

# Selection Enhancements



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Florian Weber



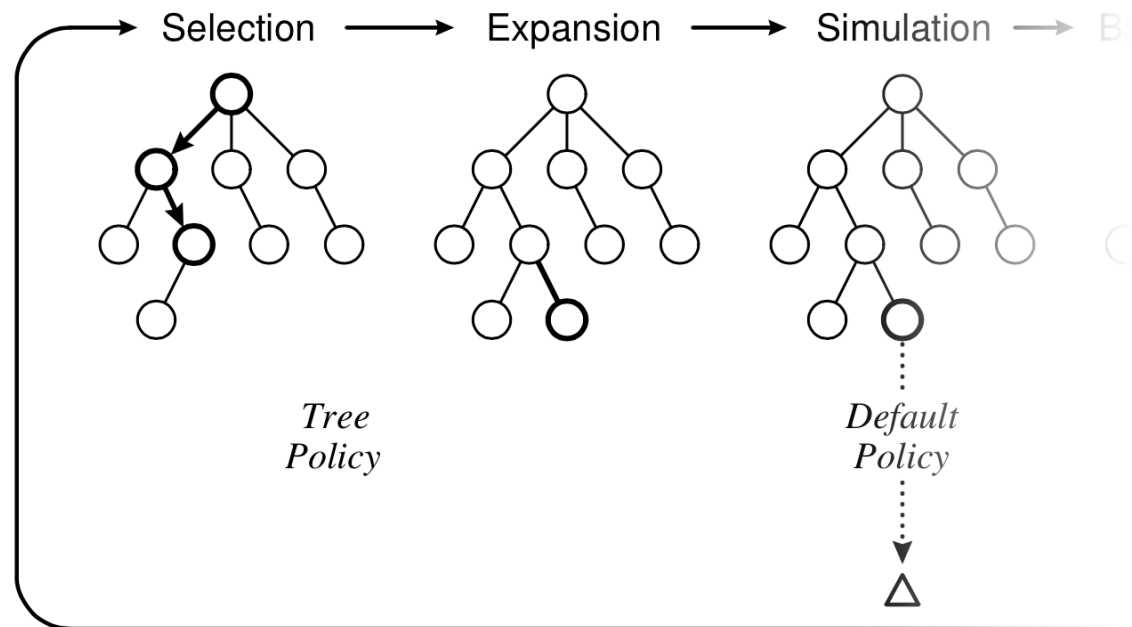
# Selection Enhancements

## Motivation

- **Selection** assigns some numeric score to each action in order to balance exploration with exploitation
  - We heard already of the use of UCB1 for node selection in UCT
  - In the following this approach is called MCTS

- **Enhancement:**

- Bias the search in
  - a specific direction,
  - certain actions,
  - reward estimates



---

# Selection Enhancements

## Agenda

---

- 1 Domain Independent Enhancements
- 2 Domain Dependent Enhancements
- 3 All Moves As First Heuristic (AMAF)
- 4 Rapid Action Value Estimation (RAVE)

# Domain Independent Enhancements

## First Play Urgency (FPU)

---

- **How MCTS does it:** Visit unexplored nodes in random order
  - Trees with large branching factors
    - => exploitation will rarely occur
- **Idea of FPU:**
  - Score unvisited nodes with a fixed value
  - Score visited ones with UCB1
    - => early exploitations are encouraged

## Domain Independent Enhancements

# Transpositions

---

- **How MCTS does it:** Moves are represented as a search tree.
  - We have a transposition when several paths are leading to the same position
  - The statistics are **dependent on the path** to the position
- **Idea of Transpositions:** A Directed Acyclic Graph (DAG) is used to store the information about transpositions
  - Identifying transpositions is easy. Mostly transposition tables are used
  - Statistics are cumulated for the position – **independently** of where it occurs in the tree

# Domain Independent Enhancements

## Move Groups

---

- **How MCTS does it:** A lot of simulation is needed to find groups with a highly correlated expected reward (high branching factor)
- **Idea of Move Groups:** Similar moves (hence correlated actions) are grouped to reduce the branching factor
  - All actions are collected into **Move Groups**
  - UCB1 decides which of them is used to select a move from

