



Universiteit  
Leiden  
The Netherlands

## Searching by learning: Exploring artificial general intelligence on small board games by deep reinforcement learning

Wang, H.

### Citation

Wang, H. (2021, September 7). *Searching by learning: Exploring artificial general intelligence on small board games by deep reinforcement learning*. Retrieved from <https://hdl.handle.net/1887/3209232>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/3209232>

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <https://hdl.handle.net/1887/3209232> holds various files of this Leiden University dissertation.

**Author:** Wang, H.

**Title:** Searching by learning: Exploring artificial general intelligence on small board games by deep reinforcement learning

**Issue Date:** 2021-09-07

# Chapter 6

## Adaptive Warm-Start AlphaZero-like Self-play

### 6.1 Introduction

Following Chapter 5, in this chapter, we will further propose an adaptive warm-start method on AlphaZero-like deep reinforcement learning framework.

The combination of online MCTS [7] in self-play and offline neural network training has been widely applied as a deep reinforcement learning technique, in particular for solving game-related problems by means of the AlphaGo series programs [10, 16, 17]. The approach of this paradigm is to use game playing records from self-play by MCTS as training examples to train the neural network, whereas this trained neural network is used to inform the MCTS value and policy. Note that in contrast to AlphaGo Zero or AlphaZero, the original AlphaGo also uses large amounts of expert data to train the neural network and a fast rollout policy together with the policy provided by neural network to guide the MCTS search.

However, although the transition from a combination of using expert data and self-play (AlphaGo) to only using self-play (AlphaGo Zero and AlphaZero) appears to have only positive results, it does raise some questions.

The first question is: ‘should all human expert data be abandoned?’ In other games we have seen that human knowledge is essential for mastering complex

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY

---

games, such as StarCraft [23]. Then when should expert data be taken into consideration while training neural networks?

The second question is: ‘should the fast rollout policy be abandoned?’ Chapter 5 has proposed to use warm-start search enhancements **at the start phase** in AlphaZero-like self-play, which improves performance in 3 small board games. Instead of only using the neural network for value and policy, in the first few iterations, classic rollout can be used (or RAVE, or a combination, or a combination with the neural network). This can improve training especially at the start phase of self-play training.

In fact, the essence of the warm-start search enhancement is to re-generate expert knowledge in the start phase of self-play training, to reduce the cold-start problem of playing against untrained agents. The method uses rollout (which can be seen as experts) instead of a randomly initialized neural network, up until a number of  $I'$  iterations, when it switches to the regular value network. In their experiments, the  $I'$  was fixed at 5. Obviously, a fixed  $I'$  may not be optimal. Therefore, in this work, we propose an adaptive switch method. The method uses an **arena** in the self-play stage (see Algorithm 6), where the search enhancement and the default MCTS are matched, to judge whether to switch or not. With this mechanism, we can dynamically switch off the enhancement if it is no longer better than the default MCTS player, as the neural network is being trained.

Our main contributions can be summarized as follows:

1. Warm-start method improves the Elo of AlphaZero-like self-play in small games, but it introduces a new hyper-parameter. Adaptive warm-start further improves performance and removes the hyper-parameter.
2. For deep games (with a small branching factor) warm-start works better than for shallow games. This indicates that the effectiveness of warm-start method may increase for larger games.

The rest of this chapter is designed as follows. An overview of the most relevant literature is given in Sect. 6.2. Before proposing our adaptive switch method in Sect. 6.4, we describe the warm-start AlphaZero-like self-play algorithm in Sect. 6.3. Thereafter, we set up the experiments in Sect. 6.5 and present their results in Sect. 6.6. Finally, we conclude the chapter and discuss future work.

## 6.2 Related Work

There are a lot of early successful works in reinforcement learning [6], e.g. using temporal difference learning with a neural network to play backgammon [70]. MCTS has also been well studied, and many variants/enhancements were designed to solve problems in the domain of sequential decisions, especially on games. For example, enhancements such as RAVE and All Moves as First (AMAF) have been conceived to improve MCTS [18, 24]. The AlphaGo series algorithms replace the table based model with a deep neural network based model, where the neural network has a policy head (for evaluating of a state) and a value head (for learning a best action) [34], enabled by the GPU hardware development. Thereafter, the structure that combines MCTS with neural network training has become a typical approach for reinforcement learning tasks and many successful applications [37, 89] of this kind model-based deep reinforcement learning [68]. Comparing AlphaGo with AlphaGo Zero and AlphaZero, the latter did not use any expert data to train the neural network, and abandoned the fast rollout policy for improving the MCTS on the trained neural network. Therefore, all training data is generated purely by self-play, which is also a very important feature of reinforcement learning. We base our work on an open reimplementaion of AlphaZero, AlphaZero General [65].

