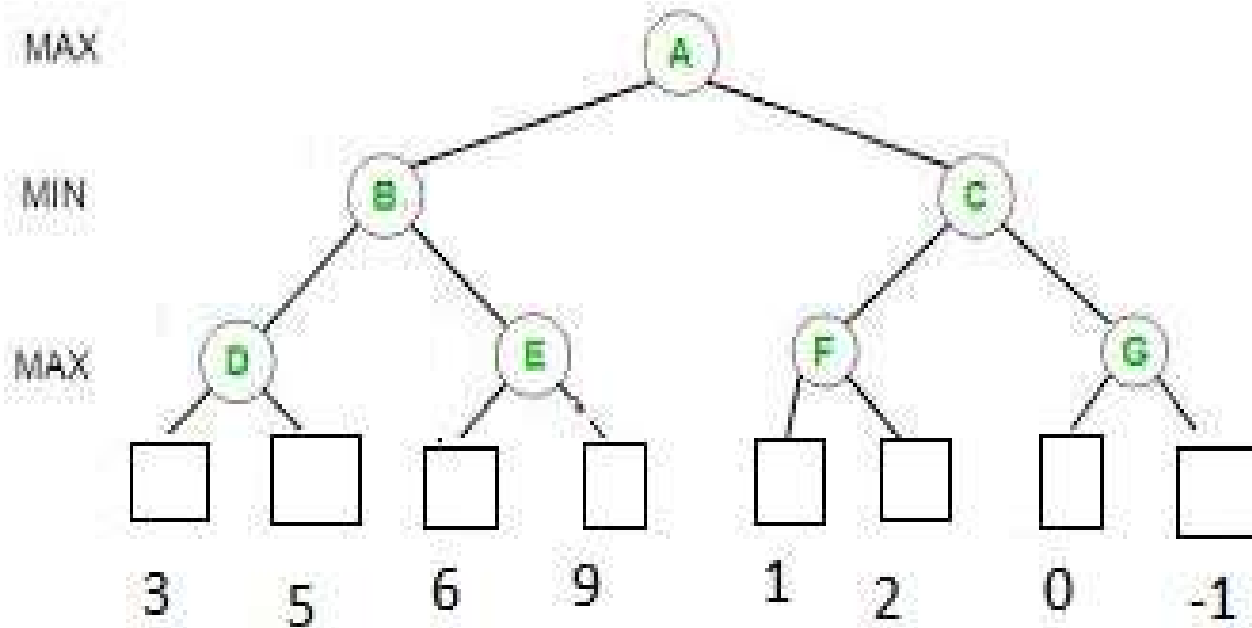




MONTE CARLO TREE SEARCH



RECAP



Rules for Alpha-beta Pruning

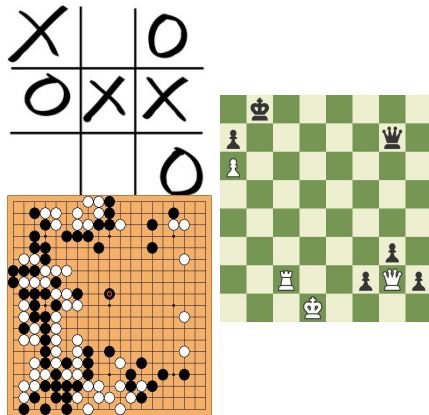
- **Alpha Pruning:** Search can be stopped below any MIN node having a beta value less than or equal to the alpha value of any of its MAX ancestors.
- **Beta Pruning:** Search can be stopped below any MAX node having a alpha value greater than or equal to the beta value of any of its MIN ancestors.

PRE-MCTS ALGORITHMS

- Deterministic, Fully Observable Games
- “Perfect information”
- Can construct a tree that contains all possible outcomes because everything is fully determined