---

## Domain Independent Enhancements

# Monte Carlo Paraphrase Generation (MCPG)

---

- **How MCTS does it:** The (average) score expectation is used for selection
- **Idea of MCPG:** The maximum reachable score for each state is used for selection
- **Benefits/where is it used?**
  - It is named after it's usage in generating paraphrases of natural language statements

---

# Selection Enhancements

## Agenda

---

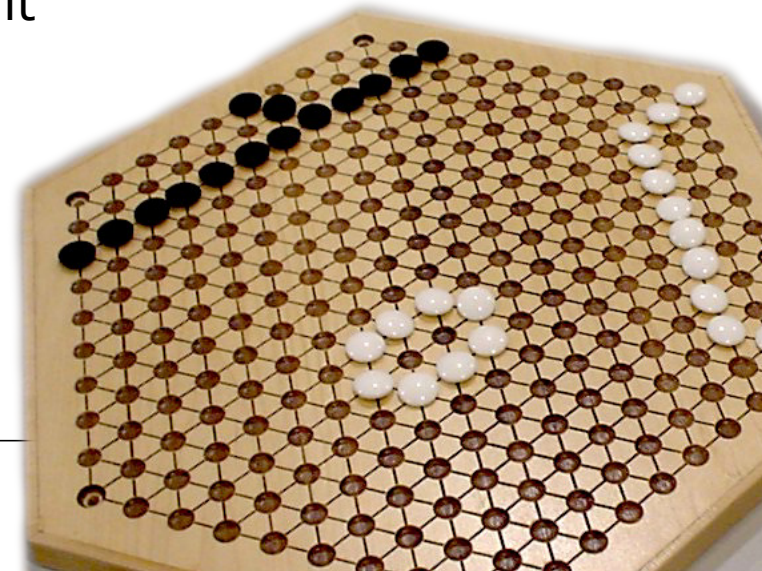
- 1 Domain Independent Enhancements
- 2 Domain Dependent Enhancements
- 3 All Moves As First Heuristic (AMAF)
- 4 Rapid Action Value Estimation (RAVE)



## Domain Dependent Enhancements

# Decisive and Anti-Decisive Moves

- In many games you can find decisive and anti-decisive moves
  - **Decisive move**
    - One that leads immediately to a win
  - **Anti-decisive move**
    - Prevents the opponent from doing a decisive move in the next turn
- **Idea:** Selection and simulation policy:
  - If either player has a decisive move, play it
  - Otherwise play standard policy



# Domain Dependent Enhancements

## Progressive Bias

- **How MCTS does it:** The information of nodes visited only a few times is not reliable
- **Idea of Progressive Bias:** Provide in this case more accurate information via a domain specific heuristic value  $H_i$ 
  - Term added to the selection formula:

$$f(n_1) = \frac{H_i}{n_i + 1} \quad \text{node with index } i, \text{ visited } n_i \text{ times}$$

- The influence of  $f(n_1)$  is high when a few games have been played
- It decreases when more games have been played

# Domain Dependent Enhancements

## Search Seeding

---

- **How MCTS does it:** Every node is initialized with zero wins & visits
- **Idea of Search Seeding:** The statistics at each node is initialized (=seeded) according to some heuristic knowledge. It's like a warm up for the nodes.
  - This extra knowledge can save time since the need for simulation may be reduced
- **Different approaches**
  - Adding virtual wins & visits, prior estimates would remain permanently
  - Some **transient** estimate blended into the regular value  
→ see Progressive Bias

# Domain Dependent Enhancements

## Opening Books

---

### 1) Using Opening Books for MCTS:

- Employing the book till an unlisted position is reached

### 2) Build your own ones:

- Generating an opening book with MCTS
  - Meta Monte-Carlo Search
    - Generate book online/offline

