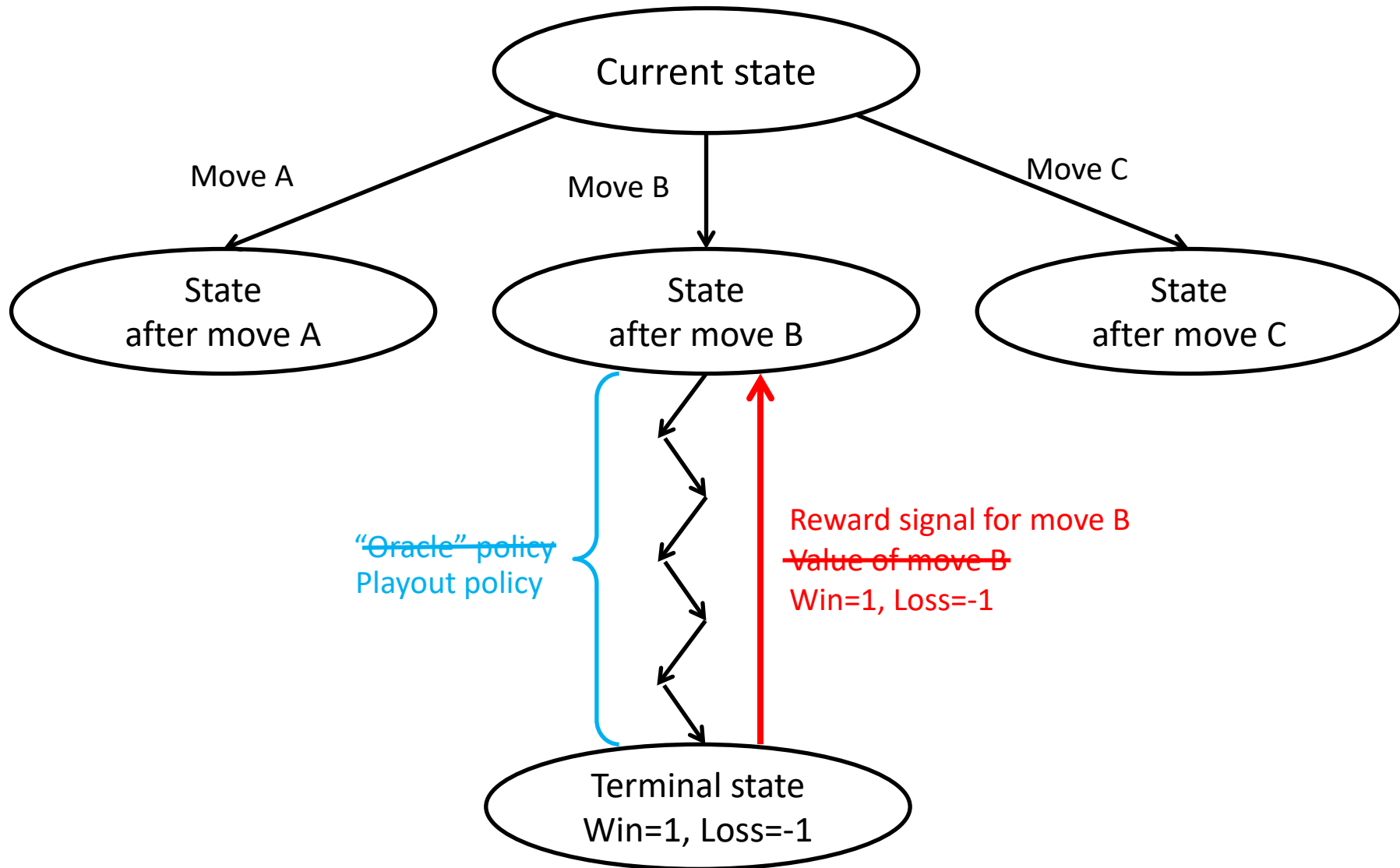
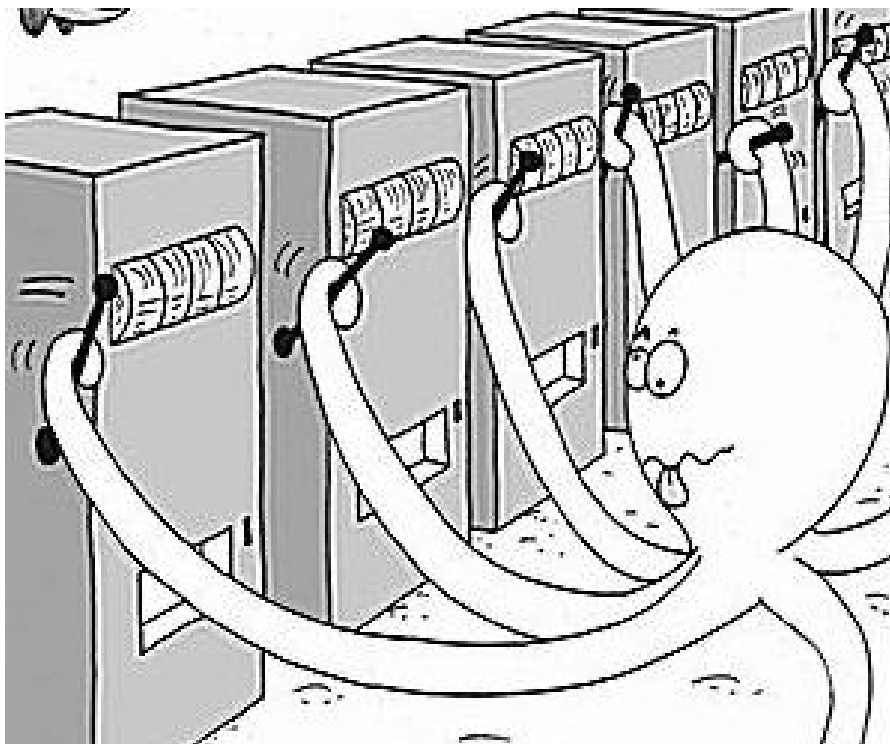