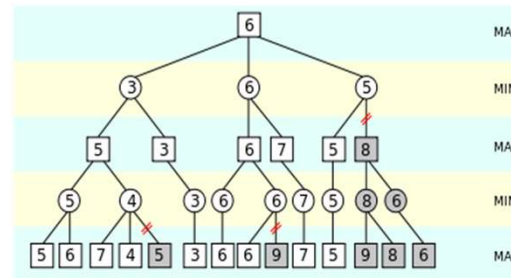
CONVENTIONAL GAME TREE SEARCH

Perfect-information games



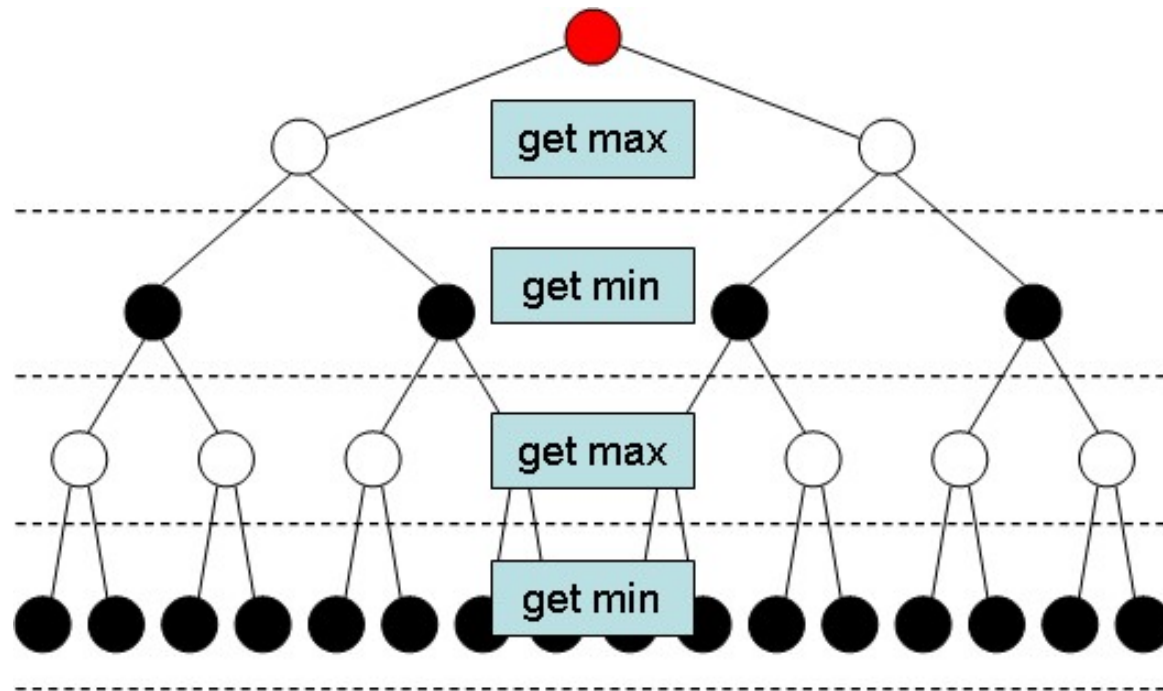
- All aspects of the state are fully observable

