# Monte Carlo Tree Search (MCTS)

Presenter: Tuomas Sandholm

#### MCTS Overview

- Iteratively building partial search tree
- Iteration
  - Most urgent node
    - Tree policy
    - Exploration/exploitation
  - Simulation
    - Add child node
    - Default policy
  - Update weights

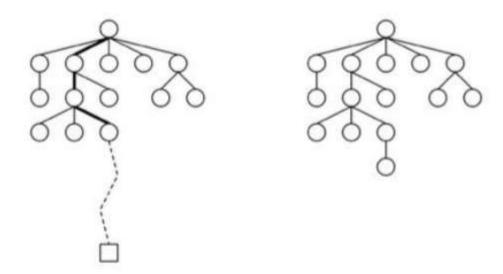


Fig. 1. The basic MCTS process [17].

### Algorithm Overview

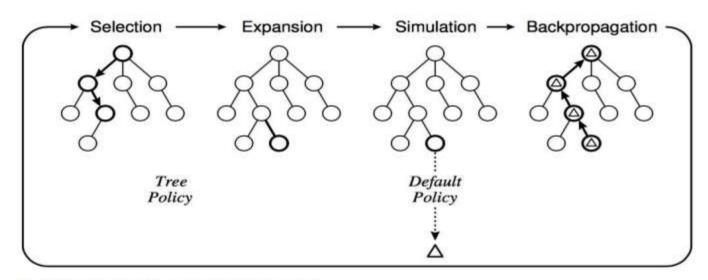


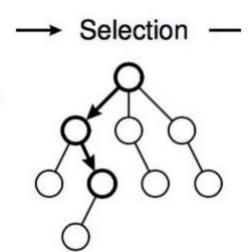
Fig. 2. One iteration of the general MCTS approach.

### **Policies**

- Policies are crucial for how MCTS operates
- Tree policy
  - Used to determine how children are selected
- Default policy
  - Used to determine how simulations are run (ex. randomized)
  - Result of simulation used to update values

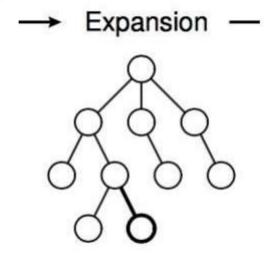
#### Selection

- Start at root node
- Based on Tree Policy select child
- Apply recursively descend through tree
  - Stop when expandable node is reached
  - Expandable -
    - Node that is non-terminal and has unexplored children



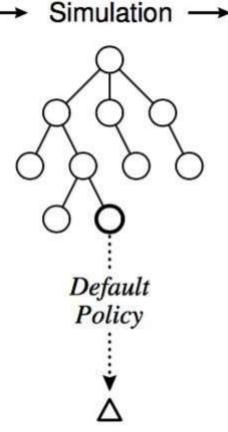
### Expansion

- Add one or more child nodes to tree
  - Depends on what actions are available for the current position
  - Method in which this is done depends on Tree Policy



### Simulation

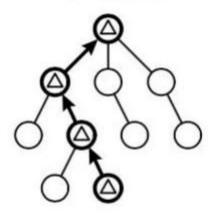
- Runs simulation of path that was selected
- Get position at end of simulation
- Default Policy determines how simulation is run
- Board outcome determines value



### Backpropagation

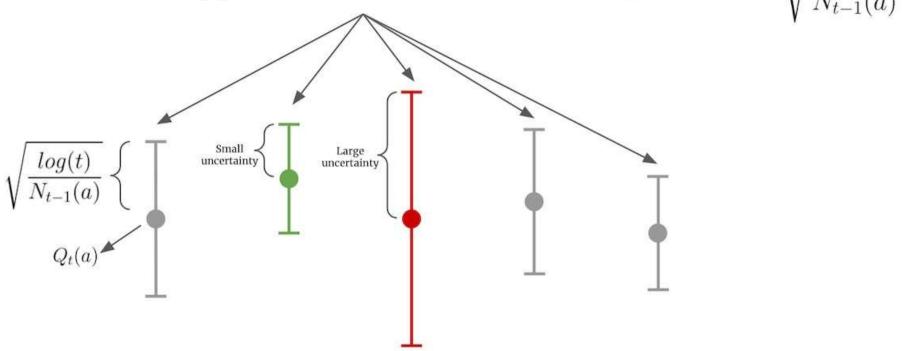
- Moves backward through saved path
- Value of Node
  - representative of benefit of going down that path from parent
- Values are updated dependent on board outcome
  - Based on how the simulated game ends, values are updated