Monte Carlo Tree Search (MCTS)

The perfect 1-ply search



The multi-armed bandit problem



At each step pull one arm

Noisy/random reward signal

In order to:

- Minimise regret

- Maximise expected return
(Find the best arm)

Upper Confidence Bound (UCB1)

- Balance **exploitation** with **exploration**
- Select the arm that maximises

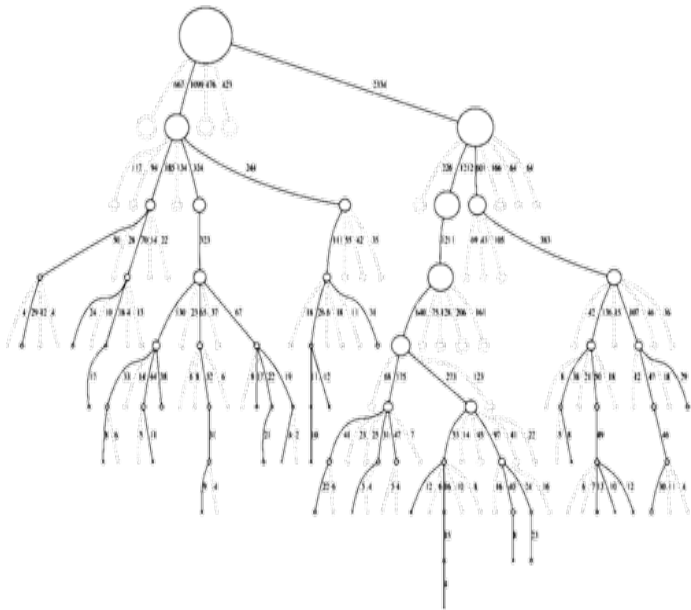
The diagram illustrates the UCB1 formula:
$$\underbrace{\frac{V_a}{n_a}}_{\text{Exploitation}} + \underbrace{k \sqrt{\frac{\log n}{n_a}}}_{\text{Exploration}}$$
 Annotations include:

- An arrow from "Total reward from this arm so far" points to V_a .
- An arrow from "Total number of trials so far" points to $\log n$.
- An arrow from "Number of times this arm has been selected so far" points to n_a in the denominator of the square root term.
- An arrow from "Tuning constant" points to k .
- Brackets below the formula label the first term as "Exploitation" and the second term as "Exploration".