## Domain Dependent Enhancements

# History Heuristic

---

- **Idea:** Using information about moves previously played
- **Grandfather heuristic:** history information is used to initialize action value estimates for new nodes
  - Expected return  $Q_{grand}(s_t, a) = Q_{UCT}(s_{t-2}, a)$  with state  $s$ , action  $a$
- **General game player CadiaPlayer:** history heuristic used for seeding node values

---

# Selection Enhancements

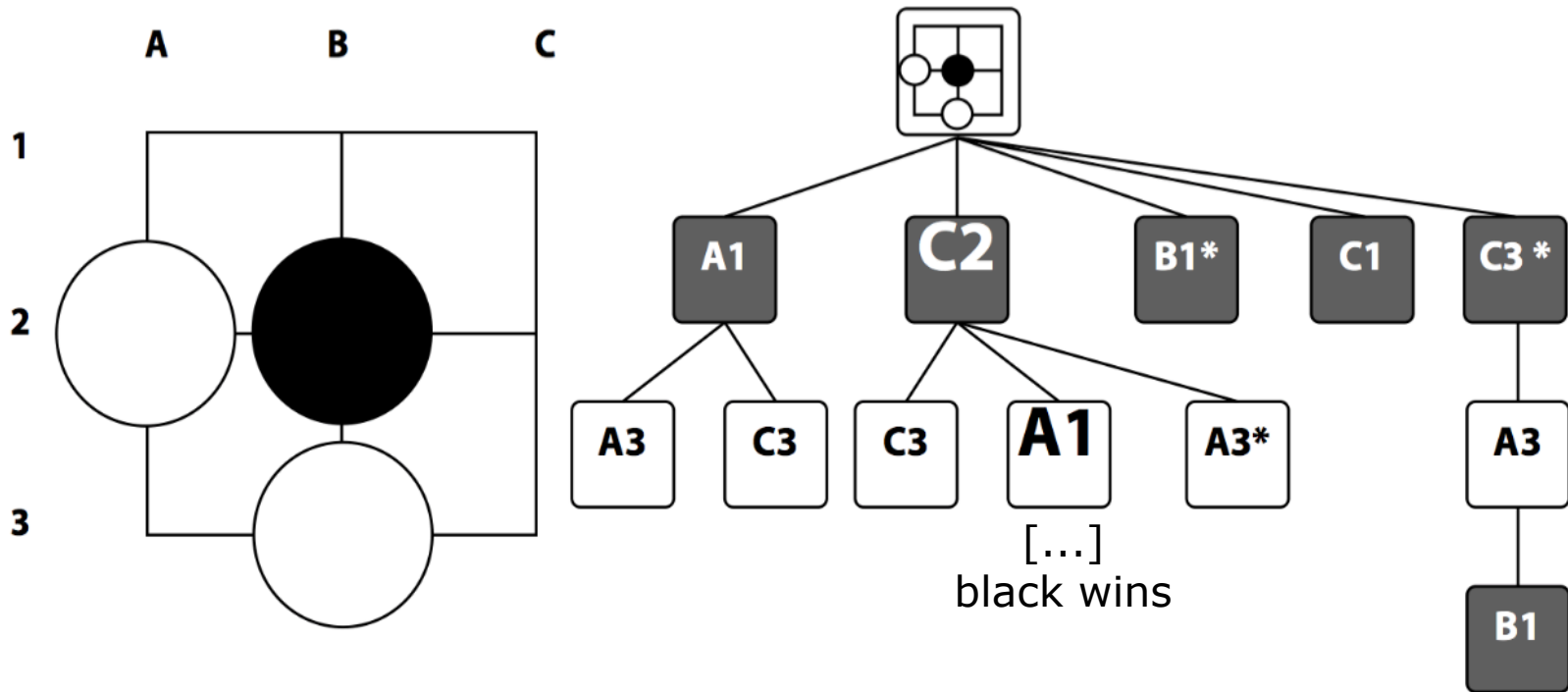
## Agenda

---

- 1 Domain Independent Enhancements
- 2 Domain Dependent Enhancements
- 3 All Moves As First Heuristic (AMAF)
- 4 Rapid Action Value Estimation (RAVE)