Minimax algorithm with alpha-beta pruning

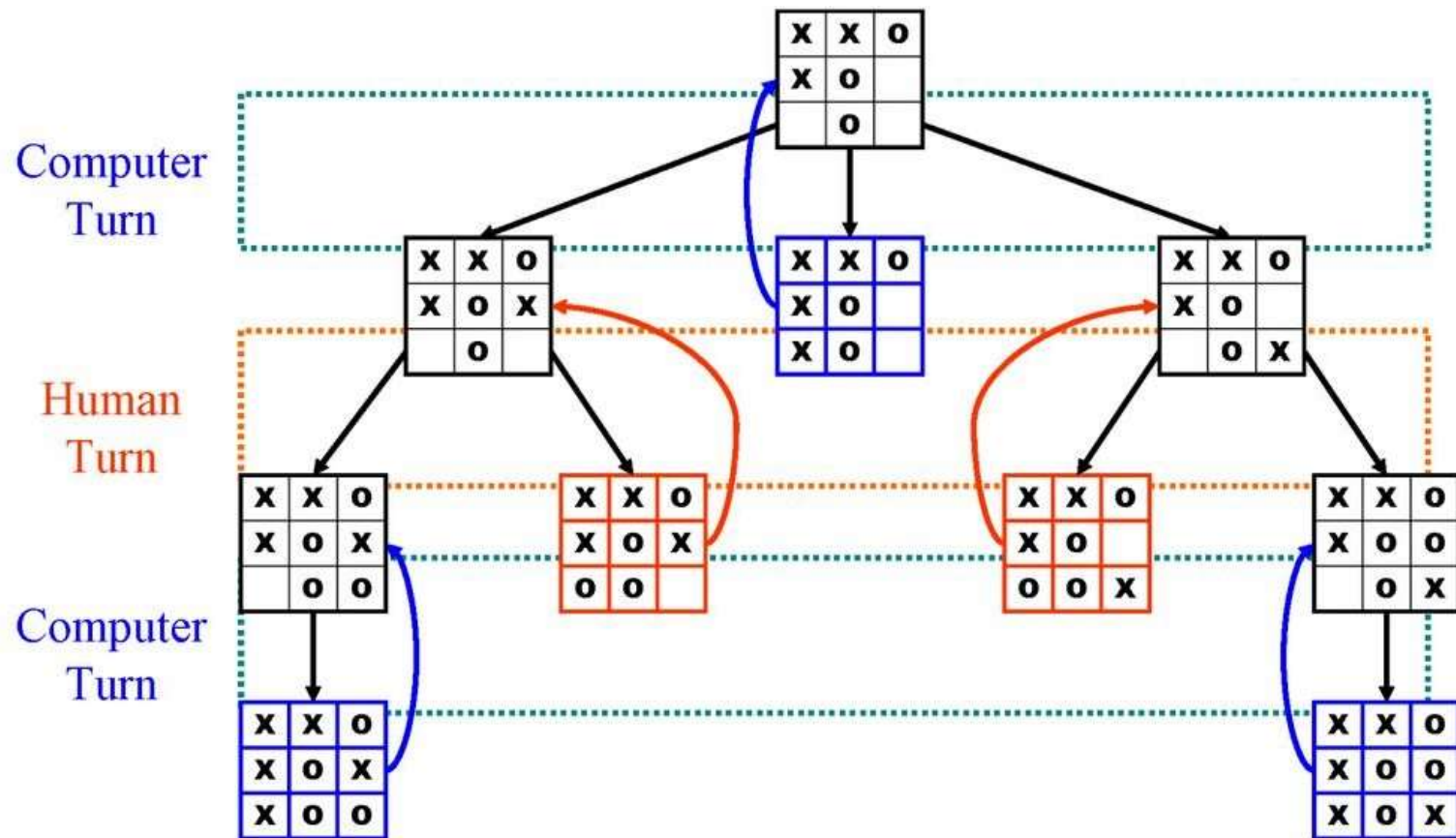


- Effective for
 - Modest branching factor
 - A good heuristic value function is known

MINIMIZE THE MAXIMUM POSSIBLE LOSS



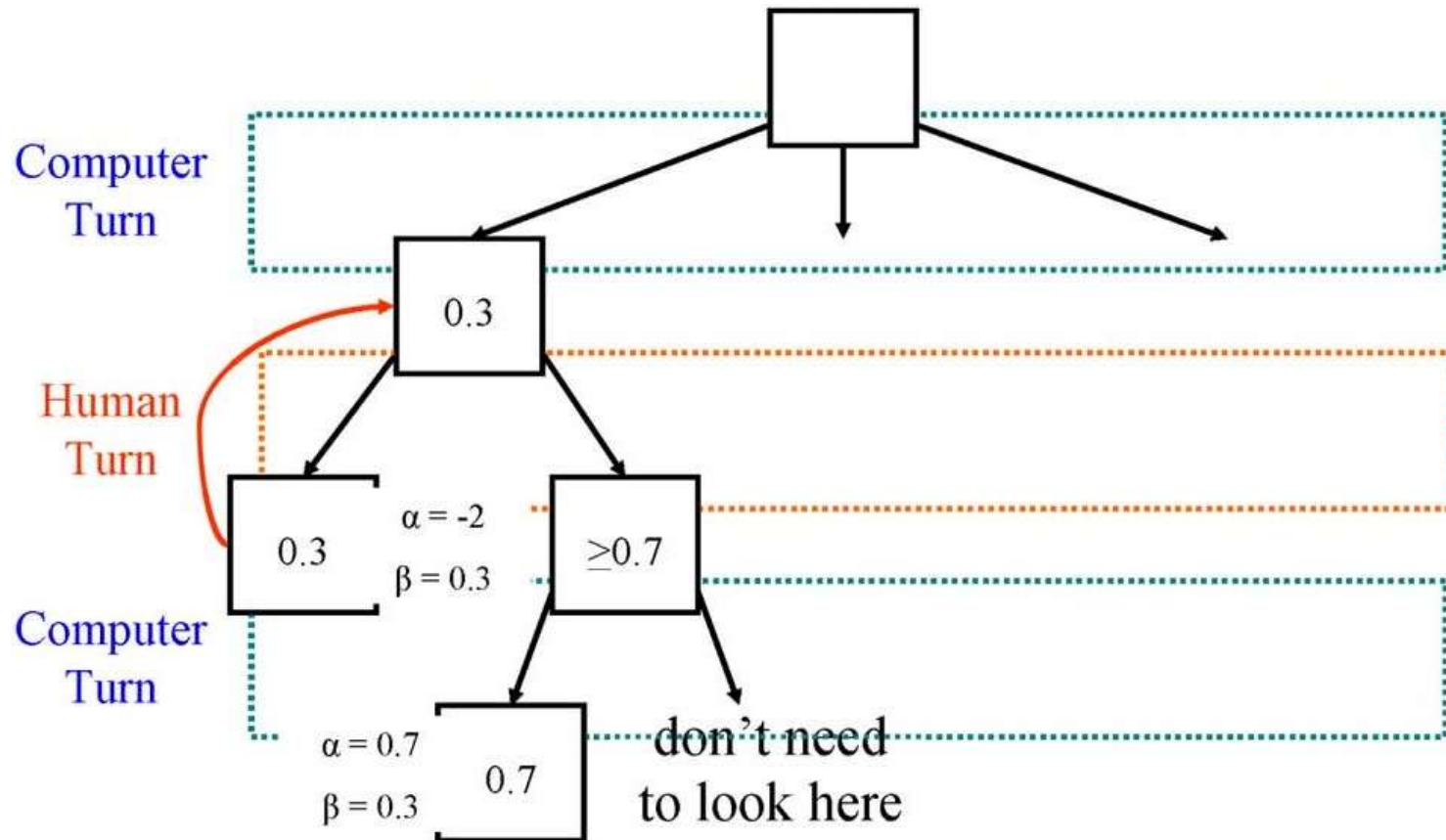
MINIMAX



ALPHA-BETA PRUNING

Prunes away branches that cannot influence the final decision

ALPHA - BETA



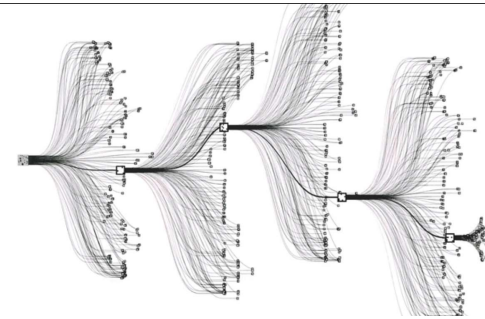
2^4 VS. 2^{250}



© Macmillan Publishers Limited, part of Springer Nature. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>.