UCB1 demo

<http://orangehelicopter.com/academic/bandits.html?ucb>

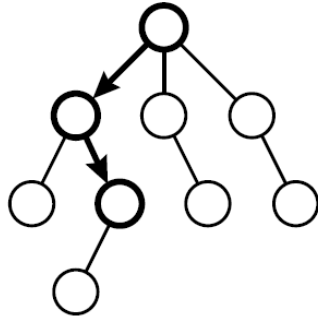
Monte Carlo Tree Search (MCTS)



- Iteratively build a partial search tree, one node per iteration
- Monte Carlo evaluation (random playouts) for nonterminal states
- **Asymmetric**
 - balance exploitation with exploration
- **Anytime**
 - can do as many (or as few) iterations as we want
- **Aheuristic**
 - Monte Carlo evaluation requires only a forward model

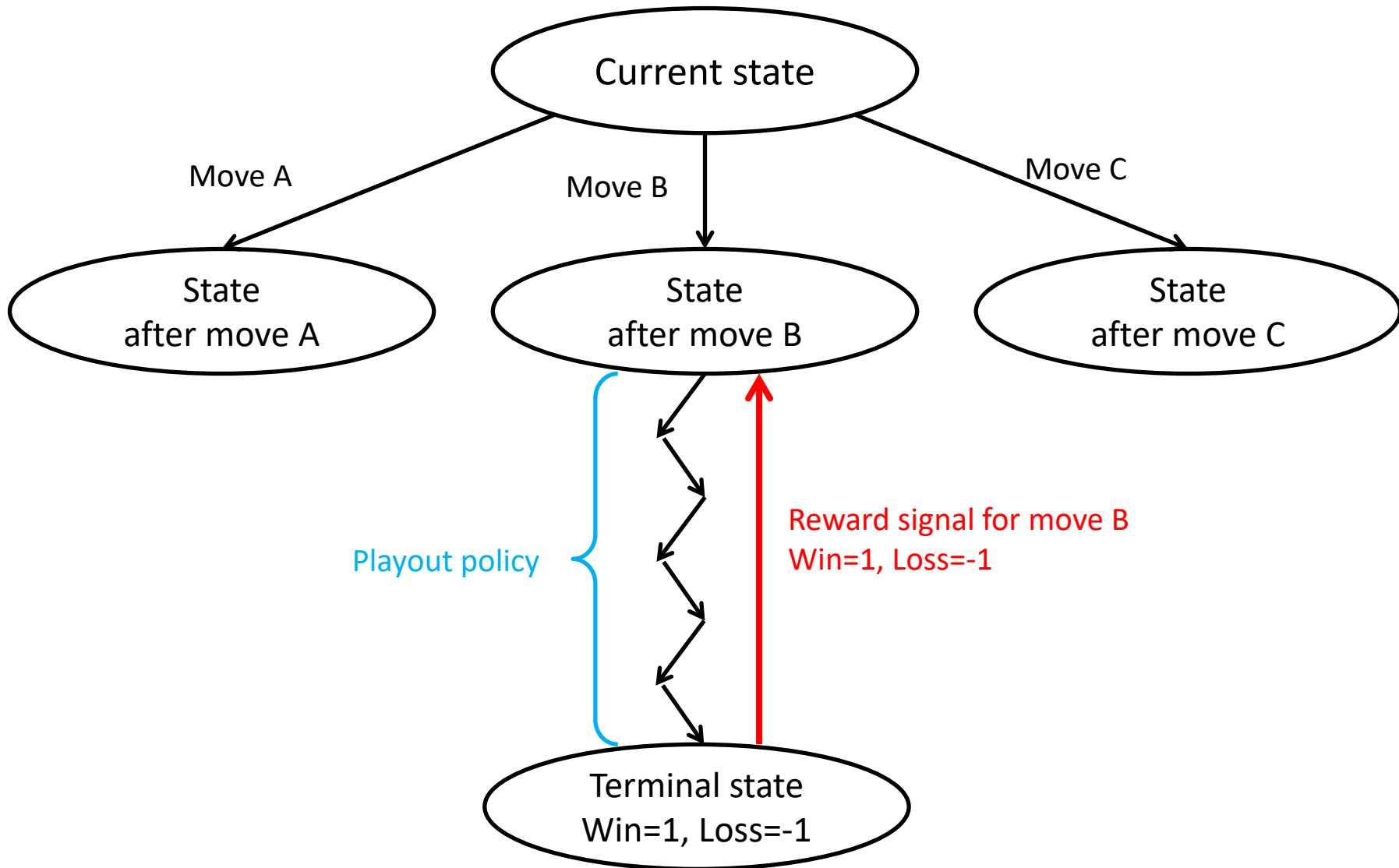
Monte Carlo Tree Search (MCTS)