# All Moves As First Heuristic (AMAF)

- **How MCTS does it:** after the play-out the estimated rewards for selected nodes are updated (here: C2, A1)

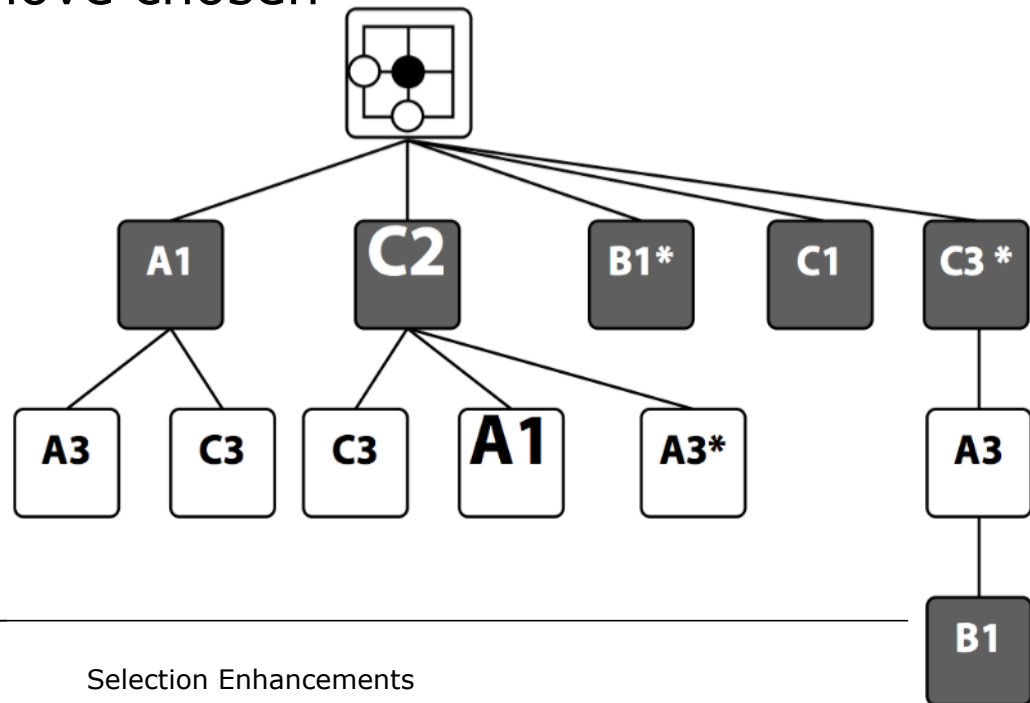


# All Moves As First Heuristic (AMAF)

- **How MCTS does it:** after the play-out the estimated rewards for selected nodes are updated (here: C2, A1)
- **Idea of AMAF:** the reward estimate for an action  $a$  from a state  $s$  is updated whenever  $a$  is encountered during a play-out, even if  $a$  was not the actual move chosen

## Play-out Path:

Black **C2** (UCT)  
White **A1** (UCT)  
Black B1 (\*)  
White A3 (\*)  
Black C3 (\*)  
Win for black





