# Asymmetric Move Selection Strategies in Monte-Carlo Tree Search: Minimizing the Simple Regret at Max Nodes

Yun-Ching Liu
Graduate School of Engineering
University of Tokyo
cipherman@logos.t.u-tokyo.ac.jp

Yoshimasa Tsuruoka
Graduate School of Engineering
University of Tokyo
tsuruoka@logos.t.u-tokyo.ac.jp

*Abstract*—The combination of multi-armed bandit (MAB) algorithms with Monte-Carlo tree search (MCTS) has made a significant impact in various research fields. The UCT algorithm, which combines the UCB bandit algorithm with MCTS, is a good example of the success of this combination. The recent breakthrough made by AlphaGo, which incorporates convolutional neural networks with bandit algorithms in MCTS, also highlights the necessity of bandit algorithms in MCTS. However, despite the various investigations carried out on MCTS, nearly all of them still follow the paradigm of treating every node as an independent instance of the MAB problem, and applying the same bandit algorithm and heuristics on every node. As a result, this paradigm may leave some properties of the game tree unexploited. In this work, we propose that max nodes and min nodes have different concerns regarding their value estimation, and different bandit algorithms should be applied accordingly. We develop the Asymmetric-MCTS algorithm, which is an MCTS variant that applies a simple regret algorithm on max nodes, and the UCB algorithm on min nodes. We will demonstrate the performance of the Asymmetric-MCTS algorithm on the game of $9 \times 9$ **Go**, $9 \times 9$ **NoGo, and Othello.**

## I. Introduction

Monte-Carlo Tree Search (MCTS) has made a significant impact on various fields in AI, especially on the field of computer Go [3]. The key factor to the success of MCTS lies in its combination with bandit algorithms, which solves the multi-armed bandit problem (MAB) [4]. The MAB problem is a problem where the agent needs to decide whether it should act optimally based on current available information (*exploitation*), or gather more information at the risk of suffering losses incurred by performing suboptimal actions (*exploration*) [1]. One of the most widely used MCTS variants is the UCT algorithm, which simply applies the UCB algorithm to every node in the tree [5]. The development of MCTS in recent years can be broadly classified into two main directions: one is the integration of knowledge learnt offline, and the other is increasing the effectiveness of the knowledge accumulated online.

The integration of offline knowledge was mainly focused on using logistic models to improve the quality of the simulations [13][14]. Recently, a lot of effort has been put into the training of convolutional neural networks and combining them with MCTS in computer Go [17][18]. A breakthrough was made by the program *AlphaGo*, which essentially combines convolutional neural networks with the PUCB bandit algorithm [16], and has beaten a top human professional player Lee Sedol in a five-game challenge match [18].

On the other hand, various investigations in increasing the effectiveness of online knowledge have also been carried out. One of them is using various bandit algorithms with MCTS, especially the bandit algorithms that solve the *pure exploration* MAB problem [6]. The pure exploration MAB problem is a variant of the MAB problem. Unlike the standard MAB problem, its objective is to identify the optimal action after a number of trials, rather than accumulate as much reward as possible during those trials. The goal of the pure exploration MAB problem can be equivalently formulated as the minimization of *simple regret*, which is defined as the difference of the expected reward between the true optimal action and the action that has been identified as the optimal action. It has been argued that *simple regret bandit algorithms* might be better suited to the task of game tree search, since the ultimate goal of game tree search is to find the best possible action [8]. Therefore, various MCTS variants have been proposed based on simple regret bandit algorithms. The SR+CR scheme [8] is an MCTS algorithm that applies a simple regret bandit algorithm on the root node, and the UCB algorithm on all other nodes. The sequential halving on trees (SHOT) algorithm combines the sequential halving algorithm [15], which is a near optimal simple regret bandit algorithm, with MCTS. The Hybrid MCTS (H-MCTS) algorithm [11] first applies the UCB algorithm on each node, and then switches to the sequential halving algorithm if the number of times a node has been visited has exceeded a predetermined threshold. The CCB-MCTS algorithm [12] uses the improved UCB algorithm to regulate the amount of exploration performed by simple regret bandit algorithms.

However, the paradigm of applying bandit algorithms in all MCTS variants is still essentially the same: viewing every node in the game tree as an independent instance of the MAB problem, and applying the same bandit algorithm and heuristics on every node. Although this approach allows MCTS to be applied in general domains other than game-play, it leaves certain properties of the game tree unexploited.

The adversarial game tree consists of two types of nodes: min nodes and max nodes. Max nodes and min nodes generally represent the decision of different players in the game tree, and

it is conventional knowledge in various games that the decision of which strategy to adopt, should be based on which player he or she is. For example, in the game of Go, a komi of 6.5 is given to the Black player, that is the Black player needs to obtain at least 6.5 points more than the White player to win the game. Therefore, the Black player needs to adopt a more aggressive strategy, while the White player can play more conservatively or defensively. The same can also be observed in the game of Chess, where White is generally considered to have the initiative from the start, and hence needs to play more actively, while Black needs to solve its passivity first. Therefore, max nodes and min nodes are intrinsically different from this high-level point of view, and it would be natural to treat them differently, rather than symmetrically.

Some methods have been proposed to reflect the min-max property of game trees in MCTS, but still essentially treat max nodes and min nodes symmetrically, and apply the same heuristic on every node [9]. The SR+CR scheme differs only the root node from other nodes, rather than between max nodes and min nodes [8].

In this paper, we propose that max nodes and min nodes should be treated differently, and one should apply different bandit algorithms for each node type in MCTS. We will develop the *Asymmetric-MCTS* algorithm, which applies the UCB$_{\sqrt{}}$ algorithm on max nodes and the UCB algorithm on min nodes. We will demonstrate its performance on the game of $9 \times 9$ Go, $9 \times 9$ NoGo, and Othello.

## II. PRELIMINARIES

A key ingredient in the success of MCTS is the application of bandit algorithms. *Bandit algorithms* are algorithms that solve the MAB problem [1].

