# Monte Carlo Tree Search (cont.)

UAlg FCT
UNIVERSIDADE DO ALGARVE
FACULDADE DE CIÊNCIAS E TECNOLOGIA

---

# Classes for next week



| | SEG. 30 | TER. 31 | QUA. 1 | QUI. 2 |
|---|---|---|---|---|
| 09 | REDES DE ... | REDES DE ... 09 - 11 | | DESENVO... |
| 10 | ANÁLISE ... | | | DESENVO... 10 - 12 |
| 11 | | INTERFAC... | | |
| 13 | | DESENVO... | | INTELIGÊ... |
| 14 | | INTERFAC... | | |
| 15 | INTERFAC... | INTERFAC... 15 - 17 | | |
| 16 | INTERFAC... 15:30 - 17:... | | | |
| 17 | | ANÁLISE ... | | |

## Monte Carlo Tree Search (Kocsis & Szepesvari, 2006)



Selection → Expansion → Simulation → Backpropagation

*Tree Policy*

*Default Policy*

## UCT Algorithm (Kocsis & Szepesvari, 2006)

Current World State



Tree Policy

What is an appropriate tree policy?

When all node actions tried once, select action according to tree policy

## Played them again…

-1€, 20€        2€, 3€        10€, 2€

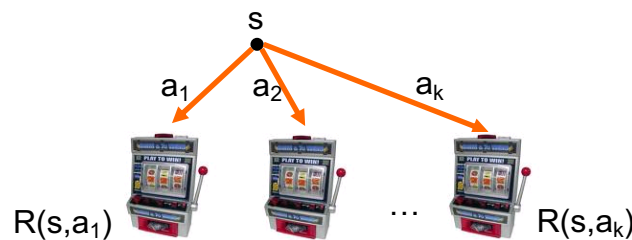- Problem: which machine should you play next?

http://tiny.cc/ia2022-23

## Multi-armed Bandit Problem, informally

- Find an arm-pulling strategy such that the expected total reward at time *n* is close to the best possible (*i.e.*, always pulling the best arm)

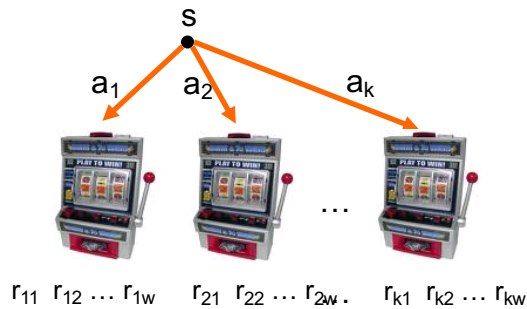  - Must balance **exploring** machines to find good payoffs and **exploiting** current knowledge

s

$a_1$   $a_2$        $a_k$

$R(s,a_1)$        …        $R(s,a_k)$

http://tiny.cc/ia2022-23 21-6

3

## Uniform (Naïve) bandit algorithm (Even-Dar et. al., 2002)

1. Pull each arm *w* times (uniform pulling).
2. Return arm with best average reward.



Uniform is a poor choice – waste of time on bad arms

## Exploration vs. exploitation

- When to explore (try a new machine)?

- When to stop exploring and start exploiting (play the apparently best machine)?

# UCT Algorithm (Kocsis & Szepesvari, 2006)

Exploration Parameter

- **Tree policy is guided by:**

$$\pi_{UCT}(s) = \arg\max_a Q(s,a) + c\sqrt{\frac{\ln n(s)}{n(s,a)}}$$

- Q(s,a): average reward received in current trajectories after taking action a in state s, Q(s,a)=wins(s,a)/n(s,a)
- n(s): number of times s encountered in the sims
- n(s,a): number of times action a taken in s

# UCT Algorithm (Kocsis & Szepesvari, 2006)

- Q(s,a)=wins(s,a)/n(s,a)
- n(s): number of times s encountered in the sims
- n(s,a): number of times action a taken in s

Exploration Parameter

$$\pi_{UCT}(s) = \arg\max_a Q(s,a) + c\sqrt{\frac{\ln n(s)}{n(s,a)}}$$

**Value Term:**
favors actions that looked good historically

**Exploration Term:**
actions get an exploration bonus that grows with ln *n(s)*

Doesn't waste much time on sub-optimal arms unlike uniform!