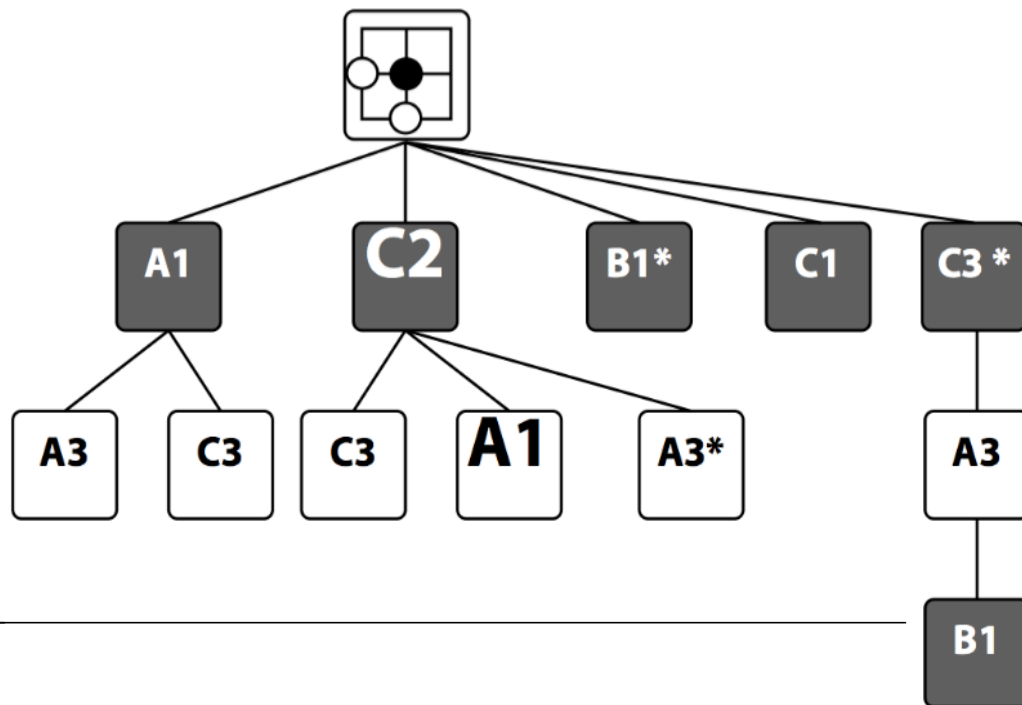
# All Moves As First Heuristic (AMAF)

## Some-First AMAF

- **Idea:** The history of the play-out moves is truncated after the first  $m$  random moves
  - For  $m=0 \Rightarrow$  standard UCT
  - For  $m > \text{play-out length} \Rightarrow$  standard AMAF

### Play-out Path:

Black **C2** (UCT)  
White **A1** (UCT)  
Black B1 (\*)  
White A3 (\*)  
Black C3 (\*)  
Win for black



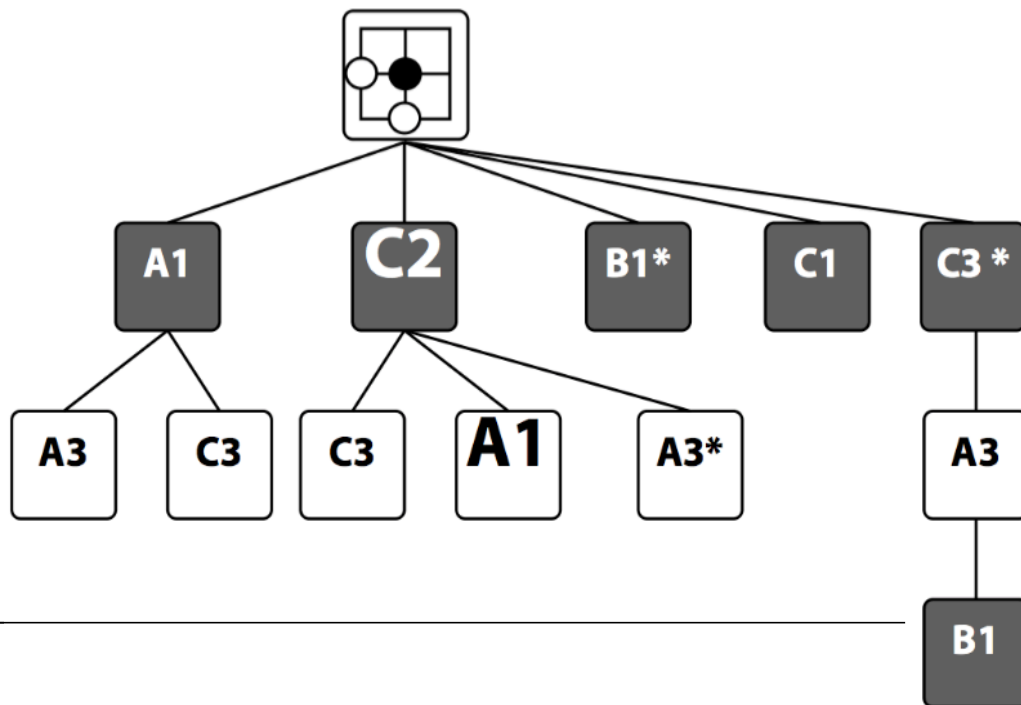
# All Moves As First Heuristic (AMAF)