→ Backpropagation -



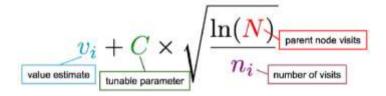
### **UCB** in Bandits

Upper Confidence Bound:  $UCB(a_t) = Q_t(a) + c\sqrt{\frac{log(t)}{N_{t-1}(a)}}$ 



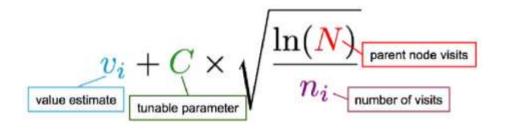
### Upper Confidence bounds applied to Trees (UCT) Algorithm

- Selecting child node: multi-armed bandit problem
  - UCB for child selection
- UCT



- v: value estimate
- C: exploration parameter
- N: number of parent node visits
- n: number of visits

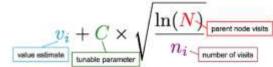
### **UCT** Algorithm

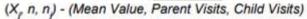


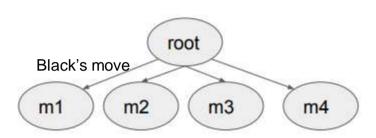
- n = 0 means infinite weight
  - Guarantees we explore each child at least once
- Each child has non-zero probability of selection
- Adjust C to change explore-exploit tradeoff

**Theorem**. MCTS with UCT action selection in the Selection phase finds an optimal policy [Kocsis and Szepesvári. ECML '06]

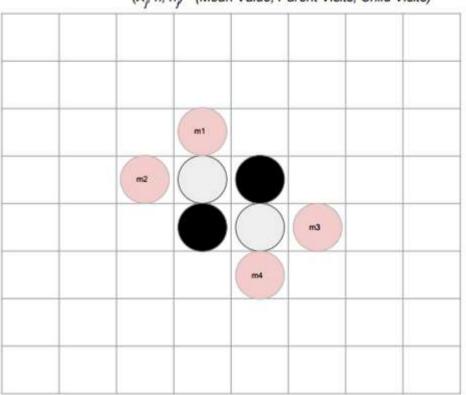
### Example - The Game of Othello



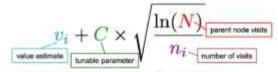




- n<sub>i</sub> initially 0
  - o all weights are initially infinity
- n initially 0
- C<sub>p</sub> some constant > 0
  - o For this example
  - o  $C = (1/2\sqrt{2})$
- X<sub>j</sub> mean reward of selecting this position
  - o [0, 1]
  - Initially N/A



### Example - The Game of Othello cont.

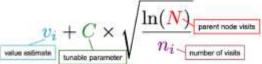


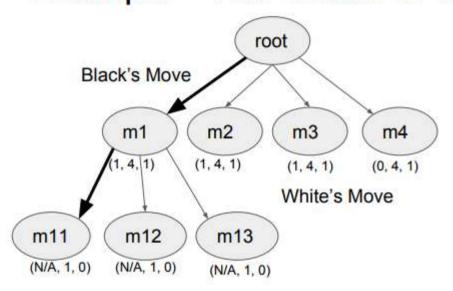
After first 4 iterations:
Suppose m1, m2, m3
black wins in simulation
and m4 white wins

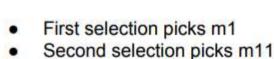
m1 m2 m3 m4

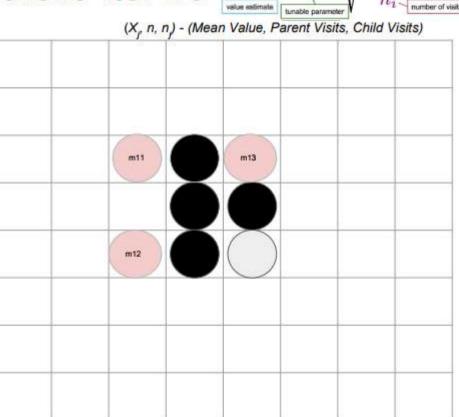
	$X_{j}$	n	<b>n</b> <sub>j</sub>
m1	1	4	1
m2	1	4	1
m3	1	4	1
m4	0	4	1