Selection

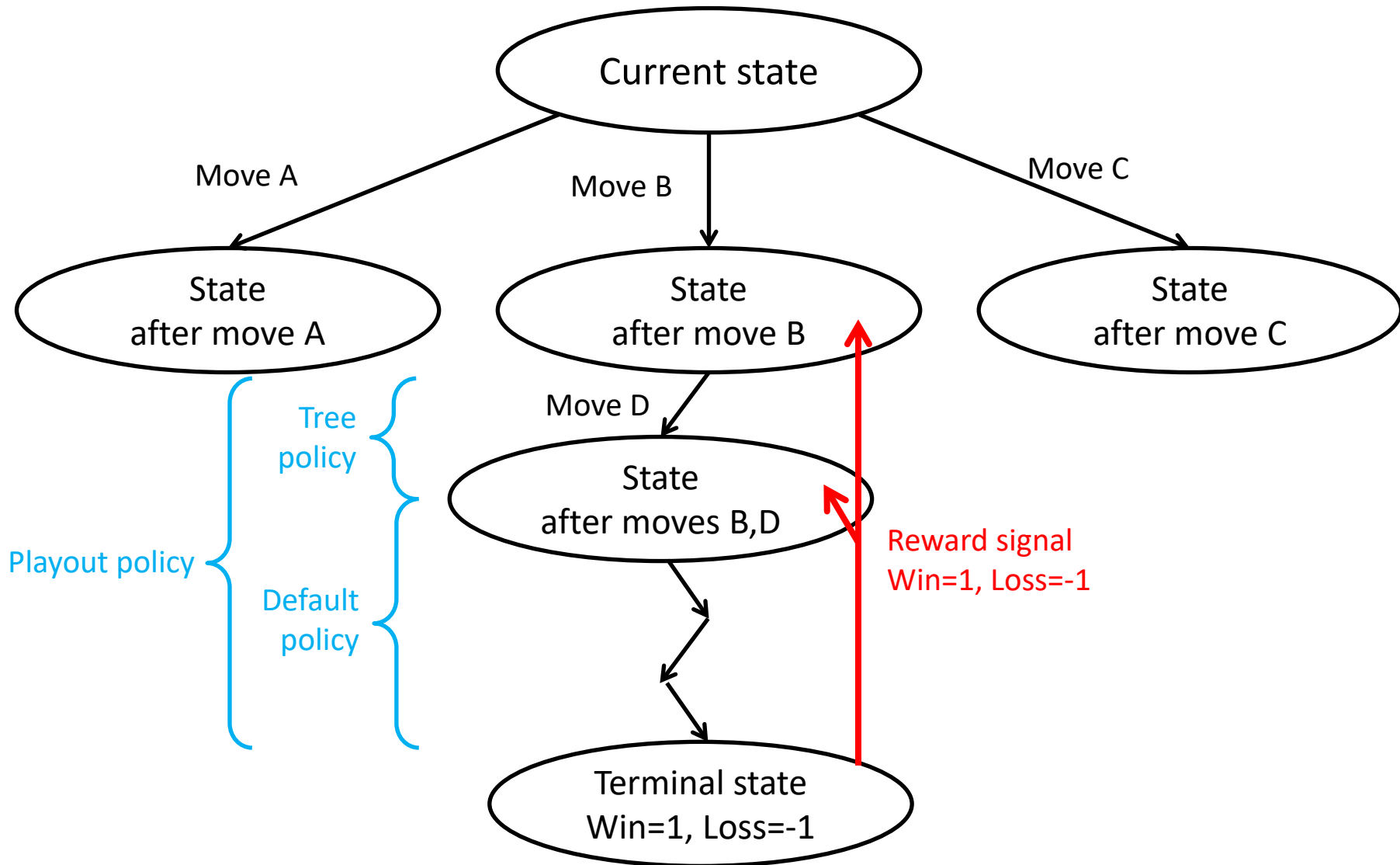


Choose a node
with unexpanded
children

Game Decisions



Monte Carlo Tree Search (MCTS)