There are many interesting works on self-play in reinforcement learning [70, 74, 107]. Temporal difference learning for acquiring position evaluation in small board Go with co-evolution has been compared to self-play [107]. These works demonstrated the impressive results for self-play and emphasized its importance.

Within a GGP framework, in order to improve training examples efficiency, [31] assessed the potential of classical Q-learning by introducing MCS enhancements. In an AlphaZero-like self-play framework, [108] used domain-specific features and optimizations, starting from random initialization and no preexisting data, to accelerate the training.

However, AlphaStar, the acclaimed algorithm for beating human professionals at StarCraft [23], went back to utilizing human expert data, thereby suggesting that this is still an option at the start phase of training. Apart from this, there are few studies on applying MCTS enhancements in AlphaZero-like self-play. Only [35] (presented in Chapter 5), which proposed a warm-start search enhancement method, pointed out the promising potential of utilizing MCTS enhancements (like a rollout policy) to re-generate expert data at the start phase of training. Our approach differs from AlphaStar, as we generate expert data using MCTS enhancements other than collecting it from humans; further, compared to

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY

---

the static warm-start in Chapter 5, we propose an adaptive method to control the iteration length of using such enhancements instead of a fixed  $I'$ .

### 6.3 Warm-Start AlphaZero Self-play

Based on Chapter 5, we will now briefly recall the warm-start enhancement method.

#### 6.3.1 The Algorithm Framework

Based on [10, 34] and Chapter 5, the core of AlphaZero-like self-play (see Algorithm 5) is an iterative loop which consists of three different stages within the single iteration as follows:

1. **self-play**: The first stage is playing several games against with itself to generate training examples.
2. **neural network training**: The second stage is feeding the neural network with training examples (generated in the first stage) to train a new model.
3. **arena comparison**: The last stage is employing a tournament to compare the newly trained model and the old model to decide whether to update or not.

The detail description of these 3 stages can be found in Chapter 5 (Sect. 5.3.1). Note that in the Algorithm 5, line 5, a fixed  $I'$  is employed to control whether to use neural network MCTS or MCTS enhancements, the  $I'$  should be set as relatively smaller than  $I$ , which is known as warm-start search. The MCTS algorithm and MCTS enhancements will be introduced in next subsections.

#### 6.3.2 MCTS

Classical MCTS has shown successful performance to solve complex games, by taking random samples in the search space to evaluate the state value. Basically, the classical MCTS algorithm can be divided into 4 stages, which are known as *selection*, *expansion*, *rollout* and *backpropagate* [7]. However, for the default MCTS in AlphaZero-like self-play (eg. our Baseline), the neural network directly informs the MCTS state policy and value to guide the search instead of running a rollout (see Algorithm 4).

---

**Algorithm 5** Warm-start AlphaZero-like Self-play Algorithm
 

---

```

1: Randomly initialize  $f_\theta$ , assign retrain buffer  $D$ 
2: for iteration=1, ...,  $I'$ , ...,  $I$  do
3:   for episode=1, ...,  $E$  do ▷ self-play
4:     for t=1, ...,  $T'$ , ...,  $T$  do
5:       if  $I \leq I'$  then  $\pi_t \leftarrow$  MCTS Enhancement
6:       else  $\pi_t \leftarrow$  default MCTS
7:       if  $t \leq T'$  then  $a_t =$  randomly select on  $\pi_t$ 
8:       else  $a_t = \arg \max_a(\pi_t)$ 
9:       executeAction( $s_t, a_t$ )
10:       $D \leftarrow (s_t, \pi_t, z_t)$  with outcome  $z_{t \in [1, T]}$ 
11:      Sample minibatch  $(s_j, \pi_j, z_j)$  from  $D$  ▷ training
12:      Train  $f_{\theta'} \leftarrow f_\theta$ 
13:       $f_\theta = f_{\theta'}$  if  $f_{\theta'}$  is better, using default MCTS ▷ arena
14: return  $f_\theta$ ;

```

---

### 6.3.3 MCTS enhancements

In this chapter, we adopt the same two individual enhancements and three combinations to improve neural network training as were used in Chapter 5. Here we briefly recall them again with a possible minor change for *weight* calculation.