(X <sub>p</sub> n, 1	n <sub>j</sub> ) - (Mear	i value, P	arent VISI	is, Chila	risits)
	mt				
m2		•			
	•	m4	m3		

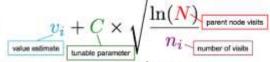


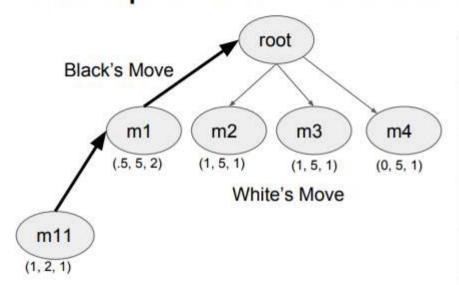




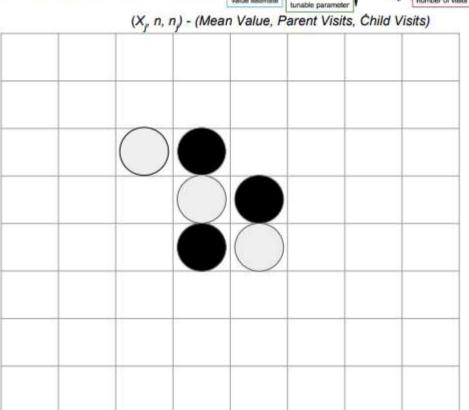


### Example - The Game of Othello Iter #5 $v_i + C \times C$

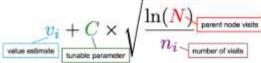


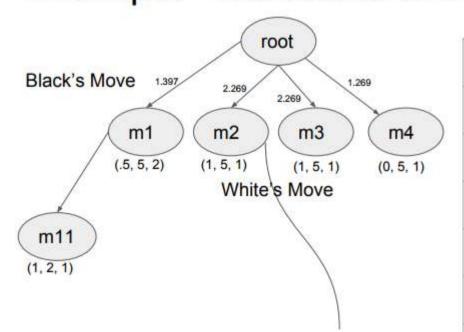


- Run a simulation
- White Wins
- Backtrack, and update mean scores accordingly.

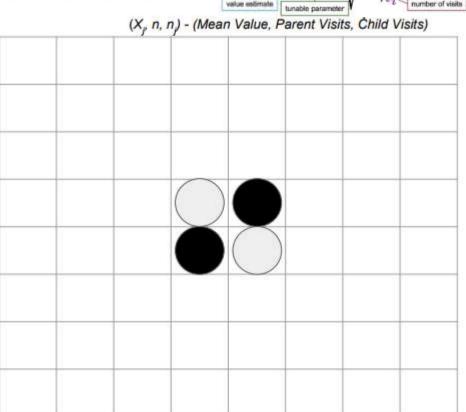


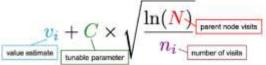
### Example - The Game of Othello Iter #6 \_\_vi+C×