In the MAB problem, an agent faces $K$ slot machines, or "one-armed bandits", and the agent can choose to pull one of the slot machines at each play. The chosen slot machine will then produce a reward $r \in [0, 1]$. The distribution of the reward of each slot machine is unknown to the agent.

There are two possible objectives in the MAB problem, and different types of bandit algorithms are required for achieving each objective.

### A. Cumulative Regret Minimization

The goal of the conventional MAB problem is to accumulate as much reward as possible over a total of $T$ plays. The objective can be equivalently formulated as the minimization of the *cumulative regret*, which is defined as

$$CR_T = \sum_{t=1}^{T}(r^* - r_{I_t}),$$

where $r^*$ is the expected reward of the optimal arm, and $r_{I_t}$ is the reward that the agent received by pulling arm $I_t$ on play $t$. A bandit algorithm is considered optimal if it can restrict the increase of cumulative regret to $O(\log T)$ [1].

The UCB algorithm [5], which is applied in the UCT algorithm [4], is an optimal bandit algorithm which restricts the growth of cumulative regret to $O(\frac{K \log T}{\Delta})$, where $\Delta$ is the difference of expected reward between a suboptimal arm and the optimal arm.

**Algorithm 1** The $UCB$ algorithm [5]

---

**Initialization**: Play each machine once.
**for** $t = 1, 2, 3, \cdots$ **do**
    play arm $a_i = \arg \max_{i \in K} w_i + c\sqrt{\frac{\log t}{t_i}}$,
where $w_i$ is the current average reward, $t_i$ is the number of times arm $a_i$ has been sampled.
**end for**

---

The UCB algorithm, shown in Algorithm 1, maintains an estimated confidence bound of the expected reward of each arm, and the algorithm simply chooses the arm that has the highest upper bound to pull at each play. The UCB algorithm estimates the confidence bound of arm $a_i$ at play $t$ as

$$UCB_i = w_i \pm c_r\sqrt{\frac{\log t}{t_i}},$$

where $w_i$ is the average reward received from $a_i$ so far, and $t_i$ is the number of times $a_i$ has been played up to play $t$, and $c_r$ is a constant. It can be observed that the confidence bound consists of the *exploitation term* $w_i$, and *the exploration term* $c\sqrt{\frac{\log t}{t_i}}$. The width of the confidence bound is determined by the exploration term, and it gradually decreases as the number of times arm $a_i$ is played increases, e.g. as $t_i$ increases, the bound becomes tighter.

### B. Simple Regret Minimization

The objective of the *pure exploration* MAB problem is to identify the arm that has the highest expected reward after a given total amount of $T$ plays [6]. This task can be formally stated as minimizing the *simple regret*, which is defined as

$$SR_T = r^* - r_T,$$

where $r^*$ is the expected reward of the optimal arm, and $r_T$ is the mean reward of the arm that is identified by the agent to be optimal after the $T$ plays.

Since the goal is to identify which arm is the optimal arm, it is more critical to gather as much information as possible about each arm, and therefore the amount of accumulated reward during these $T$ plays is irrelevant. It has been shown that the minimization of cumulative regret $CR_T$ and the minimization of simple regret $SR_T$ are two contradicting objectives, i.e., as $CR_T$ decreases, $SR_T$ will increase at the same time, and vice versa [6]. Therefore, in order to solve the pure exploration MAB problem, a different type of bandit algorithms is needed.

**Algorithm 2** The $UCB_{\sqrt{}}$ algorithm [8]

---

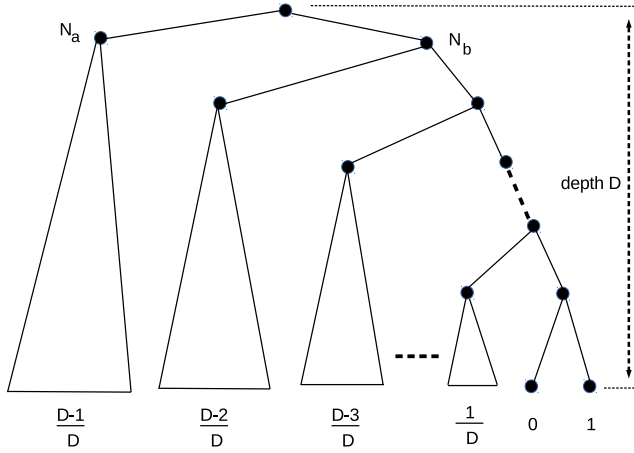**Initialization**: Play each machine once.
**for** $t = 1, 2, 3, \cdots$ **do**
    play arm $a_i = \arg \max_{i \in K} w_i + c\sqrt{\frac{\sqrt{t}}{t_i}}$,
where $w_i$ is the current average reward, $t_i$ is the number of times arm $a_i$ has been sampled.
**end for**

---

The UCB$_{\sqrt{}}$ algorithm, shown in Algorithm 2, is a bandit algorithm that restricts the growth of simple regret to

Fig. 1: An example tree for which the UCT algorithm has very poor performance [7].



Fig. 2: Asymmetric-MCTS algorithm. The gray nodes are max nodes, which the $UCB_{\sqrt{}}$ bandit algorithm are applied, and the white nodes are min nodes, which the UCB bandit algorithm applied.



$O((\Delta \exp(-\sqrt{T}))^K)$ [8].

The algorithmic aspect of the $UCB_{\sqrt{}}$ algorithm is basically the same as the UCB algorithm, as it also maintains an estimated confidence bound for each arm, and chooses the arm with the highest upper bound at each play. The $UCB_{\sqrt{}}$ algorithm defines the confidence bound for arm $a_i$ as

$$UCB_i = w_i \pm c_s \sqrt{\frac{\sqrt{t}}{t_i}},$$

where $w_i$ is the average reward received so far from arm $a_i$, $t_i$ is the number of times that $a_i$ has been played up to play $t$, and $c_s$ is a constant. As with the UCB algorithm, the confidence bound of the $UCB_{\sqrt{}}$ algorithm also consists of the *exploitation term* $w_i$ and the *exploration term* $c\sqrt{\frac{\sqrt{t}}{t_i}}$.