**Rollout** is running a classic MCTS random rollout to get a value that provides more meaningful information than a value from random initialized neural network.

**RAVE** is a well-studied enhancement to cope with the cold-start of MCTS in games like Go [18], where the playout-sequence can be transposed. The core idea of RAVE is using AMAF to update the state visit count  $N_{rave}$  and Q-value  $Q_{rave}$ , which are written as:  $N_{rave}(s_{t_1}, a_{t_2}) \leftarrow N_{rave}(s_{t_1}, a_{t_2}) + 1$ ,  $Q_{rave}(s_{t_1}, a_{t_2}) \leftarrow \frac{N_{rave}(s_{t_1}, a_{t_2}) * Q_{rave}(s_{t_1}, a_{t_2}) + v}{N_{rave}(s_{t_1}, a_{t_2}) + 1}$ , where  $s_{t_1} \in VisitedPath$ , and  $a_{t_2} \in A(s_{t_1})$ , and for  $\forall t < t_2, a_t \neq a_{t_2}$ . The P-UCT of RAVE is calculated as follows:

$$PUCT_{rave}(s, a) = (1 - \beta) * U(s, a) + \beta * U_{rave}(s, a) \quad (6.1)$$

where

$$U_{rave}(s, a) = Q_{rave}(s, a) + c * P(s, a) \frac{\sqrt{N_{rave}(s, \cdot)}}{N_{rave}(s, a) + 1} \quad (6.2)$$

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY

---

and

$$\beta = \sqrt{\frac{\textit{equivalence}}{3 * N(s, \cdot) + \textit{equivalence}}} \quad (6.3)$$

The value of equivalence is usually set to the number of MCTS simulations (i.e  $m=100$  in our experiments).

**RoRa** is the combination which simply adds the random rollout to enhance RAVE.

**WRo** introduces a weighted sum of rollout value and the neural network value as the return value to guide MCTS. In our experiments,  $v(s)$  is computed as follows:

$$v(s) = (1 - \textit{weight}) * v_{\textit{network}} + \textit{weight} * v_{\textit{rollout}} \quad (6.4)$$

**WRoRa** also employs a weighted sum to combine the value from the neural network and the value of RoRa. The  $v(s)$  for MCTS search in WRoRa is computed as follows:

$$v(s) = (1 - \textit{weight}) * v_{\textit{network}} + \textit{weight} * v_{\textit{rora}} \quad (6.5)$$

Different from Chapter 5, since there is no pre-determined  $I'$ , in our work,  $\textit{weight}$  is simply calculated as  $1/i, i \in [1, I]$ , where  $i$  is the current iteration number.

### 6.4 Adaptive Warm-Start Switch Method

The fixed  $I'$  to control the length of using warm-start search enhancements as suggested in Chapter 5, but seems to require different parameter values for different games. In consequence, a costly tuning process would be necessary for each game. Thus, an adaptive method would have multiple advantages.

We notice that the core of the warm-start method is re-generating expert data to train the neural network at the start phase of self-training to avoid learning from weak (random or near random) self-play. We suggest to stop the warm-start when the neural network is on average playing stronger than the enhancements. Therefore, in the self-play, we employ a tournament to compare the standard AlphaZero-like self-play model (Baseline) and the enhancements (see Algorithm 6). The switch occurs once the Baseline MCTS wins more than 50%. In order to avoid spending too much time on this, these arena game records will directly be used as training examples, indicating that the training data is played by the enhancements and the Baseline. This scheme enables to switch at individual points in time for different games and even different training runs.