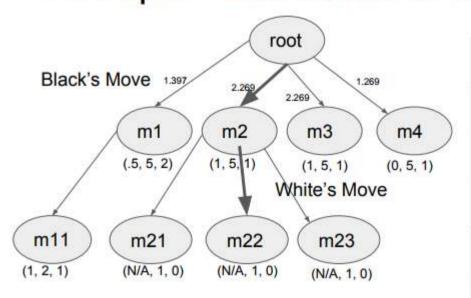




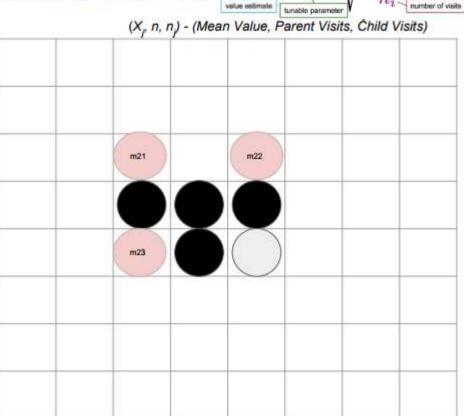
Suppose we first select m2

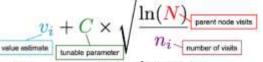


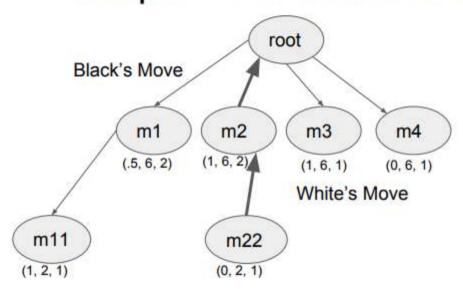


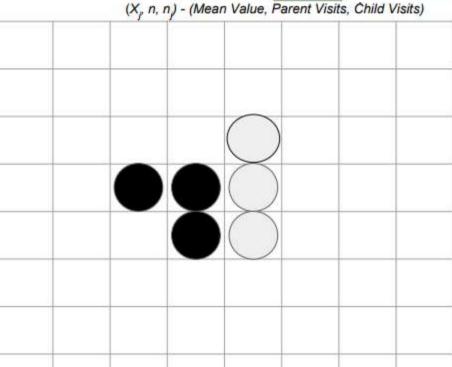


Suppose we pick m22

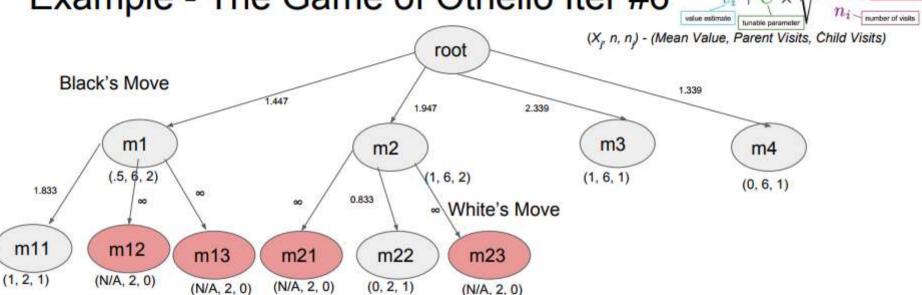








- Run simulated game from this position.
- Suppose black wins the simulated game.
- Backtrack and update values



- This is how our tree looks after 6 iterations.
- Red Nodes not actually in tree
- Now given a tree, actual moves can be made using max, robust, maxrobust, or other child selection policies.
- Only care about subtree after moves have been made

# AlphaGo

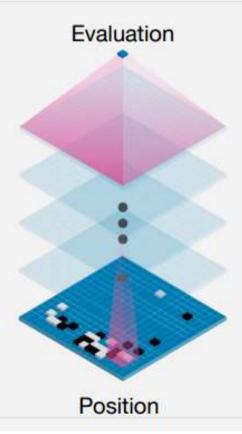


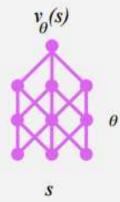
### AlphaGo

Uses a value network and policy network to augment MCTS

- 1. Policy network first trained on professional Go games and then trained further using reinforcement learning
- 2. Value network trained using self-play using the policy network
- 3. Then MCTS is run leveraging the two networks

