Monte Carlo Tree Search

(based on a tutorial by Simon Lucas and Peter Cowling)

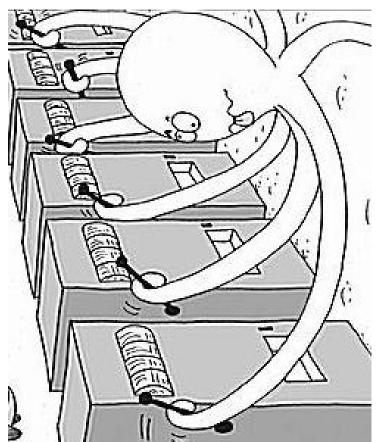
Bandit problems

At each step, pull one arm



Different average reward

Task: maximise reward (minimise regret)



Which arm to pull?

Which arm to pull?

- Pull all arms equally often?
- Only pull the arm that has given the best results so far?
- Mostly pull the "best" arm, but sometimes one of the others?
- An example of the exploration/exploitation dilemma
- Principled solution?

Which arm to pull?

Flat Monte Carlo Share trials uniformly between arms ε-Greedy P(1-ε) - Best arm so far P(ε) - Random arm

UCB1 (Auer et al (2002)). Choose arm *j* so as to maximise:

$$ar{X}_j + \sqrt{rac{2\log n}{T_j(n)}}$$
Mean Upper bound

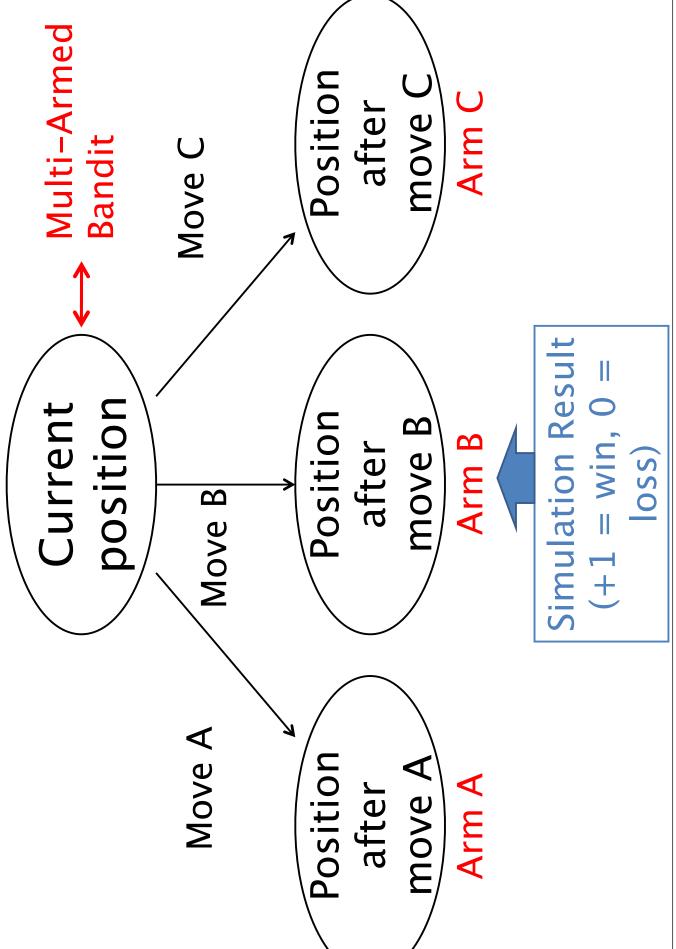
Mean Upper bound so far on variance

Tj(n) = number of times arm j was pulled n = number of plays so far

- Monte Carlo Tree Search (MCTS),
- Upper Confidence Bounds (UCB),
- Upper Confidence Bounds for Trees (UCT)

sequence of decisions with dependencies From making single decision to making a

Jame decisions



U C B

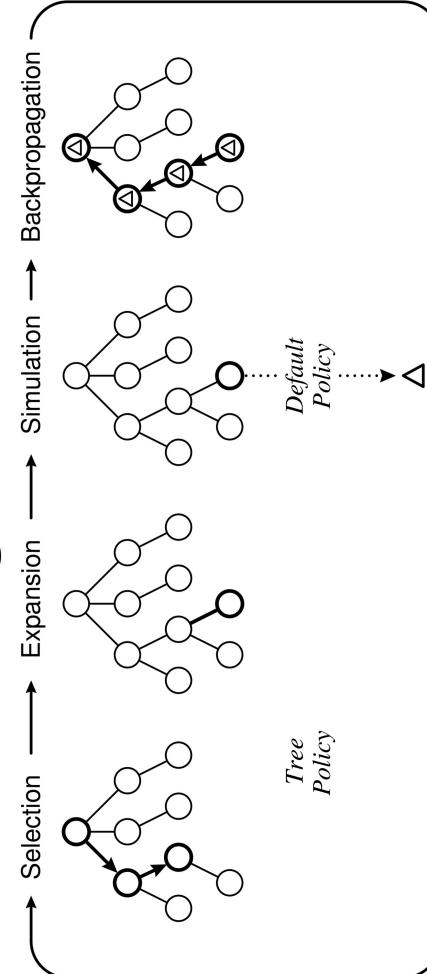
- Anytime stop whenever you like
- UCBI formula minimises regret
- Grows like log(n)
- Needs only game rules:
- Move generation
- Terminal state evaluation
- Surprisingly effective, but...
 ... doesn't look ahead