## UCT Algorithm (Kocsis & Szepesvari, 2006)

Current World State



$$\pi_{UCT}(s) = \arg\max_a Q(s,a) + c\sqrt{\frac{\ln n(s)}{n(s,a)}}$$

$$\frac{w_i}{n_i} + c\sqrt{\frac{\ln t}{n_i}}$$

Tree Policy

$w_i$- number of wins in node $i$
$n_i$ – number of sims in $i$
$t$ – number of sims in the father of $i$

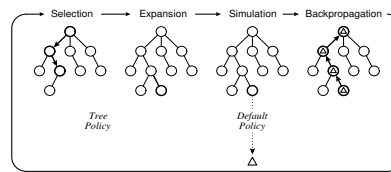When all node actions tried once, select action according to tree policy

---

## UCT Algorithm (Kocsis & Szepesvari, 2006)

Current World State



Tree Policy

In this level, corresponding to opponent turn, take the min of

$\pi_{UCT}(s)$

# UCT Recap

- To select an action at a state *s*
  - Build a tree using *N* iterations of monte-carlo tree search
    - Default policy is uniform random
    - Tree policy is based on UCT rule, i.e., select action *a* that maximizes

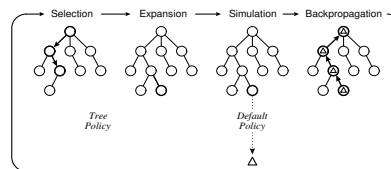$$\pi_{UCT}(s) = \arg\max_a Q(s,a) + c\sqrt{\frac{\ln n(s)}{n(s,a)}}$$

$$\frac{w_i}{n_i} + c\sqrt{\frac{\ln t}{n_i}}$$

- More simulations more accuracy

30-Oct-23                                                                 21-13
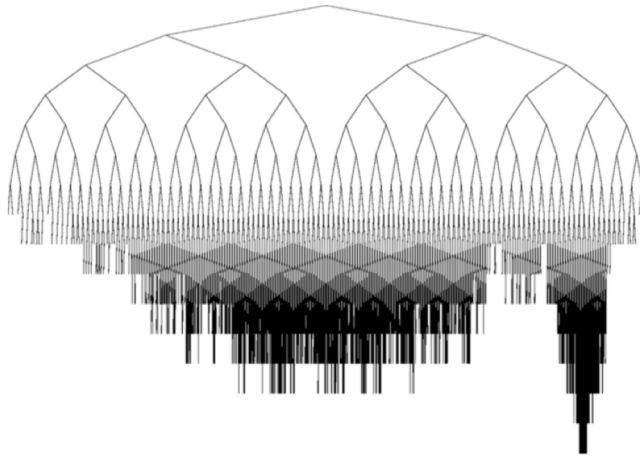
---

# UCT Recap

**Algorithm 1** General MCTS approach.

**function** MCTSSEARCH($s_0$)
    create root node $v_0$ with state $s_0$
    **while** within computational budget **do**
        $v_l \leftarrow$ TREEPOLICY($v_0$)
        $\Delta \leftarrow$ DEFAULTPOLICY($s(v_l)$)
        BACKUP($v_l, \Delta$)
    **return** $a($BESTCHILD($v_0$)$)$

30-Oct-23                                                                 21-14

## Result: Partial exploration of the game tree



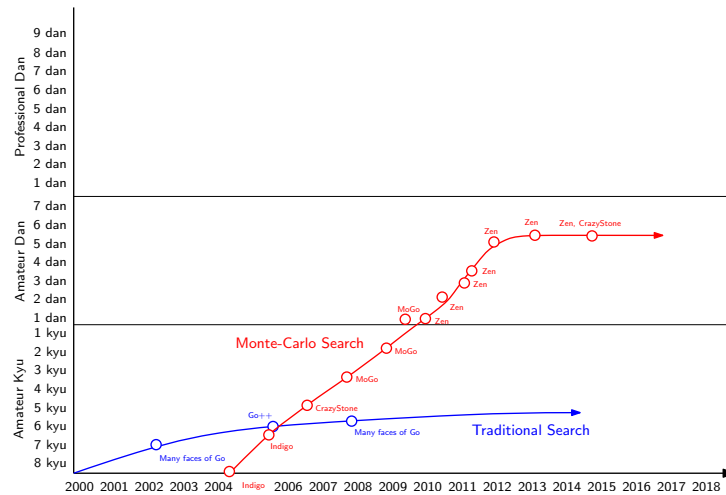21-15

## Some Improvements

- Use domain knowledge to figure out a better than random (rollout) policy
  - e.g., don't choose obviously bad actions

- Use domain knowledge to predict the opponent actions

- Learn a heuristic function to evaluate positions and use it to evaluate and select leaf nodes for further simulations