---

**Algorithm 6** Adaptive Warm-Start Switch Algorithm

---

```

1: Initialize  $f_\theta$  with random weights; Initialize retrain buffer  $D$ , Switch $\leftarrow$ False,  $r_{mcts} \leftarrow 0$ 
2: for iteration=1, ...,  $I$  do ▷ no  $I'$ 
3:   if not Switch then ▷ not switch
4:     for episode=1, ...,  $E$  do ▷ arena with enhancements
5:       for t=1, ...,  $T'$ , ...,  $T$  do
6:         if episode  $\leq E/2$  then
7:           if t is odd then  $\pi_t \leftarrow$  MCTS Enhancement
8:           else  $\pi_t \leftarrow$  default MCTS
9:         else
10:          if t is odd then  $\pi_t \leftarrow$  default MCTS
11:          else  $\pi_t \leftarrow$  MCTS Enhancement
12:          if  $t \leq T'$  then  $a_t =$  randomly select on  $\pi_t$ 
13:          else  $a_t = \arg \max_a(\pi_t)$ 
14:          executeAction( $s_t, a_t$ )
15:           $D \leftarrow (s_t, \pi_t, z_t)$  with outcome  $z_{t \in [1, T]}$ 
16:           $r_{mcts} +=$  reward of default MCTS in this episode
17:        else ▷ switch
18:          for episode=1, ...,  $E$  do ▷ purely self-play
19:            for t=1, ...,  $T'$ , ...,  $T$  do
20:               $\pi_t \leftarrow$  default MCTS
21:              if  $t \leq T'$  then  $a_t =$  randomly select on  $\pi_t$ 
22:              else  $a_t = \arg \max_a(\pi_t)$ 
23:              executeAction( $s_t, a_t$ )
24:               $D \leftarrow (s_t, \pi_t, z_t)$  with outcome  $z_{t \in [1, T]}$ 
25:          Set Switch $\leftarrow$ True if  $r_{mcts} > 0$ , and set  $r_{mcts} \leftarrow 0$ 
26:          Sample minibatch  $(s_j, \pi_j, z_j)$  from  $D$  ▷ training
27:          Train  $f_{\theta'} \leftarrow f_\theta$ 
28:           $f_\theta = f_{\theta'}$  if  $f_{\theta'}$  is better, using default MCTS ▷ arena
29: return  $f_\theta$ ;

```

---

## 6.5 Experimental Setup

Since Chapter 5 only studied the winrate of single rollout and RAVE against a random player, this can be used as a test to check whether rollout and RAVE work. However, it does not reveal any information about relative playing strength, which is necessary to explain how good training examples provided by MCTS enhancements actually are. Therefore, at first we let all 5 enhancements and the baseline MCTS (in this test, the neural network for each player is randomly

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY

---

initialized) play 100 games with each other on the same 3 games ( $6 \times 6$  Connect Four, Othello and Gobang, game description can be found in Chapter 4) in order to investigate the relative playing strength of each pair.

In the second experiment, we tune the fixed  $I'$ , where  $I' \in \{1, 3, 5, 7, 9\}$ , for different search enhancements, based on Algorithm 5 to play  $6 \times 6$  Connect Four.

In our last experiment, we use new adaptive switch method Algorithm 6 to play  $6 \times 6$  Othello, Connect Four and Gobang. We set parameters values according to Table 4.1. The parameter choices are based on [33].

Our experiments are run on a high-performance computing (HPC) server, which is a cluster consisting of 20 CPU nodes (40 TFlops) and 10 GPU nodes (40 GPU, 20 TFlops CPU + 536 TFlops GPU). We use small versions of games ( $6 \times 6$ ) in order to perform a medium number of repetitions. In the following, our figures show error bars of 8 runs, of 100 iterations of self-play. Each single run is deployed in a single GPU which takes several days for different games.

## 6.6 Results

We list results for a tournament of Baseline and enhancements (Table 6.1). Digging deeper, we also report the effect of the hyper-parameter  $I'$  (Fig 6.1). And results for the adaptive warm-start switch are shown in Table 6.2, Fig 6.2 and Fig 6.3.

### 6.6.1 MCTS vs MCTS Enhancements

Here, we compare the Baseline player (the neural network is initialized randomly which can be regarded as an arena in the first iteration self-play) to the other five MCTS enhancements players on 3 different games. Each pair performs 100 repetitions. In Table 6.1, we can see that for Connect Four, the highest winrate is achieved by WRoRa, the lowest by Rave. Except Rave, others are all higher than 50%, showing that the enhancements (except Rave) are better than the untrained Baseline. In Gobang, it is similar, Rave is the lowest, RoRa is the highest. But the winrates are relatively lower than that in other 2 games. It is interesting that in Othello, all winrates are relatively the highest compared to the 2 other games (nearly 100%), although Rave still achieves the lowest winrate which is higher than 50%.

One reason that enhancements work best in Othello is that the Othello game tree is the longest and narrowest (low branching factor). Enhancements like Rollout

**Table 6.1:** Results of comparing default MCTS with Rollout, Rave, RoRa, WRo and WRoRa, respectively on the three games with random neural network, weight as  $1/2$ ,  $T'=0$ , win rates in percent (row vs column), each pair played 100 games.

