^{bal}$  and  $\pi_{ind}^{def}$  are not a single model, but ten models each sampled at the same fixed intervals from their pools. We then conducted the same experiments for the remaining two styles; the results are given in Table IV. As shown in the table, agents trained with shared pool outperform their counterparts.

Based on the data in Table IV, the effect of using a shared pool is marginal in the case of aggressive agents. It indicates that the strategy spaces in which trainings take place are similar whether or not various opponents are provided. This is related to the nature of fighting games in which one side should fight back if the other side approaches and initiates a brawl. Thus, in the case of an aggressive agent that attacks consistently, there is

<sup>2</sup>We measured how the average game length differs for each style because game length is a good proxy for assessing the degree of defensiveness of an agent's game play. The results were as follows: 66.6 s for the aggressive, 91.7 s for the balanced, and 179.9 s for the defensive agent.

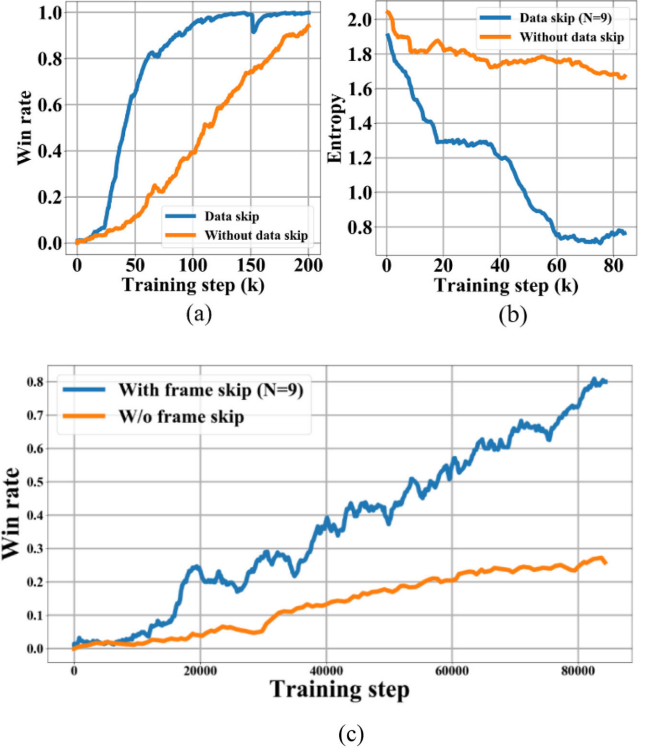


Fig. 5. Results of data skipping experiments. (a) Skipping passive “no-op.” (b) Maintain move decision. (c) Winning rates against built-in BAB AIs with and without “maintain move decision.”

a little difference in the experience regardless of the diversity of the opponent's fighting style.

### C. Effect of Discarding Passive “No-Op”

As discussed in the previous section, the “no-op” decision may be either active or passive. We conducted an experiment to investigate the effect of discarding such passive “no-op” data from learning. The sparring partner for the experiment was the built-in BAB AI, with a performance comparable to the top 20% of the players. We measured how fast agents learned to defeat it, and the results are shown in Fig. 5(a). If “no-op” ticks are discarded from the learning data, the winning rate reaches 80% after 70 k steps, whereas 170 k steps are required when “no-op” ticks are included. The amount of time steps required to reach 90% winning rate was reduced to half when passive “no-op” data was skipped. This experiment confirms that the training performance is improved by discarding passive “no-op” from the learning data.