## Some-First AMAF

- **Idea:** The history of the play-out moves is truncated after the first  $m$  random moves
  - For  $m=2$  ...

### Play-out Path:

Black **C2** (UCT)  
White **A1** (UCT)  
Black B1 (\*)  
White A3 (\*)  
~~Black C3 (\*)~~  
Win for black



# All Moves As First Heuristic (AMAF)

## Cutoff AMAF

- **Idea:** Warm up with AMAF data, then use more accurate UCT
  - Update statistics for the **first  $k$  simulations**
  - **Afterwards** standard **UCT** is used
  - For  $k=0 \Rightarrow$  standard UCT
  - For *sufficiently high*  $k \Rightarrow$  standard AMAF

# All Moves As First Heuristic (AMAF)

## Permutation AMAF

- **Idea:** Also updates nodes that can be reached by permutations of moves that preserve the state reached
  - Stone colors and player alternation are preserved
  - In our example **additional** to the AMAF updates, the nodes **A3** and **B1** are updated

### Play-out Path:

Black **C2** (UCT)

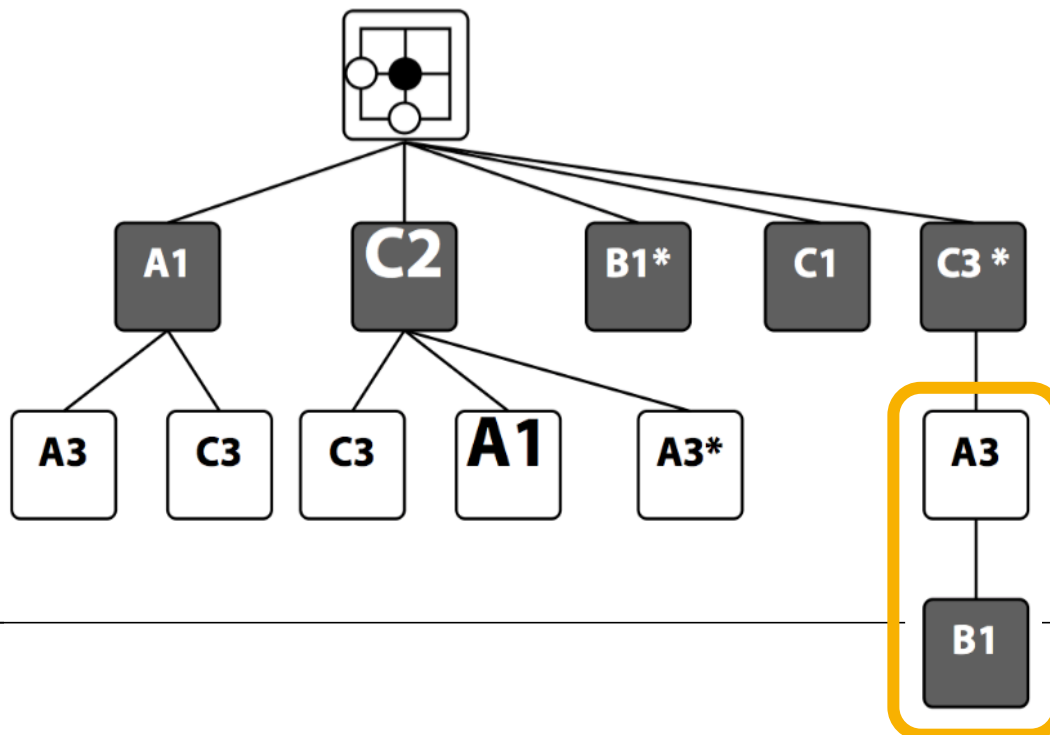
White **A1** (UCT)