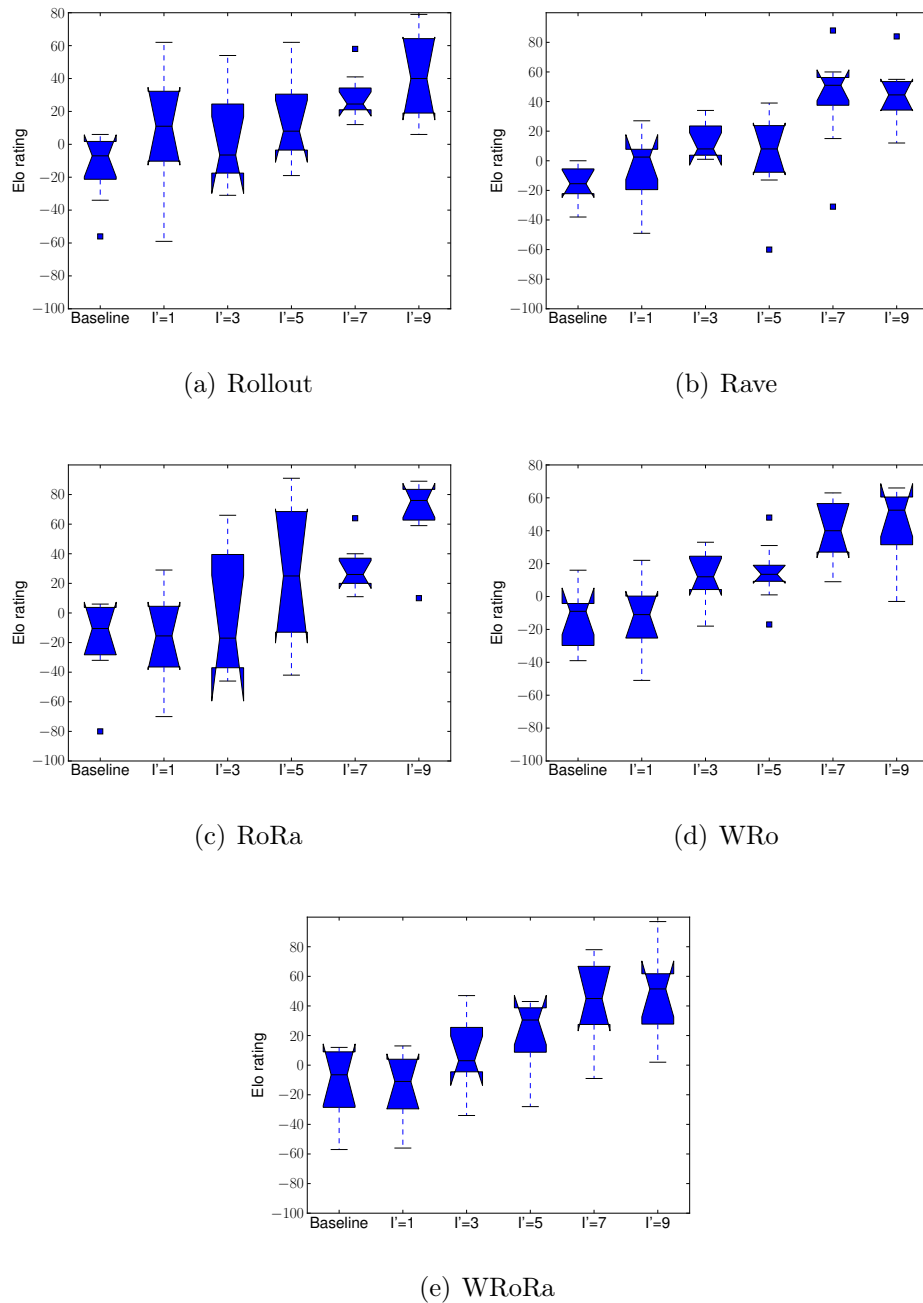
	Default MCTS		
	ConnectFour	Othello	Gobang
Rollout	64	93	65
Rave	27.5	53	43
RoRa	76	98	70
WRo	82	96	57
WRoRa	82.5	99	62

can provide relatively accurate estimations for these trees. In contrast, Gobang has the shortest game length and the most legal action options. Enhancements like Rollout do not contribute much to the search in short but wide search tree with limited MCTS simulation. As in shorter games it is more likely to reach a terminal state, both Baseline and enhancements get the true result. Therefore, in comparison to MCTS, enhancements like Rollout work better while it does not terminate too fast. Rave is filling more state action pairs based on information from the neural network, its weaknesses at the beginning are more emphasized. After some iterations of training, the neural network becomes smarter, and Rave can therefore enhance the performance as shown in Chapter 5 .