Upper Confidence bound for Trees (UCT)

- Tree policy: UCB1 $\frac{V_a}{n_a} + k \sqrt{\frac{\log n}{n_a}}$
- Default policy: uniform random

Demo



MCTS for a real-time game

Edward J. Powley, Daniel Whitehouse and Peter I. Cowling. Monte Carlo Tree Search with macro-actions and heuristic route planning for the Physical Travelling Salesman Problem. Proceedings of IEEE Conference on Computational Intelligence in Games (CIG), 234–241, 2012.

Diego Perez, Edward J. Powley, Daniel Whitehouse, Philipp Rohlfshagen, Spyridon Samothrakis, Peter I. Cowling and Simon M. Lucas. Solving the Physical Travelling Salesman Problem: tree search and macro-actions. IEEE Transactions on Computational Intelligence and AI in Games, 6(1):31–45, 2014.

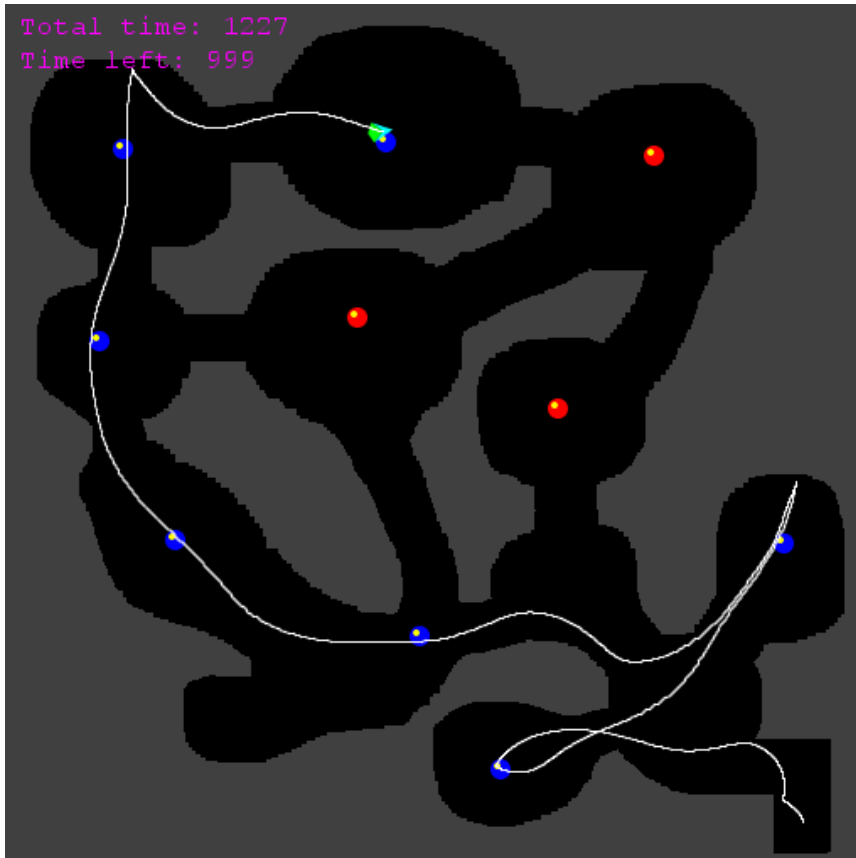
Travelling Salesman Problem (TSP)

- Classic problem in Computer Science
 - We have a **graph**
 - From starting node S , find the shortest possible path that visits every node exactly once and returns to S
 - Many real-world applications
 - Transport and logistics
 - Manufacturing
 - Playing Pac-Man
 - Pub crawls
- (<http://www.math.uwaterloo.ca/tsp/pubs/>)

Solving TSP

- TSP is **NP-complete**
- Assuming $P \neq NP$, there is no polynomial time algorithm for solving TSP perfectly
 - I.e. all algorithms scale horribly as the graph gets large
- However there are many good **heuristics** and **approximate algorithms**
- In this work we used **multiple fragment** and **3-opt**

