Monte Carlo Tree Search

tree ← root

**Until out of time**

traverse tree     root ⇝ leaf
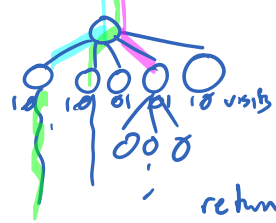
expand     if leaf expandable, add its children

simulate     play to terminal pos (from arbi. child if newly expanded)
can be random   default policy

update     backprop result up tree from start of simulation ⇝ root

return move to child w/ best stats ( highest mean reward or highest visit count )

tree policy
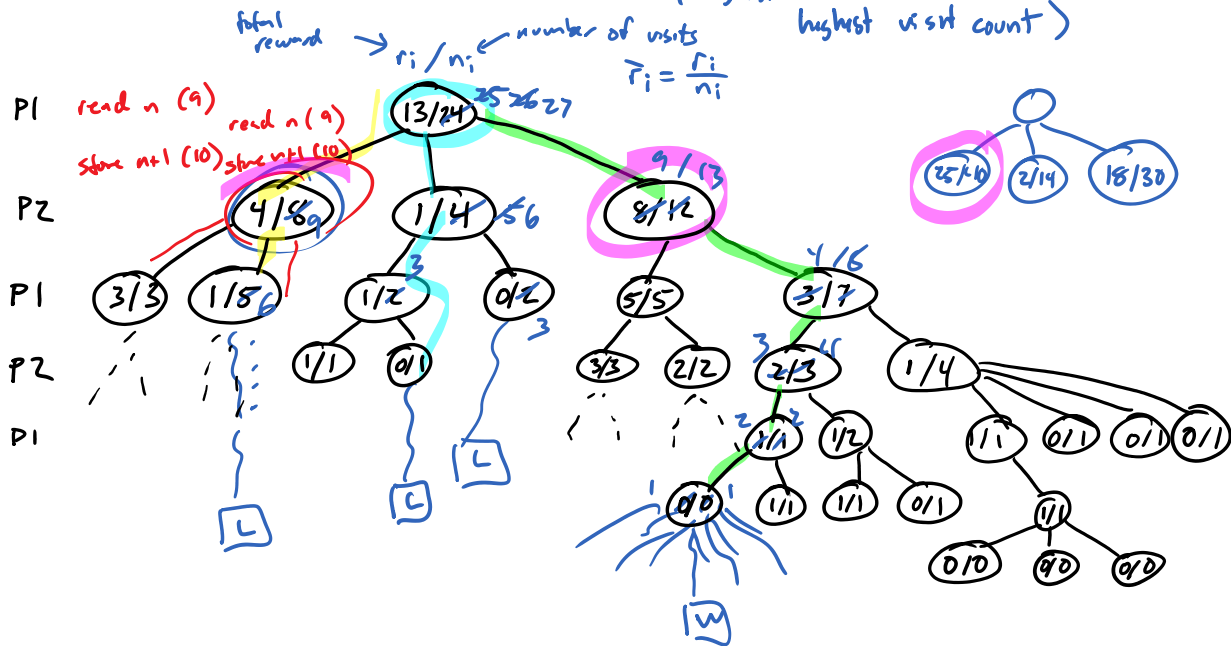UCB (choosing a child w/ 0 visits when one exists)

$$\bar{r_i} + \sqrt{\frac{2 \ln T}{n_i}}$$
exploit   exploration
non-terminal and non-zero visits   reward for P1

P2
max: $(c - \bar{r_i}) + \sqrt{\frac{2 \ln T}{n_i}}$
mm $\bar{r_i} - \sqrt{\frac{2 \ln T}{n_i}}$

UCT = MCTS + UCB



$10 \quad 10 \quad 01 \quad 01 \quad 10$ visits
$00 \quad 0$

total reward → $r_i / n_i$ ← number of visits
$\bar{r_i} = \frac{r_i}{n_i}$

P1   read n (9)   read n (9)
store n+1 (10)   store n+1 (10)