### 6.6.2 Fixed $I'$ Tuning

Taking Connect Four as an example, in this experiment we search for an optimal fixed  $I'$  value, utilizing the warm-start search method proposed in Chapter 5. We set  $I'$  as 1, 3, 5, 7, 9 respectively (the value should be relatively small since the enhancement is only expected to be used at the start phase of training). The Elo ratings of each enhancements using different  $I'$  are presented in Fig 6.1. The Elo ratings are calculated based on the tournament results using a Bayesian Elo computation system [82], same for Fig 6.3. We can see that for Rave and WRoRa, it turns out that  $I' = 7$  is the optimal value for fixed  $I'$  warm-start framework, for others, it is still unclear which value is the best, indicating that the tuning is inefficient and costly.

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY



**Figure 6.1:** Elo ratings for different warm-start phase iterations with different search enhancement on  $6 \times 6$  Connect Four

### 6.6.3 Adaptive Warm-Start Switch

In this final experiment, we apply the newly suggested adaptive switch warm-start search enhancement method and compare it to the fixed  $I'$ . We are especially interested in the averages and variances of the switching times that result from adaptive switching.

We train models with the parameters in Table 4.1 and then let them compete against each other in different games. In addition, we record the specific iteration number where the switch occurs for every training run and the corresponding self-play arena rewards of MCTS before this iteration. A statistic of the iteration number for 3 games is shown in Table 6.2.

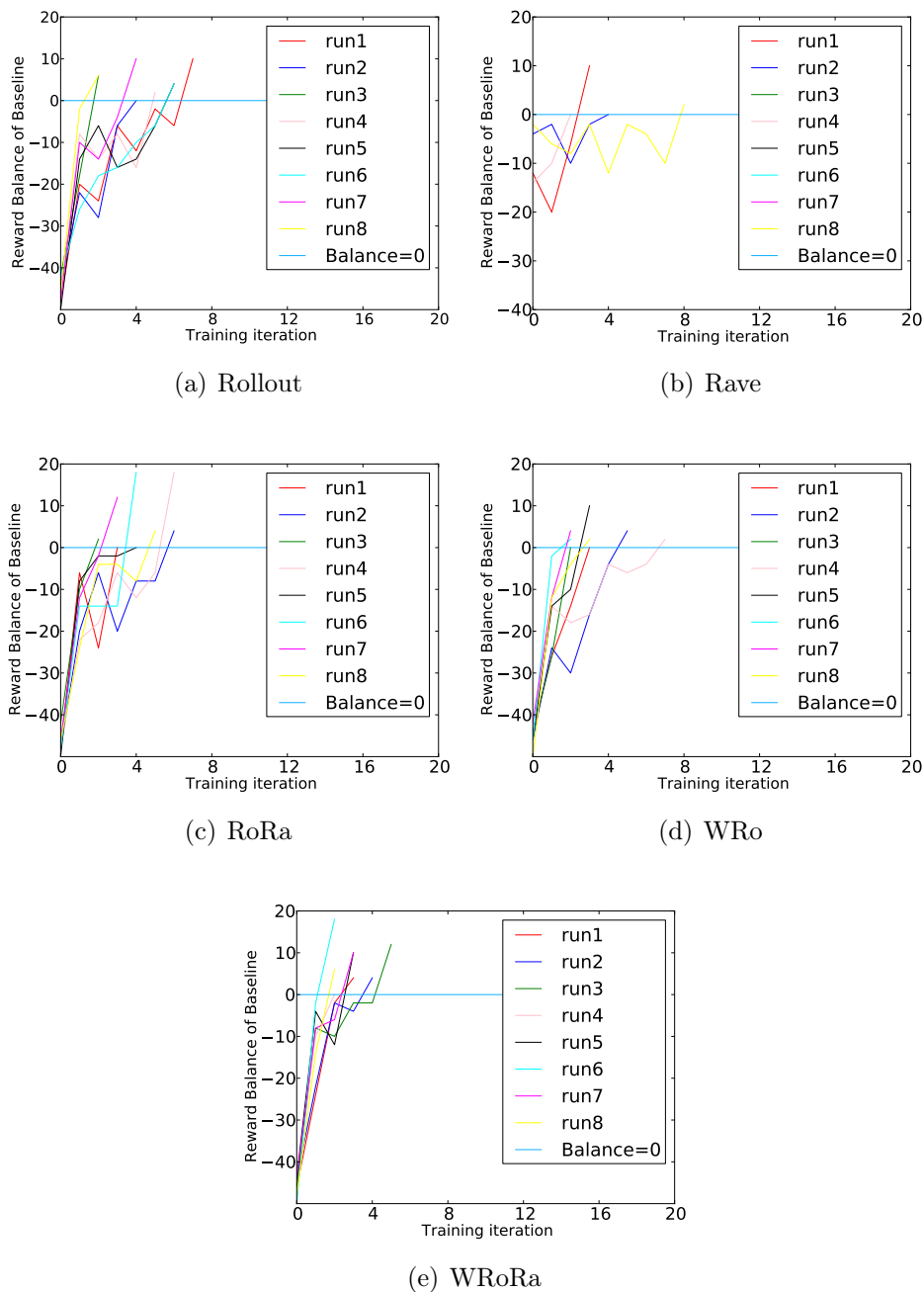
**Table 6.2:** Switching iterations for training on different games with different enhancements over 8 repetitions (average iteration number  $\pm$  standard deviation)