The difference between the UCB and the $UCB_{\sqrt{}}$ algorithm lies in the definition of the exploration term. The exploration term for the $UCB_{\sqrt{}}$ algorithm decreases more slowly than that of the UCB algorithm, and hence tends to sample more arms over time than focusing on the arm that currently seems to be optimal.

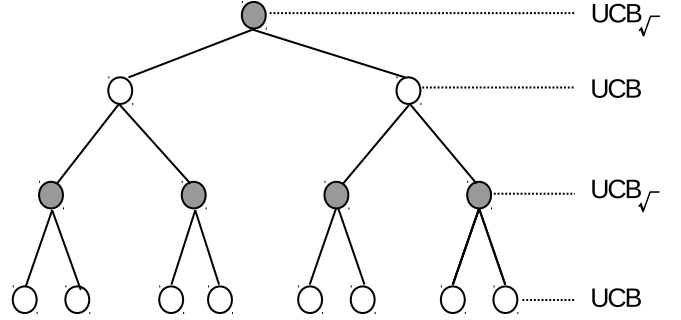## III.  ASYMMETRIC MONTE-CARLO TREE SEARCH

MCTS consists of four major steps: *selection*, *expansion*, *simulation*, and *backpropagation*. Bandit algorithms are mainly applied in the *selection* phase by viewing each node as an independent instance of the MAB problem, where each child node is a single candidate arm. Currently, the most popular variant of MCTS is the UCT algorithm, in which the UCB algorithm is the applied bandit algorithm.

Although this general MCTS paradigm allows it to be applied in a wide range of domains, it leaves a number of properties of the game tree unexploited.

### A.  Concerns on Value Estimation of Different Node Types

The role of the bandit algorithm on every node of MCTS is to estimate the value of the node and perform selection according to the estimated value. As the search progresses, the estimation value of the nodes also converges. Although the general goal is to obtain a good estimation as fast as possible, it can be observed that different node types in the game tree have different requirements to their estimated values:

- **Max node**: since the max nodes represent the view point of the current decision maker, we need to be more certain about the estimated value of each possible decision. Estimations should also be more cautious, and not overly optimistic.

- **Min node**: since the min nodes represent the reaction of the opponent, it is not necessary to determine the best possible reaction of the opponent. Just a *good enough* reaction that is sufficient to refute a decision made by the decision maker will do.

Due to the selection and expansion performed in MCTS, the reward of the MAB problem at each node is non-stationary [7]. For example, consider the binary tree used for constructing a lower bound for the UCT algorithm [7], shown in Fig. 1. The binary tree has the depth of $D$, and the rightmost path, which is from the root node to the rightmost leaf node, is the optimal path. For a node $N$ at depth $d < D$ on the optimal path, if the left action is chosen, then a reward of $\frac{D-d}{D}$ is received. In other words, all the leaf nodes of the subtree rooted at $N$ have the value of $\frac{D-d}{D}$. If the right action is chosen, the agent can proceed to expand further down the optimal path. At depth $D-1$ of the optimal path, the left action will give the reward 0, and the right action will give the reward 1. Therefore, MCTS will most likely spend the majority of its time expanding the subtrees of the left action along the optimal path, as it seems to be better. Consider the MAB problem at the root node, which has two arms node $N_a$ and node $N_b$. Since the leaf nodes of the subtree rooted at $N_a$ all have the value of $\frac{D-1}{D}$, the reward produced by $N_a$ will most likely be fixed around $\frac{D-1}{D}$. However, as the search gradually expand down the optimal path, the reward produced by $N_b$ will most likely be along the sequence

$$\{\tfrac{D-2}{D}, \cdots, \tfrac{D-2}{D}, \tfrac{D-3}{D}, \cdots, \tfrac{D-3}{D}, \cdots, \tfrac{1}{D}, \cdots, \tfrac{1}{D}, 1\},$$

**Algorithm 3** Asymmetric-MCTS Algorithm

**function** ASYMMETRIC-MCTS(Node $N$)
    $best_{ucb} \leftarrow -\infty$
    **for** all child nodes $n_i$ of $N$ **do**
        **if** $n_i.t = 0$ **then**
            $n_i.ucb \leftarrow \infty$
        **else**
            **if** $N.type$ is $MAX$ **then**
$$n_i.ucb \leftarrow n.w + c_s \cdot \sqrt{\frac{\sqrt{N.t}}{n_i.t}}$$
            **else**
$$n_i.ucb \leftarrow n.w + c_r \cdot \sqrt{\frac{\log N.t}{n_i.t}}$$
            **end if**
        **end if**
        **if** $best_{ucb} \leq n_i.ucb$ **then**
            $best_{ucb} \leftarrow n_i.ucb$
            $n_{best} \leftarrow n_i$
        **end if**
    **end for**

    **if** $n_{best}.t = 0$ **then**
        $result \leftarrow$ RANDOMSIMULATION($n_{best}$)
    **else**
        **if** $n_{best}$ is not expanded **then** EXPAND($n_{best}$)
        $result \leftarrow$ ASYMMETRIC-MCTS($n_{best}$)
    **end if**

    $n_{best}.w \leftarrow (n_{best}.w \times n_{best}.t + result)/(n_{best}.t + 1)$
    $n_{best}.t \leftarrow n_{best}.t + 1$
    $N.t \leftarrow N.t + 1$
    **return** $result$
**end function**