Black B1 (\*)

White A3 (\*)

Black C3 (\*)

Win for black



# All Moves As First Heuristic (AMAF)

## $\alpha$ -AMAF

- **Idea:** Blends UCT score  $U$  with AMAF score  $A$ 
  - The score is calculated like this:
$$\alpha A + (1 - \alpha)U$$
  - For  $\alpha = 0 \Rightarrow$  standard UCT
  - For  $\alpha = 1 \Rightarrow$  standard AMAF
- So which  $\alpha$  should we choose?

# All Moves As First Heuristic (AMAF)

## $\alpha$ -AMAF

- Alpha parameter:

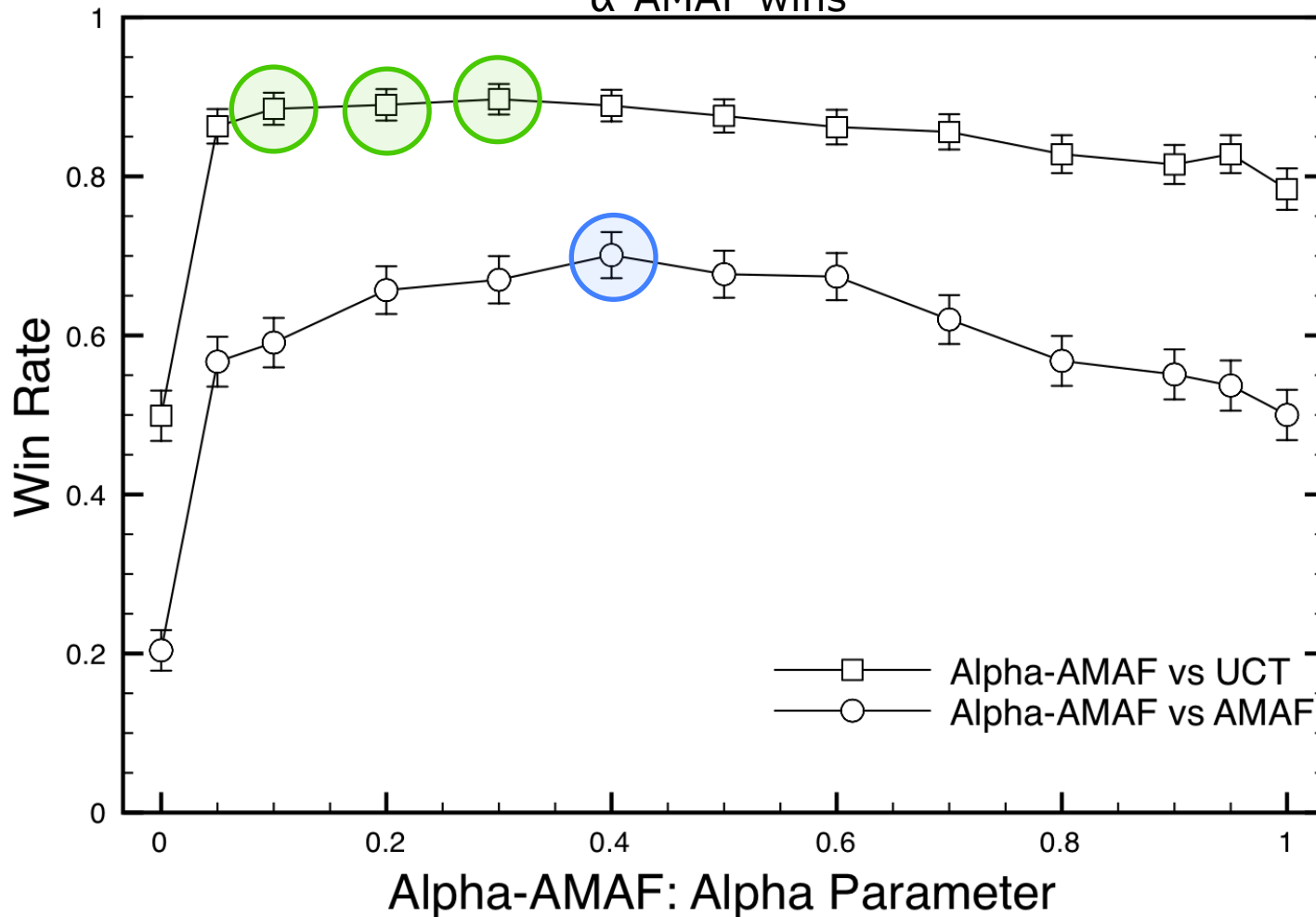
$\alpha$ -AMAF wins

$$\alpha A + (1 - \alpha)U$$

$\alpha$ -AMAF vs. standard UCT:  $\alpha=0.1, 0.2, \text{ or } 0.3$  slightly better

$\alpha$ -AMAF vs. standard AMAF:  $\alpha=0.4$  is best

performance decreases at the extremes 0 and 1



---

# Selection Enhancements

## Agenda

---

- 1 Domain Independent Enhancements
- 2 Domain Dependent Enhancements
- 3 All Moves As First Heuristic (AMAF)
- 4 Rapid Action Value Estimation (RAVE)

# Rapid Action Value Estimation (RAVE)

- **Idea:** similar to  $\alpha$ -AMAF but

- Instead of using a fixed  $\alpha$ , the  $\alpha$  value used at each node decreases with each visit:

$$\alpha = \max \left\{ 0, \frac{V - v(n)}{V} \right\} \quad \begin{array}{l} v(n) \text{ play-outs through node,} \\ \text{fixed integer } V > 0 \end{array}$$

- Exploited areas of the tree will use the accurate statistics (UCT) more than unexploited areas of the tree
- As a reminder:
  - $\alpha$ -AMAF:  $\alpha A + (1 - \alpha)U$
  - For  $\alpha = 0 \Rightarrow$  standard UCT
  - For  $\alpha = 1 \Rightarrow$  standard AMAF