	Connect Four	Othello	Gobang
Rollout	6.625 $\pm$ 3.039	5.5 $\pm$ 1.732	1.375 $\pm$ 0.484
Rave	2.375 $\pm$ 1.218	3.125 $\pm$ 2.667	1.125 $\pm$ 0.331
RoRa	7.75 $\pm$ 4.74	5.125 $\pm$ 1.364	1.125 $\pm$ 0.331
WRo	4.25 $\pm$ 1.561	4.375 $\pm$ 1.654	1.125 $\pm$ 0.331
WRoRa	4.375 $\pm$ 1.576	4.0 $\pm$ 1.0	1.25 $\pm$ 0.433

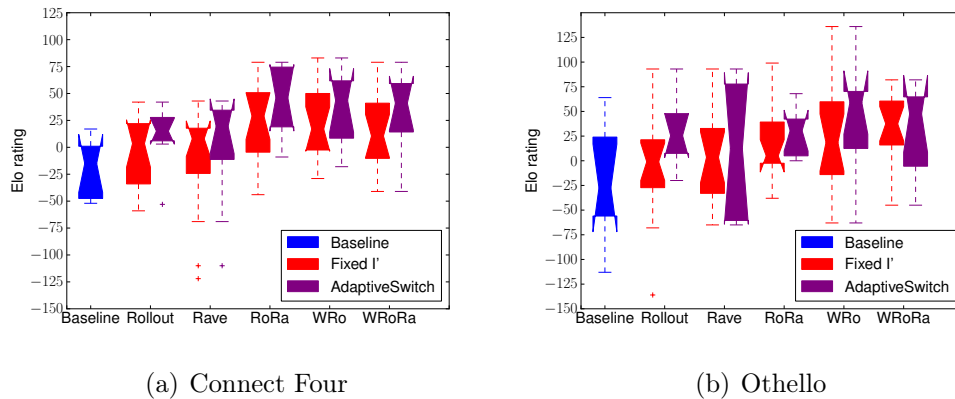
The table shows that, generally, the iteration number is relatively small compared to the total length of the training (100 iterations), and in these small games the neural network is quickly getting stronger. Besides, not only for different games, the switch iteration is different, but also for different training runs on the same game, the switch iteration also varies. This is because for different training runs, the neural network training progresses differently (we already start from different random initializations). Therefore, a fixed  $I'$  can not be used for each specific training. Note that for Gobang, a game with a large branching factor, with the default setting, it always switches at the first iteration. Therefore, we also test with larger  $m = 200$ , thereby providing more time to the MCTS. With this change, there are several runs keeping the enhancements see Table 6.2, but it still shows a small influence on this game.

In addition, we show the arena results (wins of default MCTS minus wins of enhancement) in each training iteration before switch happens in each run over 8 repetitions on Othello as an example in Fig 6.2. In most curves, we can

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY



**Figure 6.2:** Reward balances of default MCTS while competing with different enhancements in self-play arena for  $6 \times 6$  Othello. Exceeding 0 means default MCTS defeats the enhancement, switch occurs.



**Figure 6.3:** Comparison of adaptive switch method versus fixed  $I'$  based on a full tournament for  $6 \times 6$  Connect Four and Othello

see improving reward balances achieved by default MCTS since it is getting stronger.