Fig. 3: Optimal percentage of biased reward MAB problem



(a) ascending reward



(b) descending reward

instead of being more evened out. Therefore, although the distribution of the reward of the MAB problem on each node is fixed and determined by the values of the leaf node, due to the selection and expansion performed in MCTS, the reward of the MAB problem on each node is biased, and hence affect the estimation made by the bandit algorithms.
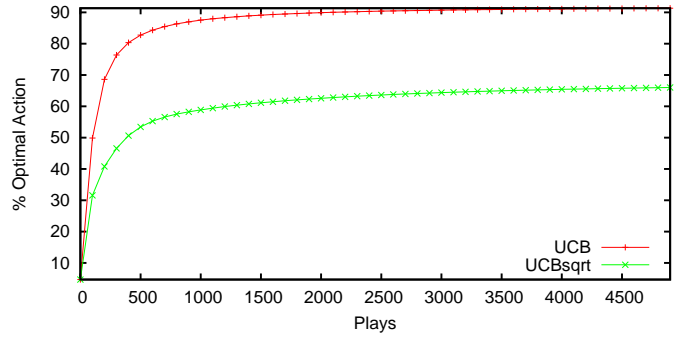
Consider the case where a sequence of rewards $(r_1, r_2, r_3, \cdots, r_n)$ are drawn from distribution $\mathcal{N}$, but due to some sampling bias, the sequence is ordered in a non-decreasing order, that is $r_i \leq r_j$ if $i < j$. Therefore, the estimated mean reward will be higher than the true mean reward in the early period of the sequence, and hence causing the agent to be too optimistic. Similarly, if the sequence is in a non-increasing order, that is $r_i \geq r_j$ if $i > j$, then the agent tends to be underestimate the mean in the early stages.

Therefore, one should choose a bandit algorithm that is most likely to resist over optimistic estimations caused by biased reward to deploy on max nodes, and a bandit algorithm that can adapt itself rapidly to provide a "good enough" estimation on min nodes.
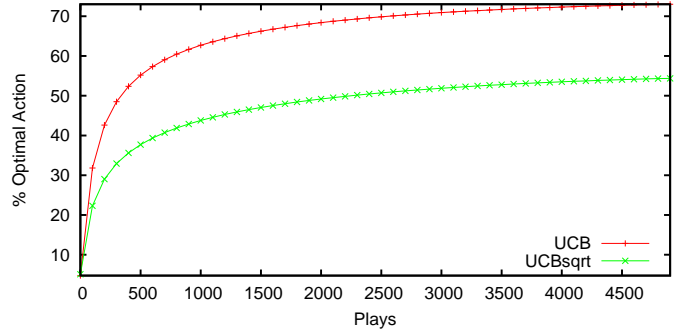
### B. Different Bandit Algorithms for Different Node Types

As simple regret and cumulative regret bandit algorithms have different properties, they can be deployed to different node types accordingly to fulfill the requirements on the estimation value of each node type:
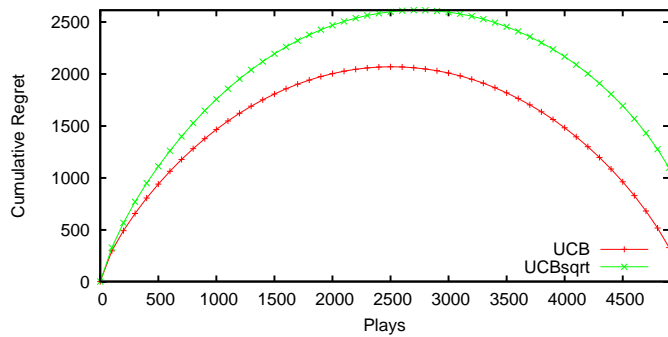
- **Max node**: *simple regret* bandit algorithms, which determine the optimal arm, have a higher level of confidence in its estimation value of each arm. Moreover, in order to provide a better estimation of the value of each arm, simple regret bandit algorithms tend to perform more exploration, and spread its sampling more evenly across the candidates, which effectively make it less likely to be too optimistic.

- **Min node**: *cumulative regret* bandit algorithms, which try to accumulate as much reward as possible, tend to focus on the current optimal arm, and adapt rapidly if the current optimal arm changes. Therefore, cumulative regret bandit algorithms seem to fit the requirement of finding a *good enough* reaction to refute a candidate decision.

The *Asymmetric-MCTS* algorithm, which is shown in Algorithm 3, still retains the four steps in conventional MCTS, namely *selection*, *expansion*, *simulation*, and *backpropagation*. The main characteristic of the Asymmetric-MCTS is that it applies the UCB$_{\sqrt{\cdot}}$ algorithm, which is a simple regret bandit algorithm, on max nodes, and the UCB algorithm, which is a cumulative regret bandit algorithm, on min nodes, as shown in Figure 2.
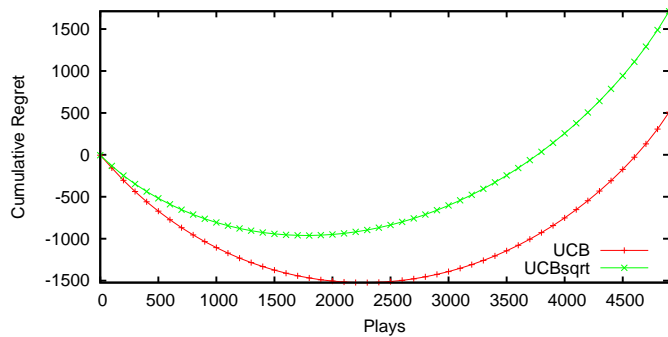
### IV. EXPERIMENTAL RESULTS

In this section, we will first demonstrate the effect of biased reward on the UCB and UCB$_{\sqrt{\cdot}}$ algorithm. We will then

Fig. 4: Cumulative regret of biased reward MAB problem.



(a) ascending reward



(b) descending reward

Fig. 5: Simple regret of biased reward MAB problem.