GO

- Level weak to intermediate with alpha-beta
- Branching factor of Go is very large
 - 250 moves on average, game length > 200 moves
 - Order of magnitude greater than the branching factor of 20 for Chess
- Lack of a good evaluation function
 - Too subtle to model: similar looking positions can have completely different outcomes



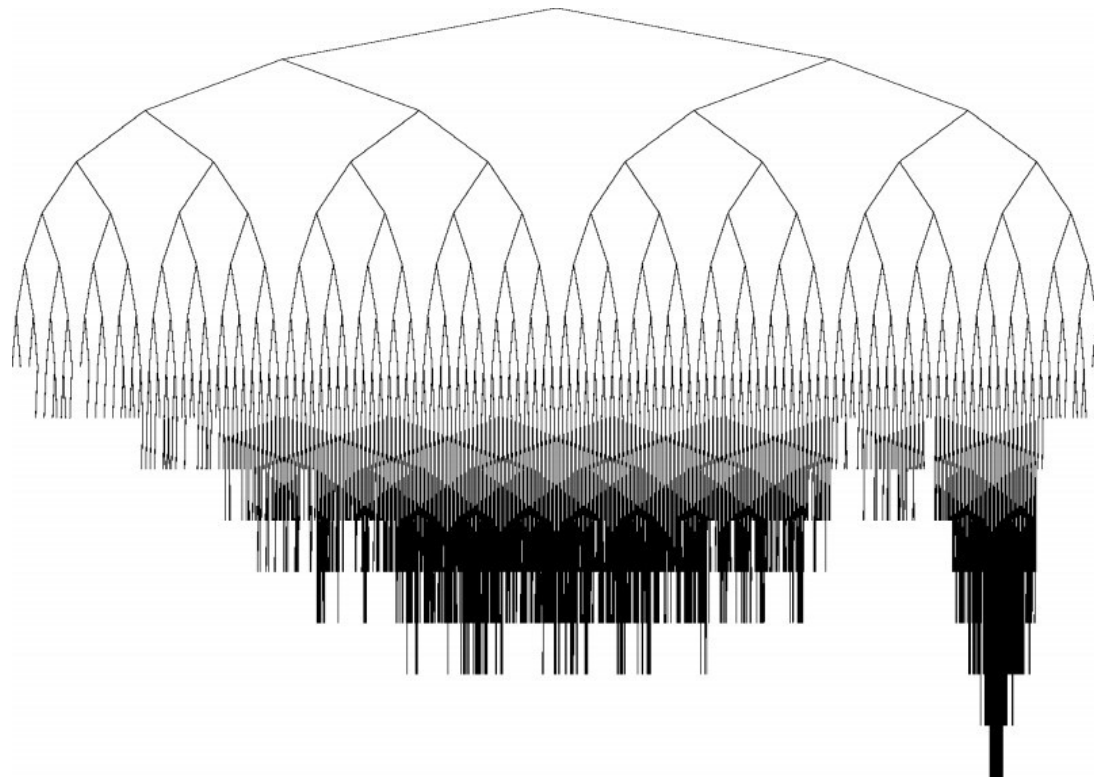
MONTE CARLO!

- u Revolutionized the world of computer Go
- u Application to deterministic games pretty recent (less than 10 years)
- u Explosion in interest, applications far beyond games
 - u Planning, motion planning, optimization, finance, energy management

ALPHAGO AND MONTE CARLO TREE SEARCH



ASYMMETRIC TREE EXPLORATION

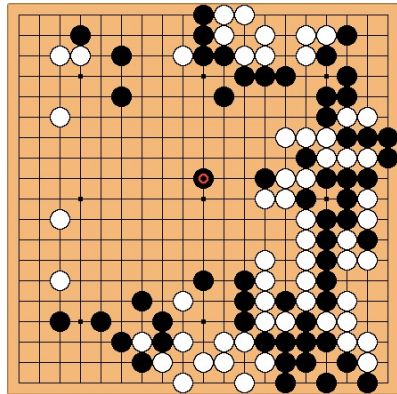


From *Bandit Algorithms for Tree Search*, Coquelin and Munos, 2007

MCTS FOR COMPUTER GO AND MINESWEEPER

Go: deterministic transitions

MineSweeper: probabilistic transitions



BASIC MONTE CARLO SIMULATION

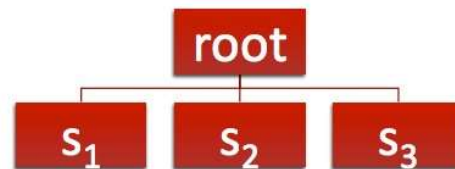
- u No evaluation function?