More importantly, we collect all trained models based on our adaptive method, and let them compete with the models trained using fixed  $I' = 5$  in a full round-robin tournament where each 2 players play 20 games.

From Fig 6.3, we see that, generally, on both Connect Four and Othello, all fixed  $I'$  achieve higher Elo ratings than the Baseline, which was also reported in Chapter 5. And all adaptive switch models also perform better than the Baseline. Besides, for each enhancement, it is important that the Elo ratings of the adaptive switch models are higher than for the fixed  $I'$  method, which suggests that our adaptive switch method leads to better performance than the fixed  $I'$  method when controlling the warm-start iteration length. Specifically, we find that for Connect Four, WRo and RoRa achieve the higher Elo Ratings (see Fig 6.3(a)) and for Othello, WRoRa performs best (see Fig 6.3(b)), which reproduces the consistent conclusion (at least one combination enhancement performs better in different games) as Chapter 5).

In addition, for Connect Four, comparing the tuning results in Fig 6.1 and the *switch iterations* by our method in Table 6.2, we find that our method generally needs a shorter warm-start phase than employing a fixed  $I'$ . The reason could be that in our method, there are always 2 different players playing the game, and they provide more diverse training data than a pure self-play player. In consequence, the neural network also improves more quickly, which is highly desired.

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY

---

Note that while we use the default parameter setting for training in the Gobang game, the *switch* occurs at the first iteration. And even though we enlarge the simulation times for MCTS, only a few training runs shortly keep using the enhancements. We therefore presume that it is meaningless to further perform the tournament comparison for Gobang.

### 6.7 Summary

Since AlphaGo Zero's results, self-play has become a default approach for generating training data tabula rasa, disregarding other information for training. However, if there is a way to obtain better training examples from the start, why not use them, as has been done recently in StarCraft (see DeepMind's AlphaStar [23]). In addition, Chapter 5 investigated the possibility of utilizing MCTS enhancements to improve AlphaZero-like self-play. They embed Rollout, RAVE and combinations as enhancements at the start period of iterative self-play training and tested this on small board games. Since the neural network and the MCTS statistics are initialized to random weights and zero, self-play suffers from a cold-start problem, and starting from scratch can lead to unstable learning at the start of the training. These problems can be cured by feeding human expert data or running MCTS enhancements or similar methods in order to generate expert data for training the neural network before switching to pure self-play. (Not unlike RAVE warm-starts the winrate statistics of the original MCTS in 2007.)

Confirming Chapter 5, we find that finding an optimal value of fixed  $I'$  is difficult, therefore, we propose an adaptive method for deciding when to switch. We also use Rollout, RAVE, and combinations with network values to quickly improve MCTS tree statistics (using RAVE) with meaningful information (using Rollout) before we switch to Baseline-like self-play training. We employed the same games, namely the 6x6 versions of Gobang, Connect Four, and Othello. In these experiments, we find that, for different games, and even different training runs for the same game, the new adaptive method generally switches at different iterations. This indicates the noise in the neural network training progress for different runs. After 100 self-play iterations, we still see the effects of the warm-start enhancements as playing strength has improved in many cases, and for all enhancements, our method performs better than the method proposed in Chapter 5 with  $I'$  set to 5. In addition, some conclusions are consistent to Chapter 5, for example, there is also at least one combination that performs better.



The new adaptive method works especially well on Othello and Connect Four, "deep" games with a moderate branching factor, and less well on Gobang, which has a larger branching factor. In the self-play arena, the default MCTS is already quite strong, and for games with a short and wide episode, the MCTS enhancements do not benefit much. Short game lengths reach terminal states early, and MCTS can use the true reward information more often, resulting in a higher chance of winning. Since, Rollout still needs to simulate, with a limited simulation count it is likely to not choose a winning terminal state but a state that has the same average value as the terminal state. In this situation, in a short game episodes, MCTS works better than the enhancement with  $T'=15$ . With ongoing training of the neural network, both players become stronger, and as the game length becomes longer,  $I' = 5$  works better than the the Baseline.

Our experiments are with small games. Adaptive warm-start works best in deeper games, suggesting a larger benefit for bigger games with deeper lines. Future work includes larger games with deeper lines, and using different but stronger enhancements to generate training examples. Beside, it is also promising to apply the adaptive warm-start idea to master single agent or multi-agent deep reinforcement learning problems.

## 6. ADAPTIVE WARM-START ALPHAZERO-LIKE SELF-PLAY