(a) ascending reward



(b) descending reward

proceed to demonstrate the performance of the Asymmetric-MCTS algorithm on the game of $9 \times 9$ Go, $9 \times 9$ Nogo, and Othello. The baseline for all experiments is the plain UCT algorithm. For a direct comparison of the effect of the bandit algorithms, all MCTS algorithms used pure random simulations, and no performance enhancement heuristics were applied. Every experimental result is the average of 2300 games, and each algorithm took turns in playing with Black and White.
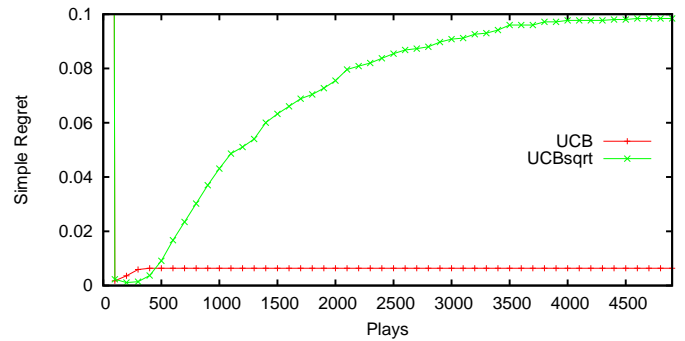
### A. Effect of Biased Reward in the MAB problem

We will first demonstrate how bias in the reward affects the performance of the UCB and UCB$_{\sqrt{\cdot}}$ algorithm. In order to enhance the effect of the biased reward, we will examine two extreme cases: the rewards are biased to be produced in ascending order, and descending order.
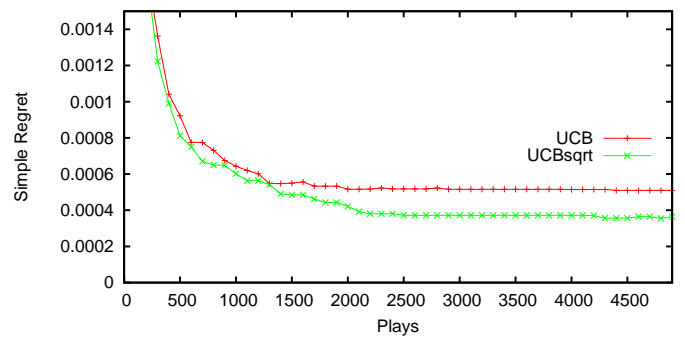
The MAB problem testbed mainly follows the settings specified in Sutton et al. [2]. The results are the average of 2000 randomly generated $K$-armed bandit problems, with $K = 20$. A total of 5000 plays were given to each problem. The rewards of each bandit are first generated from a normal (Gaussian) distribution with the mean $w_i$, $i \in K$, and variance of 1. The mean $w_i$ of the bandits were randomly selected from a normal distribution with mean 0 and variance 1. To simulate biased rewards, the rewards are then sorted in ascending order and descending order.

It can be observed from Figure 3a and Figure 3b that

regardless of the order in which the rewards are biased, the UCB algorithm has a higher percentage of pulling the optimal arm than the UCB$_{\sqrt{\cdot}}$ algorithm, and hence suggesting the UCB$_{\sqrt{\cdot}}$ algorithm tends to distribute its plays more evenly across the candidates. As a result, The UCB algorithm also has lower cumulative regret than the UCB$_{\sqrt{\cdot}}$ algorithm in both cases, as shown in Figure 4a and Figure 4b.

However, the UCB$_{\sqrt{\cdot}}$ algorithm has a lower simple regret than the UCB algorithm when the rewards are produced in descending order, as shown in Figure 5b. As the UCB$_{\sqrt{\cdot}}$ algorithm performs more exploration, it is able to obtain a better estimation of the mean reward of each candidate, and thus can make more informed recommendations, achieving lower simple regret. On the other hand, the extensive explorations performed by the UCB$_{\sqrt{\cdot}}$ algorithm cause its estimations to be too conservative and pessimistic, and hence lower the quality of the recommendations, as shown in Figure 5a.

Therefore, it can be observed that the UCB$_{\sqrt{\cdot}}$ algorithm is more conservative in its estimations, and more resistant to situations where it is more likely to make overly optimistic estimations. One the other hand, the UCB algorithm follows closely the change in the reward with high efficiency.

### B. Performance of the Asymmetric-MCTS on $9 \times 9$ Go

We will first investigate the performance of the Asymmetric-MCTS on the game of Go played on the $9 \times 9$ board, with the komi of 6.5.

TABLE I: Win Rate of SR+CR scheme [8] against plain UCT algorithm in $9 \times 9$ Go.

| $c_s$ | SR+CR Scheme |
|-----|-----|
| 0.1 | 50.00% $\pm$ 2.04% |
| 0.2 | 51.29% $\pm$ 2.04% |
| 0.3 | 51.80% $\pm$ 2.04% |
| 0.4 | 53.91% $\pm$ 2.04% |
| 0.5 | 53.50% $\pm$ 2.04% |
| 0.6 | 51.19% $\pm$ 2.04% |
| 0.7 | 52.23% $\pm$ 2.04% |
| 0.8 | 51.80% $\pm$ 2.04% |
| 0.9 | **54.83% $\pm$ 2.04%** |

TABLE II: The win rate of the Asymmetric-MCTS with $c_r = 0.5, c_s = 0.4$ against plain UCT algorithm with various constant $c$ settings on $9 \times 9$ Go. The optimal setting for plain UCT algorithm is $c = 0.3$, which can only restrict the win rate of Asymmetric-MCTS to 57.70%.

| $c$ | Win Rate |
|-----|-----|
| 0.1 | 58.96% $\pm$ 2.01% |
| 0.2 | 59.52% $\pm$ 2.01% |
| 0.3 | **57.70% $\pm$ 2.02%** |
| 0.4 | 58.61% $\pm$ 2.01% |
| 0.5 | 58.30% $\pm$ 2.01% |
| 0.6 | 60.74% $\pm$ 2.00% |
| 0.7 | 62.30% $\pm$ 1.98% |
| 0.8 | 61.91% $\pm$ 1.99% |
| 0.9 | 61.61% $\pm$ 1.99% |

*1) Performance of SR+CR scheme:* For comparison, we demonstrate the performance of SR+CR scheme on the game of $9 \times 9$ Go. The SR+CR scheme applies the $UCB_{\sqrt{\cdot}}$ bandit algorithm only on the root node, and the UCB bandit algorithm on all other nodes [8]. Table I shows the win rate of various settings for the constant $c_s$ in the $UCB_{\sqrt{\cdot}}$ algorithm in the SR+CR scheme algorithms. The best constant setting for the UCB algorithm is $c_r = 0.4$ in the SR+CR scheme and the plain UCT algorithm is $c = 0.4$. A total of 5000 playouts are given to both algorithms for each move.