### D. Effects of the “Maintain Move Decision”

To examine the effects of the “maintain move decision,” we developed two learning processes, with both processes involving learning on a self-play basis. One process makes a moving decision at every time step, while the other makes a moving decision and sends the same decision for nine more times in a row. We measured the entropy of the move policy to observe the effects. Entropy of the move policy for a given state  $s_t$  is as

follows:

$$H(s_t) = - \sum \pi_{\text{move}}(s_t) * \log \pi_{\text{move}}(s_t).$$

Generally, entropy gradually decreases as learning progresses. Fig. 5(b) shows that the entropy declines faster if the technique is applied. A noticeable difference was also observed in the quality of movement which the agent learned. Before the technique was applied, the agent did not make any improvement from random motion, but it learned to approach and retreat with data skip. In addition, Fig. 5(c) shows that the relative winning rate of the built-in *BAB* AI is made higher by applying “maintain move decision.”

The longer the decision was repeated, the agent’s reaction became less immediate, but the agent moved more consistently. In this case, we tested 1, 3, 5, and 10 ticks for maintaining time. A total of ten ticks (equivalent to 1 s) yielded the best performance.

## VI. PRO-GAMER EVALUATION

This section will address the results of both the pretest and the Blind Match, and conditions to ensure fairness for human players.

### A. Conditions for Fairness

1) *Reaction Time*: When humans confront an AI in a real-time fighting game, the most important factor that affects the result is the reaction time. Humans require some time to recognize the skill used by the opponent and to press a button by moving his/her hand. Therefore, we applied 230 ms of delay time on average to the AI’s decision-making process to add fairness against human players. The average response time of a human is approximately 270 ms [31]. However, a skilled pro-gamer’s response time will be shorter than the average. This amount of delay corresponds to the average reaction time of professional players in *BAB*.

2) *Classes and Skill Set*: There are 11 classes in *B&S*, and each class has unique characteristics. Since there exists relative superiority among classes, we fixed the class of both AI and pro-gamer as “destroyer” for all matches, training, and experiments conducted in Section V. Destroyer is a class that has an infighting style and steadily appears in the *B&S* world championship. Additionally, AI’s and pro-player’s skill trees were set as identical to ensure a fair match. The skill tree was chosen to match what the majority of users selected, based on the *BAB* user statistics.

3) *Evaluation Results*: We invited two prominent pro-gamers, Yuntae Son (GC Busan, Winner of 2017 B&S World Championship), and Shingyeom Kim (GC Busan, Winner of 2015 and 2016 B&S World Championship), to test our agents before the Blind Match. Note that the total number of games played is different for each style because the testers can play as many games as they want for each style. After the pre-test, we went for the Blind Match of 2018 World Championship. Our agents had matches against three pro-gamers: Nicholas Parkinson (EU), Shen Haoran (CHN), and Sungjin Choi (KOR). The video recording of the game highlights can be found at <https://goo.gl/7VUTzV>.

TABLE V  
FINAL SCORE OF AI VERSUS HUMAN

	Aggressive	Balanced	Defensive
Pro-Gamer 1	5-1	2-1	1-2
Pro-Gamer 2	4-0	2-4	4-1
Blind Match	2-0	1-2	0-2
Total	11-1 (92%)	5-7 (42%)	5-5 (50%)

The results of both the pre-test and the blind match are given in Table V. As can be seen from the table, the aggressive agent dominated the game, while the other two types of agents had rather intense games. Based on our interviews of pro-gamers, we concluded that the reason why the aggressive agent showed the best performance against a human player was because the aggressive agent delivered continuous attacks preventing the human player from taking short breaks in between moves required for decision-making; the AI, on the other hand, does not require these breaks.

## VII. CONCLUSION

In this article, using deep RL, we created AI agents that competed evenly with professional players in a 3-D real-time fighting game. To accomplish this, we proposed a method to guide the fighting style with reward shaping. With three styles of agents, we introduced a novel self-play curriculum to enhance generalization performance. Our self-play curriculum with three styles is effective against general or unknown opponents or when self-play training converges into an equilibrium that is not optimal. Agents will likely end up exploring only a small portion of the vast *BAB* space when they are limited to self-play. Thus, training with different styles obtained through reward shaping can be an effective way to find the optimal policy. We also proposed data-skipping techniques to improve data efficiency and enable efficient exploration. Consequently, our agents were able to compete with the best *BAB* pro-gamers in the world. The proposed training methods are generally applicable to other fighting games.