- Simulate game using random moves
- Score game at the end, keep winning statistics
- Play move with best winning percentage
- Repeat

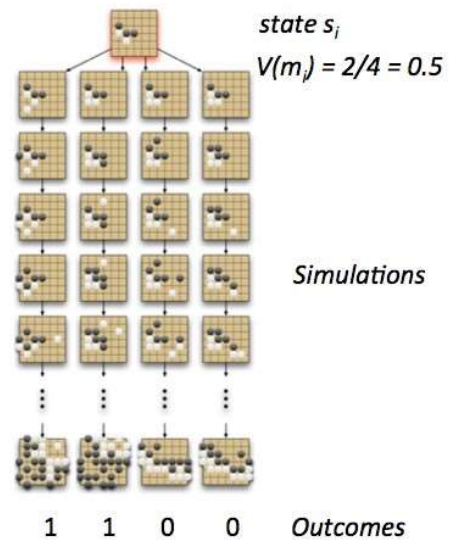
7

Use this as the evaluation function, hopefully it will preserve *some* difference between a good position and a bad position

BASIC MONTE CARLO SEARCH



1 ply tree
root = current position
 s_1 = state after move m_1
 s_2 = ...

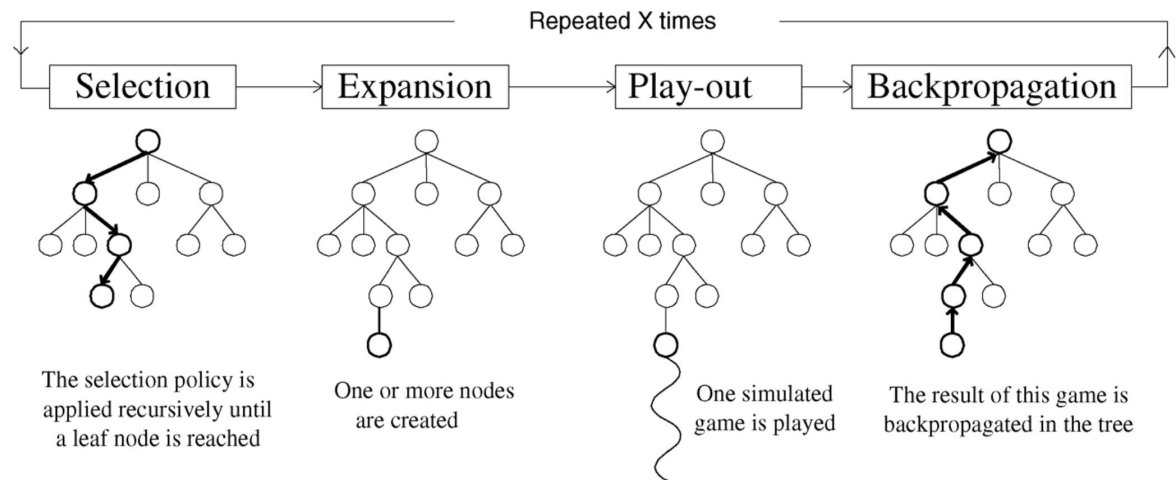


MONTE CARLO TREE SEARCH

- u MCTS builds a *statistics tree* (detailing value of nodes) that partially maps onto the entire game tree
- u Statistics tree guides the AI to focus on **most interesting nodes** in the game tree
- u Value of nodes determined by simulations

MONTÉ CARLO TREE SEARCH

- Builds and searches an asymmetric game tree to make each move



Real-Time Monte Carlo Tree Search in Ms Pac-Man. Pepels et al. IEEE Transactions on Computational Intelligence and AI in Games (

MONTE CARLO TREE SEARCH

Use results of simulations to guide growth of the game tree

Exploitation: focus on promising moves

Exploration: focus on moves where uncertainty about evaluation is high

Seems like two contradictory goals

Theory of **bandits** can help

MULTI-ARMED BANDIT PROBLEM



- u Assumptions:
 - u Choice of several arms
 - u Each arm pull is independent of other pulls
 - u Each arm has fixed, unknown average payoff
- u Which arm has the best average payoff?