Physical Travelling Salesman Problem (PTSP)



- Steer a spaceship to collect all **waypoints**
 - *Asteroids*-like controls
 - Newtonian physics
- Map is unknown in advance
- Controller has **a few seconds** of initialisation time, and then must make an input to the ship every **40ms**

Demo

Our PTSP controller in action

Challenges for tree search in real-time domains

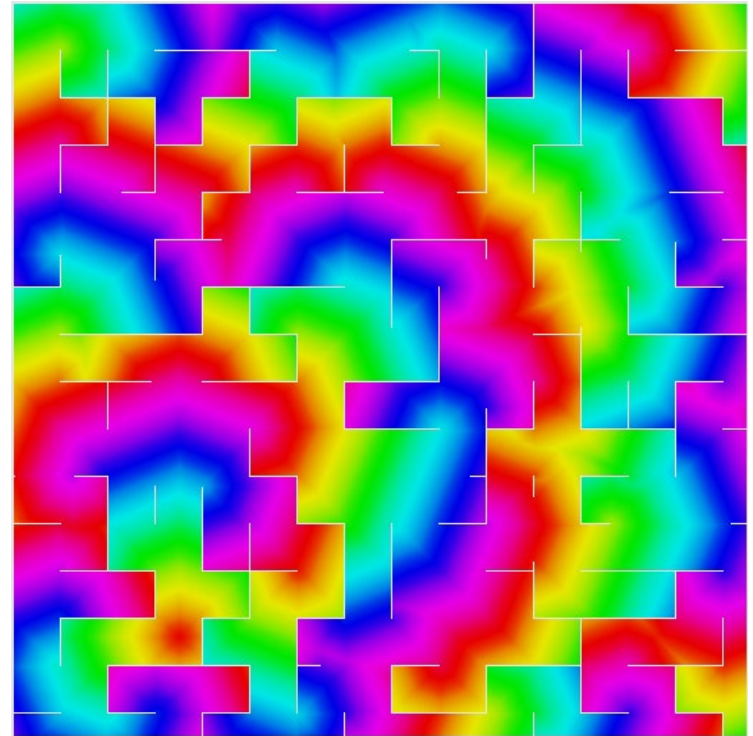
- Many more **decisions per game** than most turn-based games
- Hence state space is enormous even if branching factor is small
 - PTSP has of the order **10^{1556}** states
 - 19×19 Go has of the order **10^{171}** states
- **Time budget** is restricted (milliseconds per decision)

Key features of our PTSP controller

- **Hierarchical structure**
 - Higher-level **route planner** chooses the waypoint order
 - TSP solver with heuristics for avoiding sharp turns
 - Lower-level **steering controller** executes the route
 - MCTS with macro-actions and heuristic evaluation
- **Depth limiting and heuristic evaluation**
 - Guides the steering controller along the route
- **Macro-actions**
 - Vastly reduce the state-action space for steering