## APPENDIX

### A. State Space

A state input is a feature vector of length 606, rather than an image. Features include the following agent information: health point; skill point, distance from the opponent; distance from the arena wall in 16 directions; position and orientation; relative position from the opponent; remaining game time; remaining cooldown time and availability of each skill; abnormal condition (such as midair, stun, down, or kneel), etc. Abnormal condition refers to a state where the character’s actions are restricted. All such information is also available to human players. And all inputs are normalized into a value between 0 and 1.

### B. Feature Discretization

In the *BAB* environment, the current distance between oneself and the opponent determines whether some nontargeting skills



(i.e., skills that can be implemented regardless of the distance between oneself and the opponent) can successfully hit the opponent. For example, the destroyer's "axe sweep" skill can only hit the opponent if one's character is less than or equal to 5 m away from the opponent. Because the distance between oneself and the opponent is a feature that consists of a continuous value, the network has difficulty learning whether these non-targeting skills will be successful. Thus, in order to speed up the process, we used prior knowledge to discretize the following four crucial features.

- 1) Distance between oneself and the opponent: five levels (0–3 m, 3–5 m, 5–8 m, 8–16 m, and over 16 m).
- 2) Skill state: three levels (skill has been recently used, skill has been used but not recently, and skill has not been used).
- 3) Final six remaining ticks before crowd control loses effect: six levels (1 tick, 2 ticks, 3 ticks, 4 ticks, 5 ticks, and 6 ticks).
- 4) Elapsed time upon using "retreat" skill: six levels (1 tick, 2 ticks, 3 ticks, 4 ticks, 5 ticks, and 6 ticks).

To test the effectiveness of the feature discretization method, we designed an experiment comparing the success rates of counter attacks made against the "retreat" skill depending on the use of this technique. The "retreat" skill renders the player invincible for a certain amount of time at the cost of being unable to implement any kind of action for a slightly longer period which places them at a vulnerable position for a short time space. To be successful against *BAB* pro-gamers, an agent must learn how to seize this moment to apply crowd control against the opponent.

The agents were trained for 1M steps and tested for 20 games. Success rates were 75.0% and 13.9% with and without feature discretization, respectively. These results show that attacks against retreat skills became 5.4 times more likely to be successful by employing feature discretization.

### C. Action Space

#### 1) Skill

"Skills" refer to a number of unique actions defined in the game of *BAB*. They can be used separately or in a series to develop a myriad of strategies. The Destroyer class has 44 skills; every skill has its own function, such as dealing damage, applying crowd control, escaping, dashing, resisting, etc. Our agents make a skill decision out of 45 options including "no-op."

#### 2) Move and Target

"Move and target" refer to changing a player's position and orientation. Although human players can move in 8 cardinal directions and have 360° vision, not all combinations are necessary to achieve pro-level play. We reduced the action space into the following six options by combining move and target decisions.

- 1) Hold position and orientation.
- 2) Move forward facing the opponent.
- 3) Move to the right facing the opponent.
- 4) Move to the left facing the opponent.

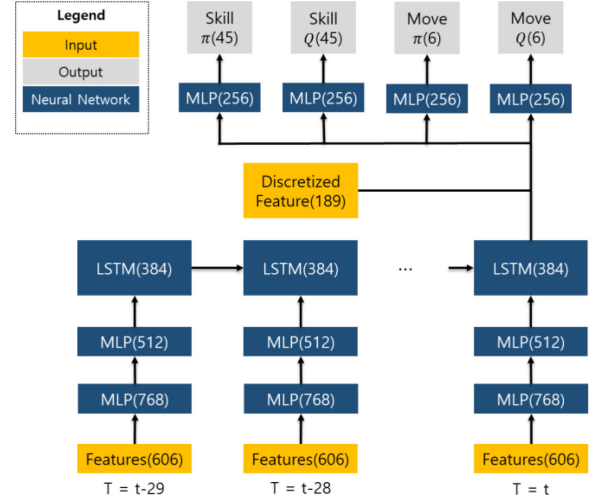


Fig. 6. Network architecture.