It can be observed that the SR+CR scheme achieves around 54% with its best constant setting, which is slightly better than the plain UCT algorithm.

*2) Tuning the C constants:* We now proceed to find the best settings for the constant $c_r$ in the UCB algorithm applied on min nodes, and the constant $c_s$ in the the $UCB_{\sqrt{\cdot}}$ algorithm applied on the max nodes, in the Asymmetric-MCTS algorithm. We have found the optimal setting as $c_r = 0.5$ and $c_s = 0.4$, and Table II shows the win rate of the Asymmetric-MCTS against various constant settings for the plain UCT algorithm. A total of 5000 playouts are given to both algorithms for each move.

It can be observed that even against the best setting $c = 0.3$ of the plain UCT algorithm, the Asymmetric-MCTS still manages to achieve a win rate of around 57.70%. In comparison to the performance of SR+CR scheme, this result suggests that applying the $UCB_{\sqrt{\cdot}}$ algorithm on the max nodes throughout the game tree, instead of only on the root node, can make a difference.

TABLE III: Scalability of the Asymmetric-MCTS on $9 \times 9$ Go.

| Playouts | Win Rate |
|-----|-----|
| 1000 | 65.22% $\pm$ 1.95% |
| 3000 | 60.00% $\pm$ 2.00% |
| 5000 | 57.70% $\pm$ 2.02% |
| 7000 | 59.39% $\pm$ 2.01% |
| 9000 | 59.57% $\pm$ 2.01% |
| 11000 | 62.61% $\pm$ 1.98% |

TABLE IV: Win Rate of $UCB_{\sqrt{\cdot}}$ MCTS and SR+CR scheme [8] against plain UCT algorithm in $9 \times 9$ NoGo. The SR+CR scheme has a best win rate of 67.21% when $c_s = 0.9$ and $c_r = 0.3$.

| $c_s$ | SR+CR Scheme |
|-----|-----|
| 0.1 | 50.23% $\pm$ 2.04% |
| 0.2 | 51.80% $\pm$ 2.04% |
| 0.3 | 52.70% $\pm$ 2.04% |
| 0.4 | 59.60% $\pm$ 2.01% |
| 0.5 | 64.35% $\pm$ 1.96% |
| 0.6 | 65.63% $\pm$ 1.94% |
| 0.7 | 65.28% $\pm$ 1.95% |
| 0.8 | 66.53% $\pm$ 1.93% |
| 0.9 | **67.21% $\pm$ 1.92%** |

*3) Scalability of Asymmetric-MCTS:* We now investigate the scalability of the Asymmetric-MCTS as the total number of playouts increases. The result is shown in Table III. The settings for Asymmetric-MCTS is $c_r = 0.5$ and $c_s = 0.4$, and that for the plain UCT algorithm is set to $c = 0.3$.

We can observe that the Asymmetric-MCTS achieves a very good win rate of around 65% over the plain UCT algorithm when 1000 playouts are given, and keeps the win rate to around 60% as more playouts are given to both algorithms. The results suggest that the Asymmetric-MCTS algorithm has very steady performance on the game of $9 \times 9$ Go.

### C. Performance of the Asymmetric-MCTS on $9 \times 9$ NoGo

We now demonstrate the performance of the Asymmetric-MCTS on the game of Nogo. Nogo is a misere variation of the game of Go, in which the first player who has no legal moves other than capturing the stones of the opponent loses.

*1) Performance of SR+CR scheme:* As in $9 \times 9$ Go, we first demonstrate the performance of the SR+CR scheme on the game of $9 \times 9$ NoGo for comparison. Table IV shows the win rate of various settings for the constant $c_s$ in the SR+CR scheme. The constant setting for the plain UCT algorithm is $c = 0.3$. A total of 5000 playouts are given to both algorithms for each move.

We can observe that SR+CR scheme did extremely well against the plain UCT algorithm, achieving a near 68% win rate against the plain UCT algorithm.

*2) Tuning the C constants:* We now proceed to find the best settings for the constants $c_r$ and $c_s$ the $UCB_{\sqrt{\cdot}}$ in the Asymmetric-MCTS algorithm. The optimal setting for the Asymmetric-MCTS algorithm is $c_r = 0.5$ and $c_s = 0.4$. Table V shows the win rate of the Asymmetric-MCTS against various

TABLE V: The win rate of the Asymmetric-MCTS with $c_r = 0.5, c_s = 0.4$ against plain UCT algorithm with various constant $c$ settings on $9 \times 9$ NoGo. The optimal setting for plain UCT algorithm is $c = 0.4$, which can only restrict the win rate of Asymmetric-MCTS to 62.43%.

| $c$ | Win Rate |
|-----|----------|
| 0.1 | 64.74% $\pm$ 1.95% |
| 0.2 | 66.48% $\pm$ 1.93% |
| 0.3 | 66.17% $\pm$ 1.93% |
| 0.4 | **62.43% $\pm$ 1.98%** |
| 0.5 | 65.65% $\pm$ 1.94% |
| 0.6 | 67.00% $\pm$ 1.92% |
| 0.7 | 67.65% $\pm$ 1.91% |
| 0.8 | 67.83% $\pm$ 1.91% |
| 0.9 | 69.83% $\pm$ 1.88% |

TABLE VI: Scalability of the Asymmetric-MCTS on $9 \times 9$ NoGo.

| Playouts | Win Rate |
|----------|----------|
| 1000 | 57.57% $\pm$ 2.02% |
| 3000 | 59.48% $\pm$ 2.01% |
| 5000 | 62.43% $\pm$ 1.98% |
| 7000 | 65.65% $\pm$ 1.94% |
| 9000 | 64.96% $\pm$ 1.95% |
| 11000 | 65.96% $\pm$ 1.94% |

constant settings for the plain UCT algorithm. A total of 5000 playouts are given to both algorithms for each move.