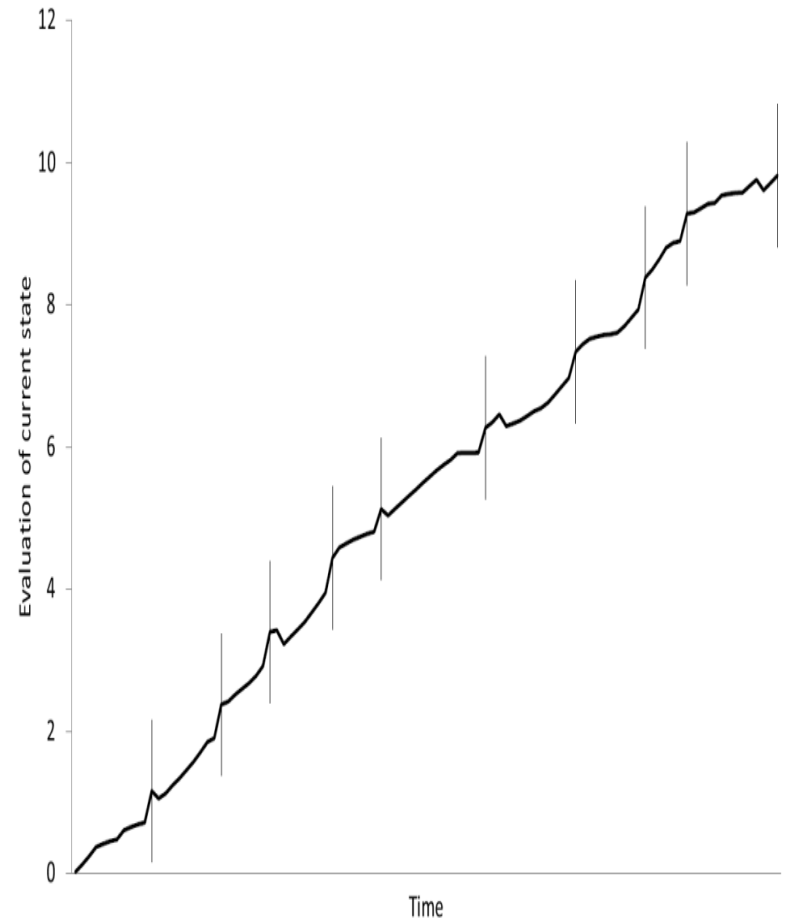
Heuristic evaluation

- Based on “A* path”
distance to next
waypoint
- Precomputed using a
flood-fill algorithm



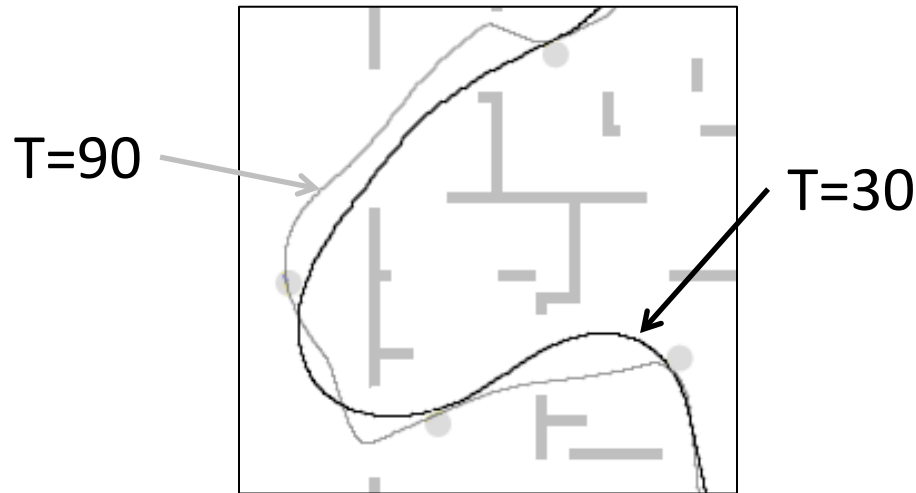
Heuristic evaluation

- Guides the ship towards the **current waypoint**
- **Sharp increase in score** for collecting the waypoint...
- ... At which point the distance to the **next waypoint** takes over



Macro-actions

- Choosing a different action every time step is an **unnecessary level of granularity**
- Instead, choose an action to be executed **for the next T time steps**



Benefits of macro-actions

- **Vastly** reduce size of state space
 - Original state space size: $\approx 10^{1556}$
 - Macro-action state space size ($T=15$): $\approx 10^{103}$
 - A reduction of **≈ 1453 orders of magnitude**
 - Size of macro-action space is comparable to 9×9 Go
- See **T times further into the future** for the same tree depth
- Have **T times longer** to make each decision
 - MCTS is not only anytime but also interruptible

Conclusion

- MCTS is a powerful **general-purpose** AI technique
 - Asymmetric, Anytime, Aheuristic
- MCTS has proven successful in several **challenging classes of games**
 - Games of imperfect information
 - Commercial mobile games
 - Real-time games
- It shows promise in **many other games and non-game applications**

Further reading

C Browne et al.

A Survey of Monte Carlo Tree Search Methods.

IEEE Transactions on
Computational Intelligence and AI
in Games, 4(1):1-43, 2012.