- 5) Move backward facing the opponent (results in slower movement but the player can react to the opponent's attack).
- 6) move forward facing away from the opponent (faster retreat but unable to react to the opponent's attack).

### D. Hyperparameters

TABLE VI  
HYPERPARAMETER SETTINGS

Name	Value
initial learning rate	1e-4
decay steps	33,000
decay rates	0.96
final learning rate	2e-5
optimizer	Adam
batch size	8 (240 transitions)
LSTM sequence length	30
send model to shared pool every C step	15,000
recent k models	10
replay size	800,000
entropy bonus coefficient	0.01
gradient clipping	20
Total HP reward weight	10
Win reward weight	10
Time reward weight per tick	-0.008~0.0
Distance reward weight per tick	-0.002~0.0

### E. Network Architecture

The network architecture employed in the *BAB* training process is described in Fig. 6. Rectified linear unit is used as the activation function in all multilayer perceptrons (MLPs) except



the final MLP. The agent selects a skill by sampling from the pointwise multiplication result of (the Destroyer's) skill probability vector and skill availability vector. The discretized features are constructed as shown in Appendix B and are concatenated with the LSTM output.

## REFERENCES

- [1] N. Brown and T. Sandholm, "Superhuman AI for multiplayer poker," *Science*, vol. 365, no. 6456, pp. 885–890, 2019.
- [2] L. Espeholt *et al.*, "IMPALA: Scalable distributed Deep-RL with importance weighted actor-learner architectures," in *Proc. 35th Int. Conf. Mach. Learn.*, 2018, vol. 80, pp. 1407–1416, 2018.
- [3] V. Firoiu, W. F. Whitney, and J. B. Tenenbaum, "Beating the world's best at Super Smash Bros. with deep reinforcement learning," 2017. [Online]. Available: <http://arxiv.org/pdf/1702.06230>
- [4] J. Heinrich and D. Silver, "Deep reinforcement learning from self-play in imperfect-information games," 2016. [Online]. Available: <http://arxiv.org/pdf/1603.01121>
- [5] M. Hessel, J. Modayil, H. Van Hasselt, T. Schaul, G. Ostrovski, and Dabney, "Rainbow: Combining improvements in deep reinforcement learning," in *Proc. 32nd AAAI Conf. Artif. Intell.*, 2018, pp. 3215–3222.
- [6] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [7] D. Horgan *et al.*, "Distributed prioritized experience replay," 2018, *arXiv:1803.00933*.
- [8] R. A. Howard *Dynamic Programming and Markov Processes*. Cambridge, MA, USA: MIT Press, 1960.
- [9] M. Ishihara, S. Ito, R. Ishii, T. Harada, and R. Thawonmas, "Monte-Carlo tree search for implementation of dynamic difficulty adjustment fighting game AIs having believable behaviors," in *Proc. IEEE Conf. Comput. Intell. Games*, 2018, pp. 1–8.
- [10] M. Jaderberg, W. M. Czarnecki, I. Dunning, L. Marris, G. Lever, and A. G. Castaneda, "Human-level performance in first-person multiplayer games with population-based deep reinforcement learning," *Science*, vol. 364, no. 6443, pp. 859–865, May 2019.
- [11] M. J. Kim and K. J. Kim, "Opponent modeling based on action table for MCTS-based fighting game AI," in *Proc. IEEE Conf. Comput. Intell. Games*, 2017, pp. 178–180.
- [12] Y. J. Li, H. Y. Chang, Y. J. Lin, P. W. Wu, and Y. C. FrankWang, "Deep reinforcement learning for playing 2.5-D fighting games," in *Proc. 25th IEEE Int. Conf. Image Process.*, 2018, pp. 3778–3782.
- [13] F. Lu, K. Yamamoto, L. H. Nomura, S. Mizuno, Y. Lee, and R. Thawonmas, "Fighting game artificial intelligence competition platform," in *Proc. IEEE 2nd Glob. Conf. Consum. Electron.*, 2013, pp. 320–323.