# Rapid Action Value Estimation (RAVE)

## Variants

---

(1) **Killer RAVE**: Only the most important moves are used for the RAVE updates

- **Important moves:**

- Those which have been flagged as strong in other parts of the tree

- **Benefits/where is it used?**

- More beneficial than standard RAVE for game Havannah

(2) **RAVE-max**: More robust variant of RAVE

- Replacing the RAVE value of move  $j$ ,  $Y_j$  with

- $\max(X_j, Y_j)$ ,  $X_j = \text{mean}(j)$

- Able to correct extreme cases of RAVE underestimation

- **Benefits/where is it used?**

- Good for degenerated cases for game Sum of Switches

- Less successful for Go

# Rapid Action Value Estimation (RAVE)

## PoolRAVE

- **Idea:** Modify MCTS as follows:
  - Build a pool of the  $k$  best moves according to RAVE
  - Choose one move  $m$  from the pool
  - Play  $m$  with a probability  $p$ , else the default policy
- **Benefits/where is it used?**
  - The probability  $p$  is used to prevent being too deterministic
  - It helps to keep the diversity of the simulations when biasing Monte-Carlo



	<b>Go</b>	<b>Havannah</b>	<b>Arimaa</b>	<b>Sum of Switches</b>	<b>Cadia Player</b>
<b>UCT</b>	X	X	X	X	X
First Play Urgency	X				
Move Groups	X				
Transpositions	X		X		
<b>MCPG</b>					
(Anti)Decisive Moves	X	X			
Progressive Bias	X				
Opening Books	X				
Search Seeding	X				
History Heuristic	X		X		X
<b>AMAF</b>	X	X			
RAVE	X	X	X	X	X
Killer RAVE		X			
RAVE-max				X	
PoolRAVE	X	X			

---

Thank you very much for you attention!

**QUESTIONS?**