30-Oct-23          http://tiny.cc/ia2022-23          21-16
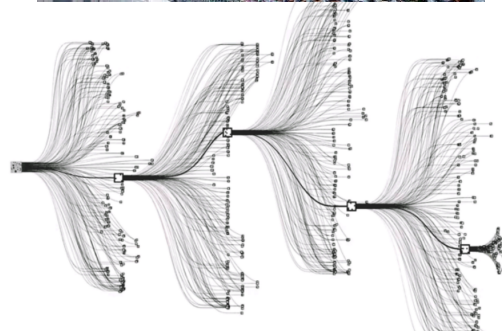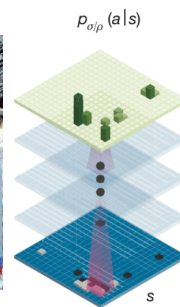
8

# Evo of computer GO



http://tiny.cc/ia2022-23

# Meet AlphaGo
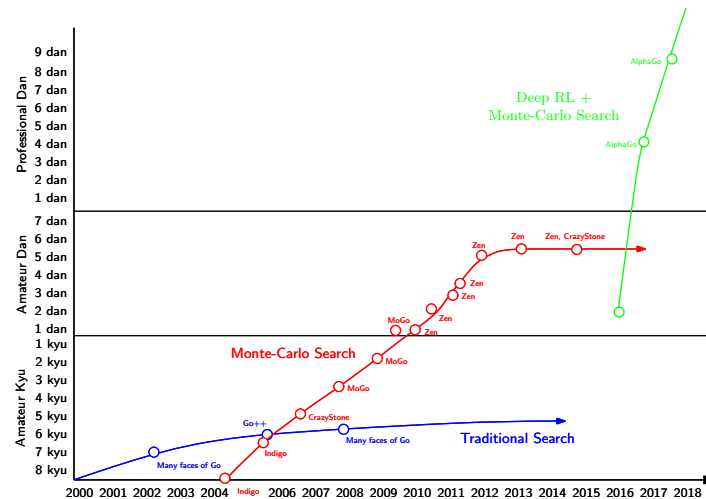
# Evo of computer GO

# MCTS/UCT take home points

- State of the Art any-time look-ahead probabilistic search for sequential decisions in the presence of uncertainty.
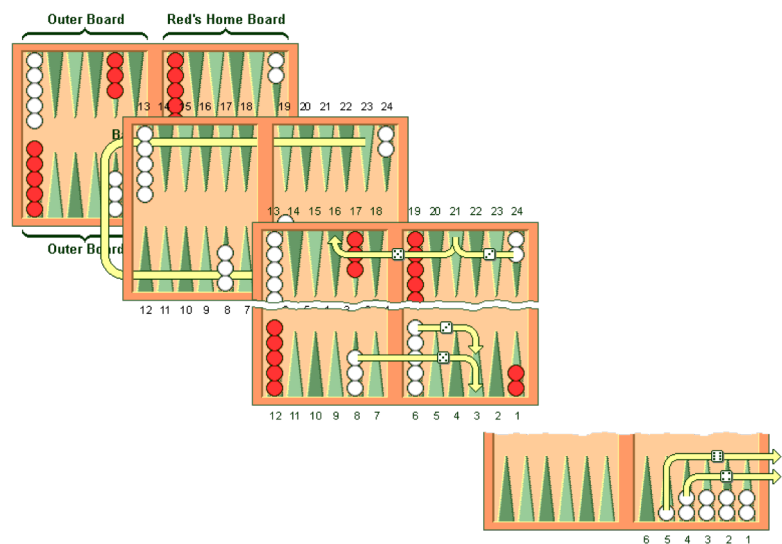
# Previously on AI: A classification of games

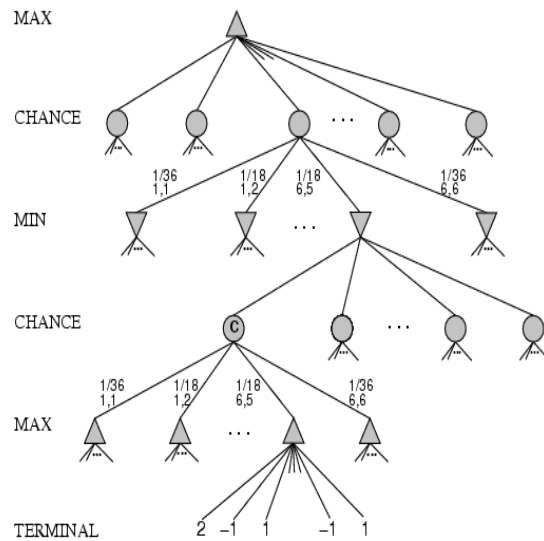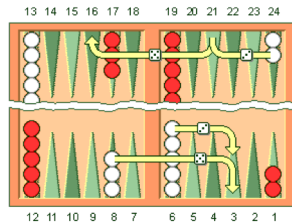|  | deterministic | chance |
|---|---|---|
| perfect information | chess, checkers, go, othello | backgammon monopoly |
| imperfect information |  | bridge, poker, scrabble nuclear war |