- [14] M. Mladenov, O. Meshi, J. Ooi, D. Schuurmans, and C. Boutilier, "Advantage amplification in slowly evolving latent-state environments," 2019. [Online]. Available: <http://arxiv.org/pdf/1905.13559>
- [15] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, and M. G. Bellemare, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [16] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Harley, and T. P. Lillicrap, "Asynchronous methods for deep reinforcement learning," in *Proc. 33rd Int. Conf. Int. Conf. Mach. Learn.*, 2016, vol. 48, pp. 1928–1937.
- [17] A. Y. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," in *Proc. 16th Int. Conf. Mach. Learn.*, 1999, vol. 99, pp. 278–287, 1999.
- [18] OpenAI, San Francisco, CA, USA, 2018. OpenAI fiv. [Online]. Available: <https://blog.openai.com/openai-five>
- [19] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," 2017. [Online]. Available: <http://arxiv.org/pdf/1707.06347>
- [20] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, and G. Van Den Driessche, "Mastering the game of Go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, 2016.
- [21] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, and A. Guez, "Mastering the game of Go without human knowledge," *Nature*, vol. 550, no. 7676, 2017, pp. 354–359.
- [22] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, and A. Guez, "Mastering Chess and Shogi by self-play with a general reinforcement learning algorithm," 2017. [Online]. Available: <http://arxiv.org/abs/1712.01815>
- [23] R. S. Sutton and A. G. Barto *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 2018.
- [24] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, and M. Yeo, "Starcraft II: A new challenge for reinforcement learning," 2017. [Online]. Available: <http://arxiv.org/pdf/1708.04782>
- [25] O. Vinyals, I. Babuschkin, J. Chung, M. Mathieu, M. Jaderberg, and W. M. Czarnecki, "Alphastar: Mastering the real-time strategy game Starcraft II. DeepMind blog" Feb. 2019.
- [26] Z. Wang *et al.*, "Sample efficient actor-critic with experience replay," in *Proc. Int. Conf. Learn. Representations*, 2017.
- [27] S. Yoshida, M. Ishihara, T. Miyazaki, Y. Nakagawa, T. Harada, and R. Thawonmas, "Application of Monte-Carlo tree search in a fighting game AI," in *Proc. IEEE 5th Glob. Conf. Consum. Electron.*, 2016, pp. 1–2.
- [28] J. Li *et al.*, "Suphx: Mastering Mahjong with deep reinforcement learning" 2020, *arXiv:2003.13590*.
- [29] O. Vinyals *et al.*, "Grandmaster level in Starcraft II using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.
- [30] R. Gibbons *A Primer Game Theory*, vol. 56, 1992.
- [31] Reaction Time Statistics. Accessed: Jan. 15, 2021. [Online]. Available: <https://humanbenchmark.com/tests/reactiontime/statistics>
- [32] Crowd Control. Accessed: Jan. 15, 2021, [Online]. Available: [https://en.wikipedia.org/wiki/Crowd\\_control\\_\(video\\_games\)](https://en.wikipedia.org/wiki/Crowd_control_(video_games))
- [33] J. Harrington, *Games, Strategies and Decision Making*. New York, NY, USA: Macmillan, 2009.
- [34] J. Heinrich, M. Lanctot, and D. Silver, "Fictitious self-play in extensive-form games," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 805–813.
- [35] S. Sukhbaatar, Z. Lin, I. Kostrikov, G. Synnaeve, A. Szlam, and R. Fergus, "Intrinsic motivation and automatic curricula via asymmetric self-play," in *Proc. Int. Conf. Learn. Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=SkT5Yg-RZ>