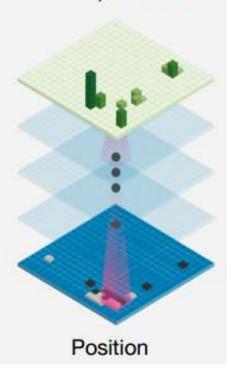
### Value network

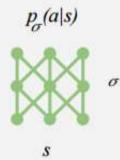




### Policy network

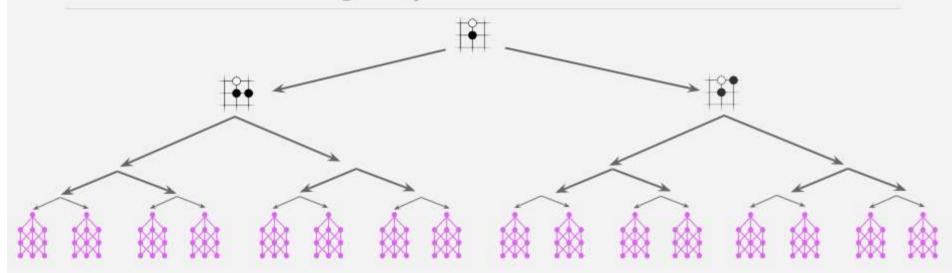
#### Move probabilities





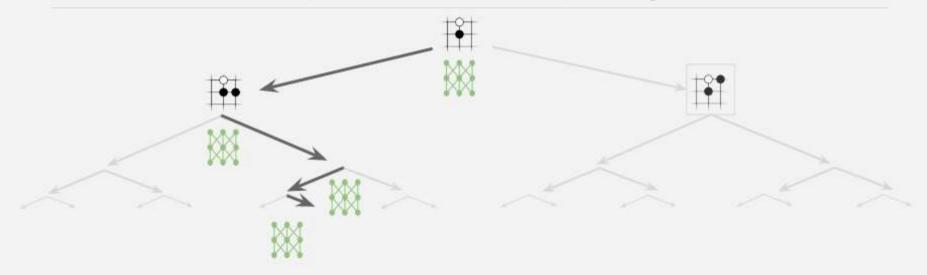
High-level idea 1:

### Reducing depth with value network

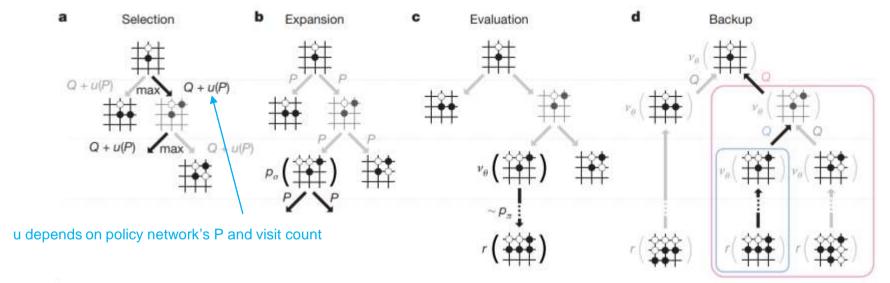


#### High-level idea 2:

### Reducing breadth with policy network



### AlphaGo's MCTS



**Figure 3** | **Monte Carlo tree search in AlphaGo. a**, Each simulation traverses the tree by selecting the edge with maximum action value Q, plus a bonus u(P) that depends on a stored prior probability P for that edge. **b**, The leaf node may be expanded; the new node is processed once by the policy network  $p_{\sigma}$  and the output probabilities are stored as prior probabilities P for each action. **c**, At the end of a simulation, the leaf node

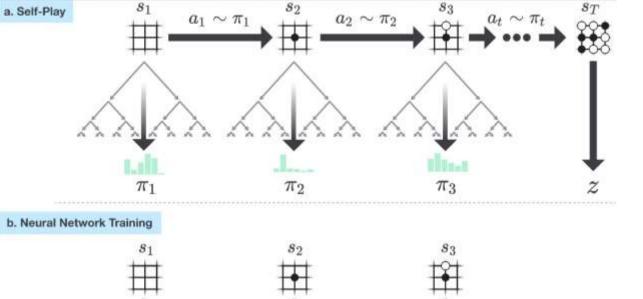
is evaluated in two ways: using the value network  $v_\theta$ ; and by running a rollout to the end of the game with the fast rollout policy  $p_\pi$ , then computing the winner with function r.  $\mathbf{d}$ , Action values Q are updated to track the mean value of all evaluations  $r(\cdot)$  and  $v_\theta(\cdot)$  in the subtree below that action.