It can be observed that the Asymmetric-MCTS algorithm achieves at least a win rate of 62.43% against the plain UCT algorithm. This result suggests that differentiating max nodes and min nodes also produces very good performance, although the SR+CR scheme might be a better choice on the game of $9 \times 9$ NoGo.

*3) Scalability of Asymmetric-MCTS:* We now investigate the scalability of the Asymmetric-MCTS as the total number of playouts increases when applied on $9 \times 9$ Nogo. The results are shown in Table VI. The settings for Asymmetric-MCTS is $c_r = 0.5$ and $c_s = 0.4$, and the constant for the plain UCT algorithm is set to $c = 0.4$.

It can be observed that the Asymmetric-MCTS algorithm dominates the plain UCT algorithm from a total of 1000 playouts to 11000 playouts, and the win rate gradually increases to near 66% when 11000 playouts are given to both algorithms. This result suggests that the effect of differentiating max nodes and min nodes will gradually increase with the number of total playouts.

### D. Performance of the Asymmetric-MCTS on Othello

Finally, we proceed to demonstrate the performance of the Asymmetric-MCTS algorithm on the game of Othello.

*1) Performance of SR+CR scheme:* We will first investigate the performance of the SR+CR scheme on Othello for comparison. Table VII shows the win rate of various settings for the constant $c_s$ in the UCB$_{\sqrt{}}$ algorithm of the SR+CR scheme. The constant setting for the UCB algorithm in the SR+CR scheme is $c_r = 0.6$ and the plain UCT algorithm is

TABLE VII: Win Rate of the SR+CR scheme [8] against plain UCT algorithm on Othello. The SR+CR scheme has a best win rate of 53.87%

| $c_s$ | SR+CR Scheme |
|-------|--------------|
| 0.1 | 37.74% $\pm$ 1.98% |
| 0.2 | 51.34% $\pm$ 2.04% |
| 0.3 | 53.26% $\pm$ 2.04% |
| 0.4 | **53.87% $\pm$ 2.04%** |
| 0.5 | 52.35% $\pm$ 2.04% |
| 0.6 | 50.74% $\pm$ 2.04% |
| 0.7 | 50.30% $\pm$ 2.04% |
| 0.8 | 49.43% $\pm$ 2.04% |
| 0.9 | 49.78% $\pm$ 2.04% |

TABLE VIII: The win rate of the Asymmetric-MCTS with $c_r = 0.7, c_s = 0.4$ against plain UCT algorithm with various constant $c$ settings on Othello. The optimal setting for plain UCT algorithm is $c = 0.6$, which the Asymmetric can only achieve win rate of 50.47%.

| $c$ | Win Rate |
|-----|----------|
| 0.1 | 88.86% $\pm$ 1.29% |
| 0.2 | 81.74% $\pm$ 1.58% |
| 0.3 | 70.48% $\pm$ 1.86% |
| 0.4 | 57.61% $\pm$ 2.02% |
| 0.5 | 53.87% $\pm$ 2.04% |
| 0.6 | **50.47% $\pm$ 2.04%** |
| 0.7 | 53.39% $\pm$ 2.04% |
| 0.8 | 52.13% $\pm$ 2.04% |
| 0.9 | 53.22% $\pm$ 2.04% |

$c = 0.6$. A total of 5000 playouts are given to both algorithms for each move.

It can be observed that the SR+CR scheme can produce a best win rate of around 53%, which is slightly better but still around the same level of the plain UCT algorithm.

*2) Tuning the C constants:* We will now proceed to find the best settings for the constants $c_r$ and $c_s$ the UCB$_{\sqrt{}}$ in the Asymmetric-MCTS algorithm. The optimal setting for the Asymmetric-MCTS algorithm is $c_r = 0.7$ and $c_s = 0.4$. Table V shows the win rate of Asymmetric-MCTS against various constant settings for the plain UCT algorithm. A total of 5000 playouts are given to both algorithms for each move.

It can be observed that the Asymmetric-MCTS algorithm can only achieve a win rate of around 50% against the plain UCT algorithm. This result suggests that differentiating max nodes and min nodes is not effective on the game of Othello, and is around the same level of performance as the plain UCT algorithm.

*3) Scalability of Asymmetric-MCTS:* We will now investigate the scalability of the Asymmetric-MCTS as the total number of playouts increases when applied on Othello. The results are shown in Table IX. The settings for Asymmetric-MCTS is $c_r = 0.7$ and $c_s = 0.4$, and the plain UCT algorithm is set to $c = 0.4$.

It can be observed that the performance of the Asymmetric-MCTS algorithm does not change with the increase of the number of playouts. The win rate of Asymmetric-MCTS

TABLE IX: Scalability of the Asymmetric-MCTS on Othello.

| Playouts | Win Rate |
|---|---|
| 1000 | 52.37% $\pm$ 2.04% |
| 3000 | 52.04% $\pm$ 2.04% |
| 5000 | 53.43% $\pm$ 2.04% |
| 7000 | 50.87% $\pm$ 2.04% |
| 9000 | 51.22% $\pm$ 2.04% |
| 11000 | 53.43% $\pm$ 2.04% |

algorithm holds steady around 50%, which is around the same performance level as the plain UCT algorithm.