CONSIDER A ROW OF THREE SLOT MACHINES

A



$P(A \text{ wins}) = 60\%$

B



$P(B \text{ wins}) = 55\%$

C



$P(C \text{ wins}) = 40\%$

Each pull of an arm is either

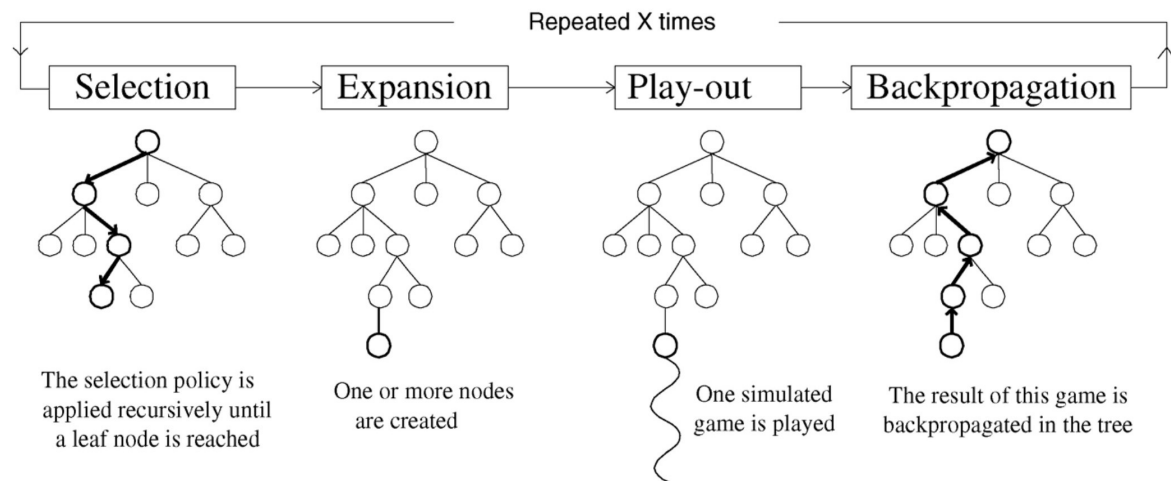
- u A win (payoff 1) or
- u A loss (payoff 0)
- u A is the best arm \rightarrow but we don't know that

EXPLORATION VS. EXPLOITATION



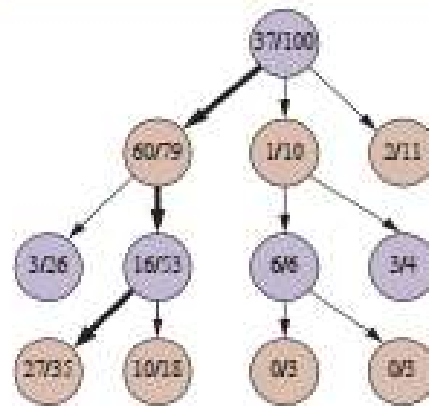
- u Want to **explore** all arms
 - u BUT, if we explore too much, may sacrifice reward we could have gotten
- u Want to **exploit** promising arms more often
 - u BUT, if we exploit too much, can get stuck with sub-optimal values
- u Want to minimize regret = loss from playing non-optimal arm
- u Need to balance between exploration and exploitation

MONTÉ CARLO TREE SEARCH



1) SELECTION

Selection policy is applied recursively until a leaf node is reached



(a) Selection

Figure 6.10 One iteration of the process of choosing a move with Monte Carlo tree search (MCTS) using the upper confidence bounds applied to trees (UCT) selection metric, shown after 100 iterations have already been done. In (a) we select moves, all the way down the tree, ending at the leaf node marked 27/35 (for 27 wins for black out of 35 playouts). In (b)

2) EXPANSION & SIMULATION

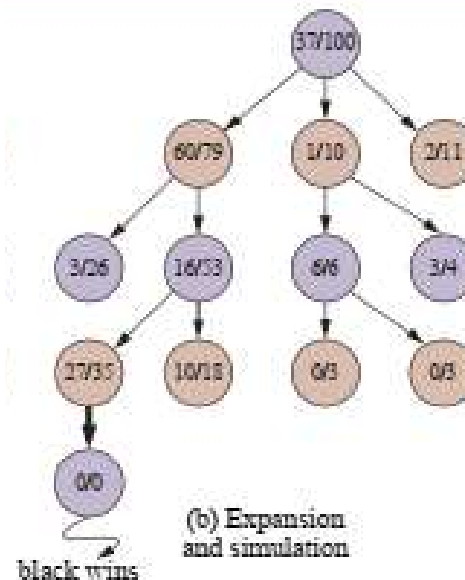


Figure 6.10 One iteration of the process of choosing a move with Monte Carlo tree search (MCTS) using the upper confidence bounds applied to trees (UCT) selection metric, shown after 100 iterations have already been done. In (a) we select moves, all the way down the tree, ending at the leaf node marked 27/35 (for 27 wins for black out of 35 playouts). In (b) we expand the selected node and do a simulation (playout), which ends in a win for black. In

Result is backpropagated up the t

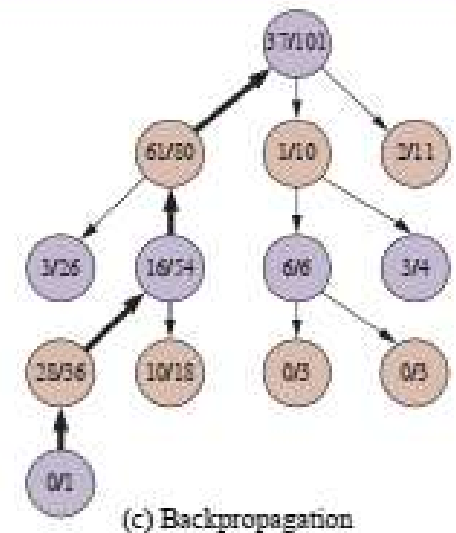


Figure 6.10 One iteration of the process of choosing a move with Monte Carlo tree search (MCTS) using the upper confidence bounds applied to trees (UCT) selection metric, shown after 100 iterations have already been done. In (a) we select moves, all the way down the tree, ending at the leaf node marked 27/35 (for 27 wins for black out of 35 playouts). In (b) we expand the selected node and do a simulation (playout), which ends in a win for black. In (c) the results of the simulation are back-propagated up the tree.