- Monte Carlo Tree Search (MCTS),
- Upper Confidence Bounds (UCB),
- Upper Confidence Bounds for Trees (UCT)

UCT (UCB on trees)

- Anytime
- Scalable
- Tackle complex games better than before
- May be logarithmically better with increased CPU
- No need for heuristic function
- Though often better with one

MCTS general idea



- Tree policy: choose which node to expand (not necessarily leaf of tree)
- Default (simulation) policy: random playout until end of game

MCTS algorithm

- Tree policy
- Expand
- Best Child (UCT Formula)
- Default Policy
- Back-propagate

Algorithm 1 General MCTS approach.

function MCTSSEARCH(s_0)

create root node v_0 with state s_0

while within computational budget do

 $v_l \leftarrow \text{TREEPOLICY}(v_0)$

 $\Delta \leftarrow \text{DEFAULTPOLICY}(s(v_l))$

 $\mathrm{BACKUP}(v_l,\Delta)$

return $a(BESTCHILD(v_0))$

Tree policy

if v not fully expanded then $v \leftarrow \text{BESTCHILD}(v, Cp)$ while v is nonterminal do return Expand(v)function TREEPOLICY (v) $\mathbf{return} \ v$ else

Note that node selected for expansion does nonterminal test refers to the game state) not need to be a leaf of the tree (the

choose $a \in \text{untried actions from } A(s(v))$ Tree expansion with s(v') = f(s(v), a)add a new child v' to vfunction Expand(v)and a(v') = areturn v'

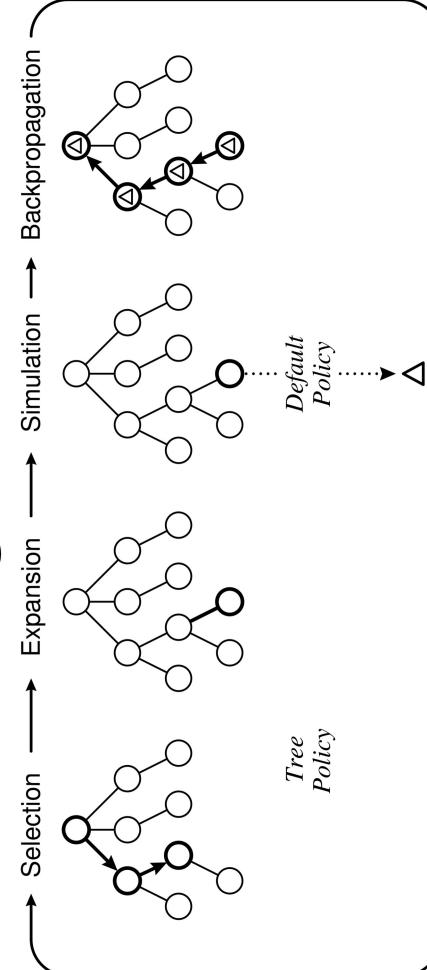
Best child (UCT)

function Bestchild
$$(v, c)$$

eturn arg max
$$\frac{Q(v')}{V(c)} + c\sqrt{\frac{2 \ln N(v)}{N(v')}}$$

- Standard UCT equation (compare UCB)
- Higher values of c lead to more exploration

MCTS general idea



- Tree policy: choose which node to expand (not necessarily leaf of tree)
- Default (simulation) policy: random playout until end of game

Default policy (rollout)

choose $a \in A(s)$ uniformly at random while s is non-terminal do function DefaultPolicy(s) return reward for state s $s \leftarrow f(s, a)$