## V. CONCLUSION

MCTS has made quite an impact on various fields, and the key to its success lies in the application of bandit algorithms, which solve the MAB problem. In most MCTS variants, the same bandit algorithm and heuristics are applied to every node in the game tree by viewing each node as an independent instance of the MAB problem. The current most dominate variant of MCTS is the UCT algorithm, which applies the UCB bandit algorithm on every node. Although this paradigm has the advantage of allowing MCTS to be applied in a wide spectrum of fields, it leaves a number of properties of the game tree unexploited.

In this work, we have proposed that max nodes and min nodes should be treated differently by applying different bandit algorithms according to its intrinsic nature, rather than using the same bandit algorithm throughout the whole tree. We have observed that different node types have different concerns in their estimation value, and the simple regret bandit algorithms seem to fit the requirements of max nodes, and cumulative regret bandit algorithms seem to fulfill the requirement of min nodes.

The Asymmetric-MCTS algorithm, which applies the UCB$_{\sqrt{\cdot}}$ algorithm on max nodes, and the UCB algorithm on min nodes is proposed based on this observation. The experimental results show that the Asymmetric-MCTS algorithm has a really good performance and scalability on the games of $9 \times 9$ Go. The Asymmetric-MCTS also did well on the game of $9 \times 9$ NoGo, but the SR+CR scheme seems to be a better choice. However, the Asymmetric-MCTS performed only on par with the UCT algorithm on the game of Othello.

As the main difference between the Asymmetric-MCTS algorithm and the UCT algorithm lies in the application of the UCB$_{\sqrt{\cdot}}$ algorithm on max nodes, and hence the effectiveness of the Asymmetric-MCTS algorithm seems to depend on whether the UCB algorithm is more likely to be too optimistic in its estimations on max nodes. Therefore, it can be suggested from the experimental results that the UCB algorithm may make too optimistic estimations on max nodes in the game of $9 \times 9$ Go, and on the root node in the game of $9 \times 9$ Nogo. On the other hand, situations where the UCB algorithm is likely to be too optimistic rarely occurs in Othello.

Applying bandit algorithms other than the UCB and the UCB$_{\sqrt{\cdot}}$ algorithm would be a natural direction for further investigation. Apart from bandit algorithms, most performance enhancement methods and heuristics in MCTS, also treats each node in the game tree as equal [13][14][18]. Therefore, it would be interesting to further investigate the possibility of developing enhancement heuristics according to node types as well.

## REFERENCES

[1] T. L. Lai and H. Robbins, "Asymptotically efficient adaptive allocation rules," *Advances in applied mathematics*, vol. 6, issue 1, pp. 4-22, 1985.

[2] R. S. Sutton and A. G. Barto, "Reinforcement learning: an introduction," *MIT Press*, Cambridge, MA, 1998.

[3] C.B. Browne, E. Powley, D. Whitehouse, S.M. Lucas, P.I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A survey of monte-carlo tree search methods," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 4, issue 1, pp.1-43, 2012.

[4] L. Kocsis, and C. Szepesvári, "Bandit based monte-carlo planning," *Proceedings of the 17th European Conference on Machine Learning (ECML'06)*, pp. 282-293, 2006.

[5] P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-time analysis of the multiarmed bandit problem," *Machine Learning*, vol. 47, issue 2, pp. 235 - 256, 2002.

[6] S. Bubeck, R. Munos, and G. Stoltz, "Pure exploration in multi-armed bandits problems", *Proceedings of the 20th International Conference on Algorithmic Learning Theory (ALT 2009)*, pp. 27-37, 2009.

[7] P.-A. Coquelin and R. Munos, " Bandit Algorithms for Tree Search", *Proceedings of the 23rd Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 67-74, 2007.

[8] D. Tolpin, and S.E. Shimony, "MCTS based on simple regret," *Proceedings of the 26th AAAI Conference on Artificial Intelligence (AAAI)*, pp. 570-576, 2012.

[9] M. Lanctot, M. H.M. Winands, T. Pepels, and N. R. Sturtevant, "Monte Carlo Tree Search with Heuristic Evaluations using Implicit Minimax Backups", *Proceedings of the 2014 IEEE Conference on Computational Intelligence and Games (CIG 2014)*, pp. 341-348, 2014.

[10] T. Cazenave, " Sequential halving applied to trees," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 7, issue 1, pp. 102-105, 2015.

[11] T. Pepels, T. Cazenave, M.H.M. Winands, M.H.M., and M. Lanctot, "Minimizing simple and cumulative regret in monte-carlo tree Search," *Proceedings of Computer Games Workshop at the 21st European Conference on Artificial Intelligence*, pp. 1-15, 2014.

[12] Y.-C. Liu and Y. Tsuruoka ,"Regulation of Exploration for Simple Regret Minimization in Monte-Carlo Tree Search", *Proceedings of the 2015 IEEE Conference on Computational Intelligence and Games (CIG 2015)*, pp.35-42, 2015.

[13] D. Silver, G. Tesauro, "Monte-Carlo Simulation Balancing," *Proceedings of the 26th Annual International Conference on Machine Learning (ICML'09)*, pp. 945-952, 2009.

[14] R. Coulom, "Computing "elo ratings" of move patterns in the game of go," *ICGA Journal*, vol. 30, issue 4, pp. 198-208, 2007.

[15] Z. Karnin, T.Koren, and S. Oren, "Almost optimal exploration in multi-armed bandits," *Proceedings of the 30th International Conference on Machine Learning (ICML'13)*, pp. 1238-1246, 2013.

[16] C. Rosin, "Multi-armed bandits with episode context," *Annals of Mathematics and Artificial Intelligence*, vol. 61, issue 3, pp. 203-230, 2011.

[17] Y. Tian and Y. Zhu, "Better Computer Go Player with Neural Network and Long-term Prediction", *Proceedings of International Conference on Learning Representations (ICLR)*, 2016.

[18] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, "Mastering the game of Go with deep neural networks and tree search," *Nature*, 529, pp. 484-489, 2016.