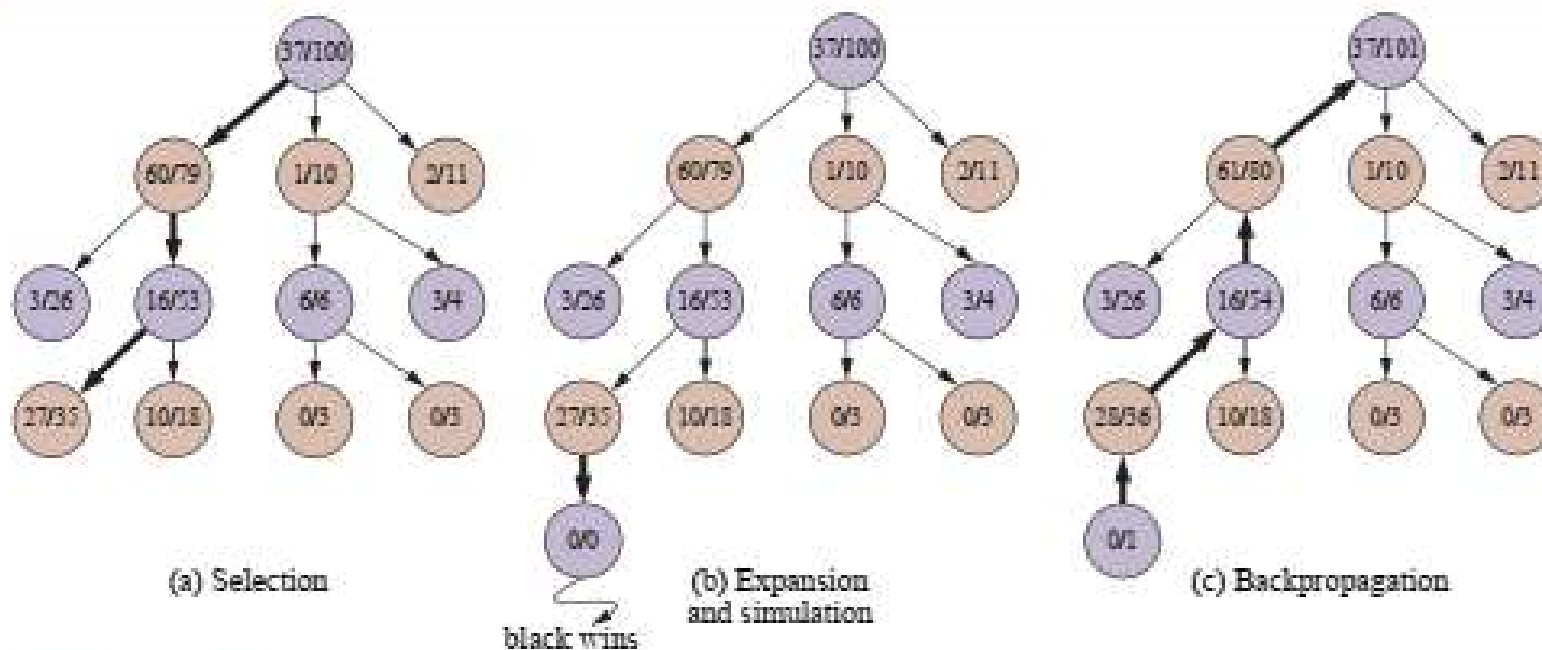
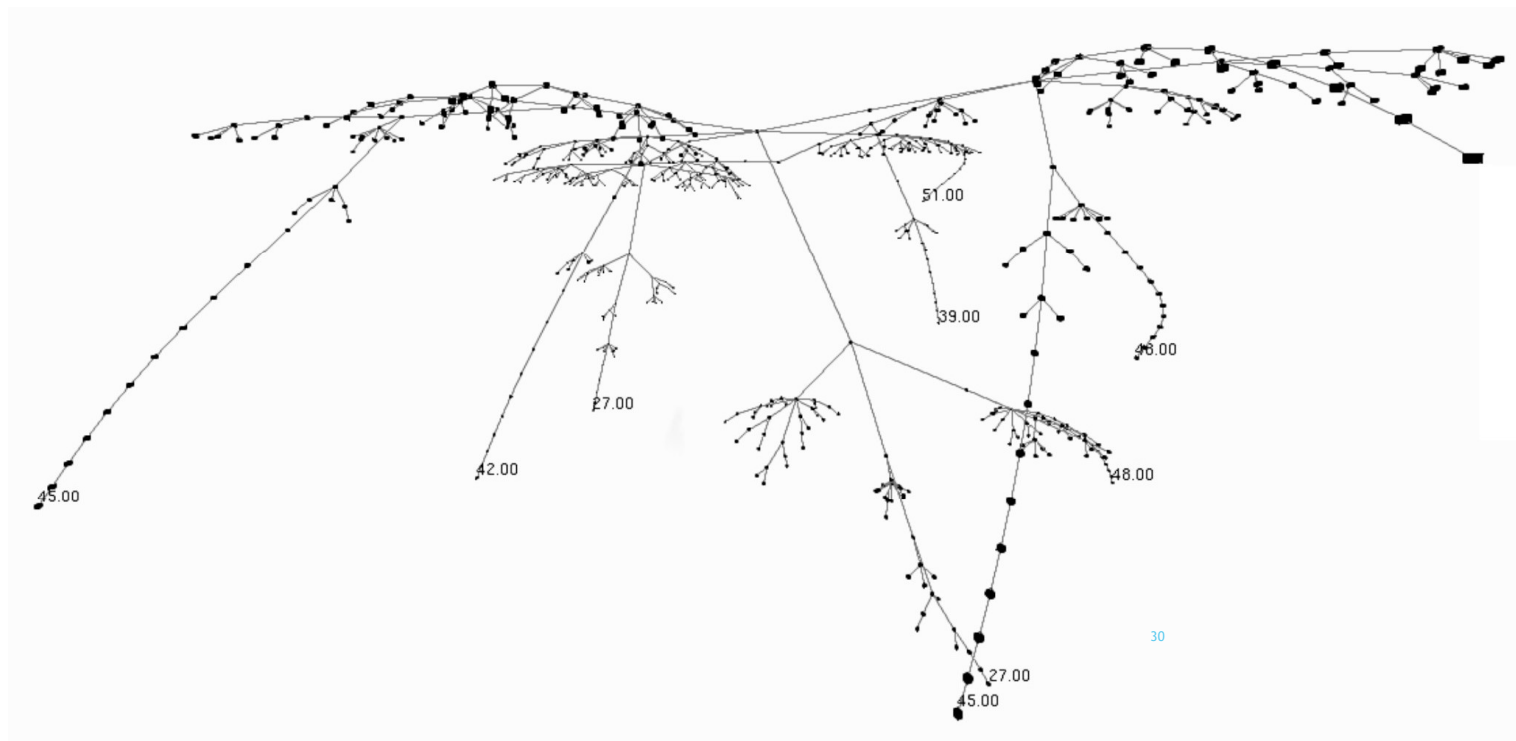