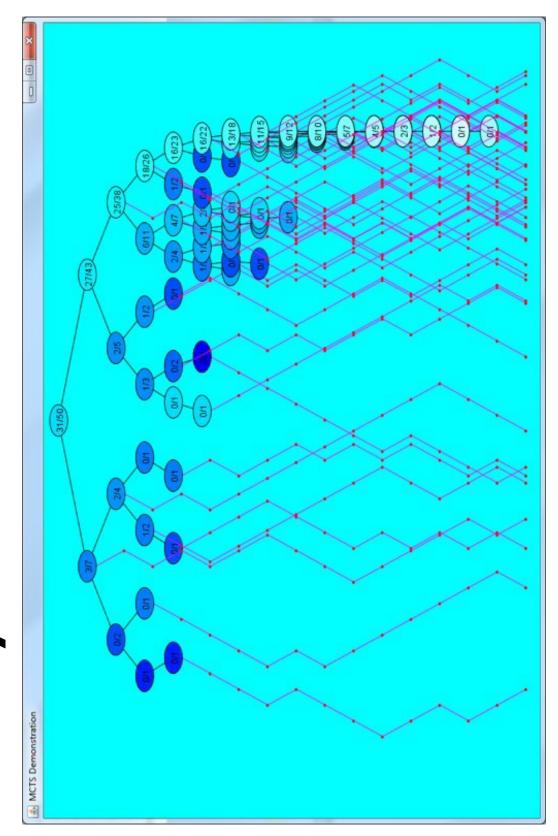
- Each time a node is added to the tree, the default policy plays out until the terminal state of the game
- The standard is to do this uniformly randomly

Backup

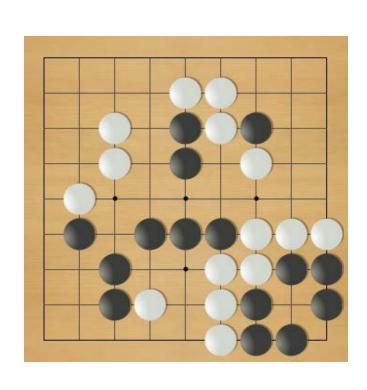
 $Q(v) \leftarrow Q(v) + \Delta(v, p)$ while v is not null do $N(v) \leftarrow N(v) + 1$ function Backup (v, Δ) $v \leftarrow \text{parent of } v$

- v is the new node added to the tree by the tree policy
- Back up the values from the added node up the tree to the root

asymmetric trees MCTS builds



Computer Go

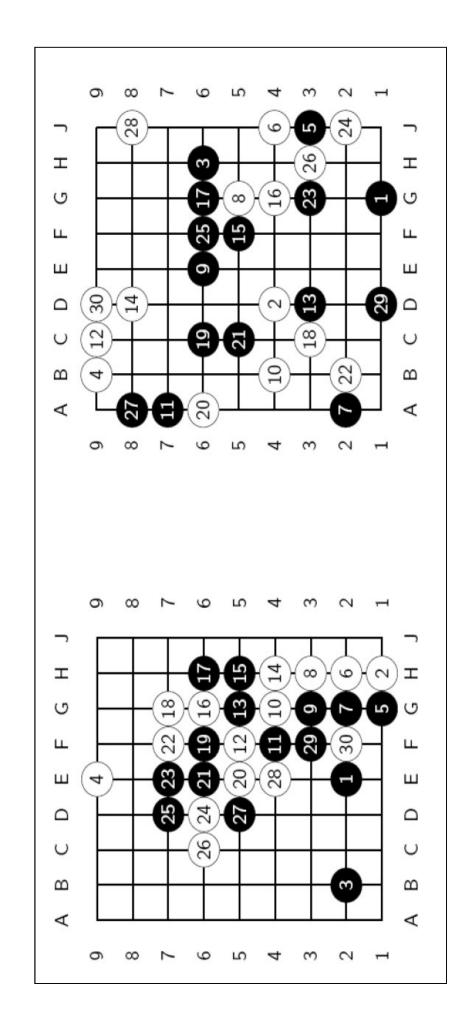


Default Policy

- The default policy is crucial to UCT
- ❖Better default policy => better UCT (?)
 - As hard as the overall problem
- Default policy must also be fast

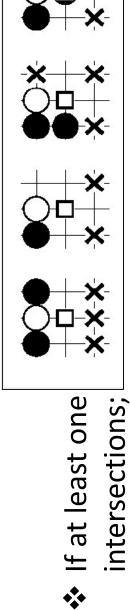
Educated simulations: Sequence-like simulations

Sequences matter!



How it works in MoGo

- Look at the 8 intersections around the previous move
- * For each such intersection, check the match of a pattern (including symetries)



matching

Else play uniformly among legal moves

Random mode	Win. Rate	Win. rate	Total
	for B. Games	for W. Games	Win. Rate
Uniform	46% (250)	36% (250)	$41.2\% \pm 2.2\%$
Sequence-like	77% (400)	82% (400)	$80\%\pm1.4\%$

Default Policy (continued)

- The default policy is crucial to UCT
- ❖Better default policy => better UCT (?)
 - *As hard as the overall problem
- Default policy must also be fast