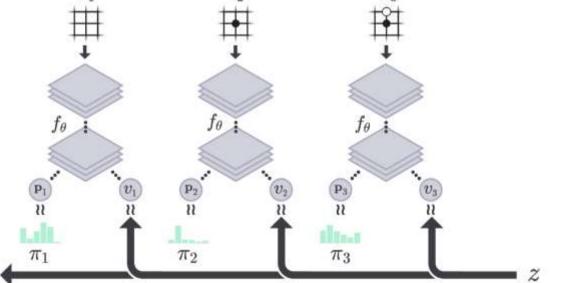
Once search is complete, the algorithm selects the most visited move from the root.

## AlphaGo Zero

### AlphaGo Zero

- No human data besides rules of the game
- The value and policy network are trained in self-play in the context of MCTS instead of human data or without search
  - MCTS as a policy improvement operator!
- Trained on 4 TPUs for 70 days
  - Compared to tens of thousands of TPUs for Gemini



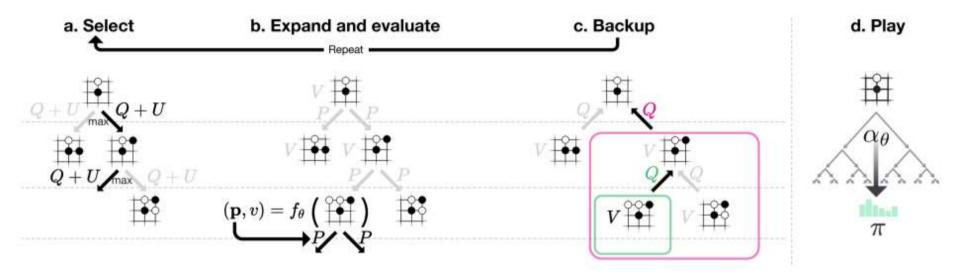


### **Neural Network Loss Function**

$$(p, v) = f_{\theta}(s)$$
 and  $l = (z - v)^2 - \pi^T \log p + c \|\theta\|^2$ 

Maximise similarity of the neural network move probabilities p to the search probabilities  $\pi$ 

### Search Algorithm



Once the search is complete, search probabilities  $\pi$  are returned proportional to N<sup>1/ $\mu$ </sup>, where N is the visit count of each move from the root state and  $\mu$  is a parameter controlling temperature

### Search Algorithm

- Each node s in the search tree contains edges (s, a) for all legal actions
- Each edge stores a set of statistics, {N(s, a), W(s, a), Q(s, a), P(s, a)}
  - N: number of visits to that edge
  - W: Total value
  - Q: Average value
  - P: Policy output

$$a_t = \operatorname{argmax}(Q(s_t, a) + U(s_t, a))$$

$$U(s, a) = c_{\text{puct}} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$$

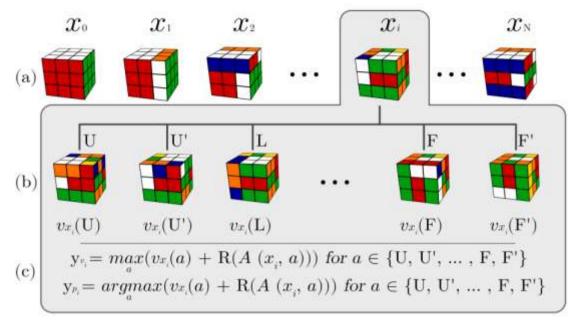
### Expand and Evaluate

- When we reach a leaf node, we run the state through the neural network to get a value estimate and policy estimate
- Each edge (N, W, Q) is initialized to 0
- Backup value

### Backup

- We update N, W, Q with the value that the neural network proposes
- N(s, a) = N(s, a) + 1
- W(s,a) = W(s,a) + v
- Q(s, a) = W(s, a) / N(s, a)

### These Techniques are Useful Also in Single-Agent Settings



E.g.1: Rubik's cube

McAleer et al. "Solving the Rubik's cube with approximate policy iteration." *ICLR*. 2018.

Agostinelli et al. "Solving the Rubik's cube with deep reinforcement learning and search." *Nature Machine Intelligence*. 2019

E.g.2: Edge test selection in kidney exchange McElfresh, Curry, Sandholm, Dickerson, "Improving Policy-Constrained Kidney Exchange via Pre-Screening", NeurIPS-20