Figure 6.10 One iteration of the process of choosing a move with Monte Carlo tree search (MCTS) using the upper confidence bounds applied to trees (UCT) selection metric, shown after 100 iterations have already been done. In (a) we select moves, all the way down the tree, ending at the leaf node marked 27/35 (for 27 wins for black out of 35 playouts). In (b)

```
function MONTE-CARLO-TREE-SEARCH(state) returns an action
  tree  $\leftarrow$  NODE(state)
  while IS-TIME-REMAINING() do
    leaf  $\leftarrow$  SELECT(tree)
    child  $\leftarrow$  EXPAND(leaf)
    result  $\leftarrow$  SIMULATE(child)
    BACK-PROPAGATE(result, child)
  return the move in ACTIONS(state) whose node has highest number of playouts
```

Figure 6.11 The Monte Carlo tree search algorithm. A game tree, *tree*, is initialized, and then we repeat a cycle of SELECT / EXPAND / SIMULATE / BACK-PROPAGATE until we run out of time, and return the move that led to the node with the highest number of playouts.

SAMPLE MCTS TREE



(fig from CadiaPlayer, Bjornsson and Finsson, IEE T-CIAIG 2009)

IMPACT — STRENGTHS OF MCTS

- u Very general algorithm for decision making
- u Works with very little domain-specific knowledge
 - u Need simulator of the domain
- u Can take advantage of knowledge when present
- u Anytime algorithm
 - u Can stop the algorithm and provide answer immediately, though improves answer with more time

IMPROVING SIMULATIONS

- Default roll-out policy is to make uniform random moves
- Goal is to find strong correlations between initial position and result of a simulation
- Game independent techniques
- If there is an immediate win, take it
- Last Good Reply
- Using prior knowledge

LAST GOOD REPLY

- u Last Good Reply (Drake 2009), Last Good Reply with Forgetting (Baier et al 2010)
- u Idea: after winning simulation, store (opponent move, our answer) move pairs
 - u Try same reply in future simulations
 - u Forgetting: delete move pair if it fails
- u Evaluation: worked well for Go program with simpler playout policy
 - u Trouble reproducing success with stronger Go programs
 - u Simple form of adaptive simulations

USING PRIOR KNOWLEDGE

- u (Silver 2009) machine-learned 3x3 pattern values
- u Mogo and Fuego: hand-crafted features
 - u Weights and interaction weights trained by Latent Feature Ranking
- u AlphaGo
 - u Neural networks trained over human expert games

Learn better knowledge

- u E.g. patterns, features of a domain