# Backgammon: A game that includes chance

# Chance nodes

MAX

CHANCE

MIN

CHANCE

MAX

TERMINAL

# Expected minimax value

EXPECTED-MINIMAX-VALUE($n$) =
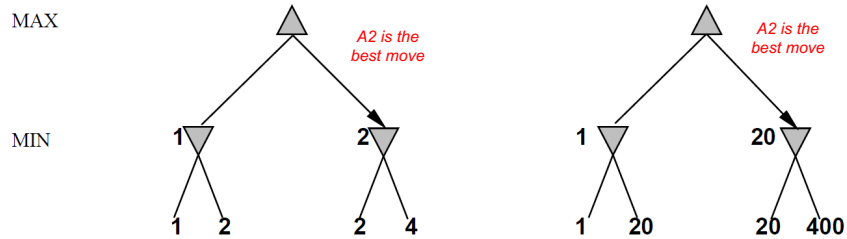
$$
\begin{cases}
\text{UTILITY}(n) & \text{If } n \text{ is a terminal} \\
\max_{s \in successors(n)} \text{MINIMAX-VALUE}(s) & \text{If } n \text{ is a max node} \\
\min_{s \in successors(n)} \text{MINIMAX-VALUE}(s) & \text{If } n \text{ is a min node} \\
\sum_{s \in successors(n)} P(s) * \text{EXPECTEDMINIMAX}(s) & \text{If } n \text{ is a chance node}
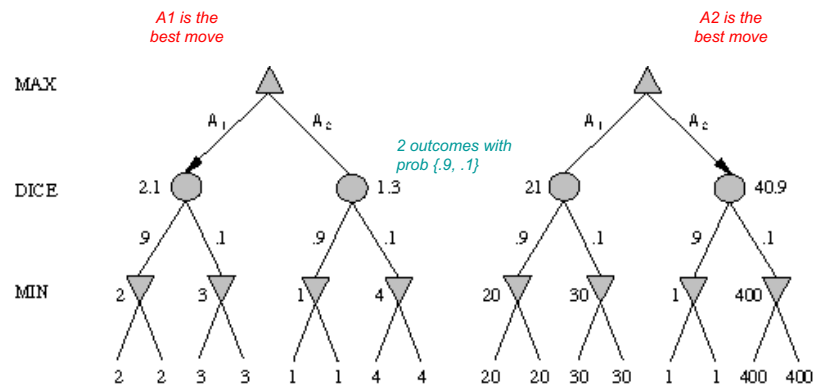\end{cases}
$$

## Minimax: exact values of Eval are irrelevant

MAX

*A2 is the best move*

MIN

1      2

1    2     2    4

*A2 is the best move*

1      20

1    20    20    400

- Behavior is preserved under any monotonic transformation of Eval function.

## Expected minimax value, example of

*A1 is the best move*

*A2 is the best move*

MAX

$A_1$    $A_2$

*2 outcomes with prob {.9, .1}*

DICE    2.1      1.3      21      40.9

.9   .1    .9   .1    .9   .1    .9   .1

MIN    2    3    1    4     20    30    1    400

2   2   3   3    1   1   4   4    20   20   30   30    1   1   400   400

## Backgammon in practice

- Dice rolls increase the branching factor: 21 possible outcomes with 2 dices

- α-β pruning becomes much less effective.

- TDGammon: uses depth-2 search + very good Eval (neural-network with reinforcement learning);

  ranked among the top three players in the world

## Classes for next week

| | SEG. 30 | TER. 31 | QUA. 1 | QUI. 2 |
|---|---|---|---|---|
| 08 | | | | |
| 09 | REDES DE ... | REDES DE ... 09 - 11 | | DESENVO... |
| 10 | ANÁLISE ... | | | DESENVO... 10 - 12 |
| 11 | | INTERFAC... | | |
| 12 | | | | |
| 13 | | DESENVO... | | INTELIGÊ... |
| 14 | | INTERFAC... | | |
| 15 | INTERFAC... | INTERFAC... 15 - 17 | | |
| 16 | INTERFAC... 15:30 - 17:... | | | |
| 17 | | ANÁLISE ... | | |
| 18 | | | | |