P2

P1

P2

P1



UCT = MCTS

advantages : convergent     converges to minimax (given enough time)

anytime     returns a move at any point

no domain knowledge     no heuristic needed

easily parallelized     leaf — multiple parallel playouts from leaves
tree — parallel traversals through tree
root — separate tree for each thread
combine results ( majority rule combining stats )

disadvantages : no domain knowledge

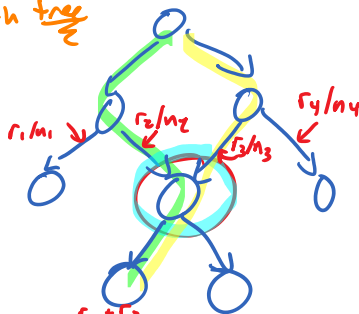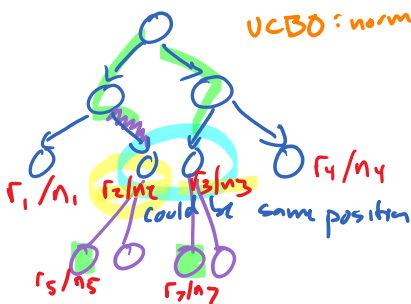some games not amenable
trap state; most moves bad

some games not amenable

trap state; most moves bad
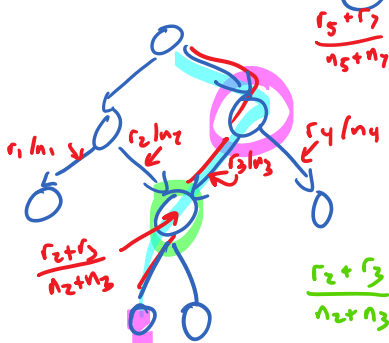one good move



MCTS adapted for games that aren't trees

UCB0: normal MCTS with $\frac{tree}{z}$

$r_1/n_1$  $r_2/n_2$  $r_3/n_3$  $r_4/n_4$

could be same position

$r_5/n_5$  $r_7/n_7$

UCB1: merge nodes but
use stats on
outgoing edges

$r_1/n_1$  $r_2/n_2$  $r_4/n_4$  $r_3/n_3$

implicit tree:
edges not stored anywhere

$\frac{r_5 + r_7}{n_5 + n_7}$

UCB2: combines observed reward
over all paths to node
(exploit)

not for exploration term

UCB3: backprop
along all paths
leaf → root
instead of one
you came down

$r_1/n_1$  $r_2/n_2$  $r_4/n_4$  $r_3/n_3$

$\frac{r_2 + r_3}{n_2 + n_3}$

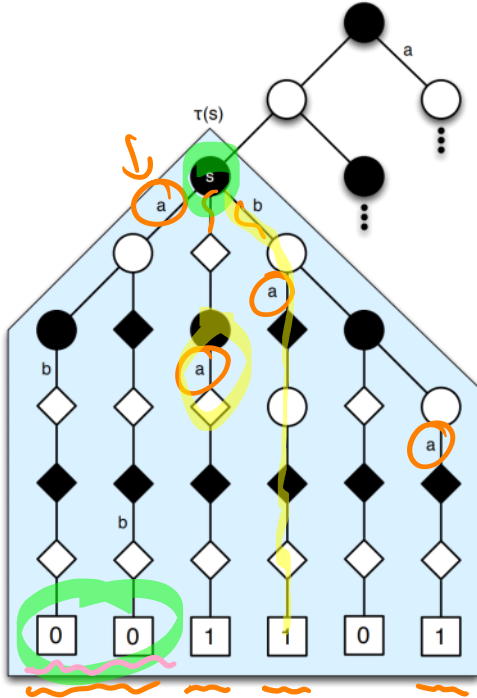$$\frac{r_2 + r_3}{n_2 + n_3} + \sqrt{\frac{2 \ln\left[n_3 + n_4\right]}{n_3}}$$

$$\frac{r_2 + r_3}{n_2 + n_3} + \sqrt{\frac{2 \ln\left[(n_2 + n_3) + n_4\right]}{n_2 + n_3}}$$

breaks
convergence

# MC-RAVE

↙ rapid averaging value estimation



$Q(s,a) = 0/2$

$Q(s,b) = 2/3$

$\tilde{Q}(s,a) = 3/5$

$\tilde{Q}(s,b) = 2/5$

From Gelly and Silver, Monte-Carlo tree search and rapid action value estimation in computer Go. Artif. Intell. 175, 1856-1875, 2011

all moves as first (AMAF)

$Q(s,a)$ : observed reward of action a in state s

$= \frac{0}{2}$

$\tilde{Q}(s,a)$ : uses AMAF heuristic
obs reward for a over
entire subtree rooted at s

weight of $Q(s,a)$ vs $\tilde{Q}(s,a)$

$$Q_*(s,a) = (1 - \beta(s,a)) \, Q(s,a) + \beta(s,a) \tilde{Q}(s,a)$$

↓

$\sqrt{\dfrac{k}{3N(s)+k}}$   for Go  (k ≈ 1000)  Fig 5

# visits for s