A Survey of Monte Carlo Tree Search Methods

Cameron B. Browne, *Member, IEEE*, Edward Powley, *Member, IEEE*,

Daniel Whitehouse, *Graduate Student Member, IEEE*, Simon M. Lucas, *Senior Member, IEEE*,

Peter I. Cowling, *Member, IEEE*, Philipp Rohlfschagen, *Member, IEEE*, Stephen Tavenor, Diego Perez,

Spyridon Samothrakis, *Graduate Student Member, IEEE*, and Simon Colton

Abstract—Monte Carlo tree search (MCTS) is a recently proposed search method that combines the precision of tree search with the generality of random sampling. It has received considerable interest due to its spectacular success in the difficult problem of computer Go, but has also proved beneficial in a range of other domains. This paper is a survey of the literature to date, intended to provide a snapshot of the state of the art after the first five years of MCTS research. We outline the core algorithm's derivation, impart some structure on the many variations and enhancements that have been proposed, and summarize the results from the key game and nongame domains to which MCTS methods have been applied. A number of open research questions indicate that the field is ripe for future work.

Index Terms—Artificial intelligence (AI), bandit-based methods, computer Go, game search, Monte Carlo tree search (MCTS), upper confidence bounds (UCB), upper confidence bounds for trees (UCT).

I. INTRODUCTION

MONTE CARLO TREE SEARCH (MCTS) is a method for finding optimal decisions in a given domain by taking random samples in the decision space and building a search tree according to the results. It has already had a profound impact on artificial intelligence (AI) approaches for domains that can be represented as trees of sequential decisions, particularly games and planning problems.

In the five years since MCTS was first described, it has become the focus of much AI research. Spurred on by some prolific achievements in the challenging task of computer Go, researchers are now in the process of attaining a better understanding of when and why MCTS succeeds and fails, and of extending and refining the basic algorithm. These developments are greatly increasing the range of games and other decision applications for which MCTS is a tool of choice, and pushing its

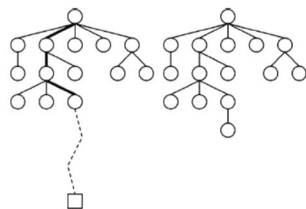


Fig. 1. The basic MCTS process [17].

performance to ever higher levels. MCTS has many attractions: it is a statistical anytime algorithm for which more computing power generally leads to better performance. It can be used with little or no domain knowledge, and has succeeded on difficult problems where other techniques have failed. Here we survey the range of published work on MCTS, to provide the reader with the tools to solve new problems using MCTS and to investigate this powerful approach to searching trees and directed graphs.

A. Overview

The basic MCTS process is conceptually very simple, as shown in Fig. 1 (from [17]). A tree¹ is built in an incremental and asymmetric manner. For each iteration of the algorithm, a tree policy is used to find the most urgent node of the current tree. The tree policy attempts to balance considerations of exploration (look in areas that have not been well sampled yet) and exploitation (look in areas which appear to be promising). A simulation² is then run from the selected node and the search tree updated according to the result. This involves the addition of a child node corresponding to the action taken from the selected node, and an update of the statistics of its ancestors. Moves are made during this simulation according to some default policy, which in the simplest case is to make uniform random moves. A great benefit of MCTS is that the values of intermediate states do not have to be evaluated, as for depth-limited minimax search, which greatly reduces the amount of domain knowledge required. Only the value of the terminal state at the end of each simulation is required.

¹Typically a game tree.

²A random or statistically biased sequence of actions applied to the given state until a terminal condition is reached.

Manuscript received October 22, 2011; revised January 12, 2012; accepted January 30, 2012. Date of publication February 03, 2012; date of current version March 13, 2012. This work was supported by the Engineering and Physical Sciences Research Council (EPSRC) under Grants EP/I001964/1, EP/H048588/1, and EP/H049061/1, as part of the collaborative research project UCT for Games and Beyond being undertaken by Imperial College London, the University of Essex, and the University of Bradford.

C. Browne, S. Tavenor, and S. Colton are with the Department of Computing, Imperial College London, London SW7 2RH, U.K. (e-mail: camb@doc.ic.ac.uk; sc1110@doc.ic.ac.uk; sgc@doc.ic.ac.uk).

S. M. Lucas, P. Rohlfschagen, D. Perez, and S. Samothrakis are with the School of Computer Science and Electronic Engineering, University of Essex, Colchester, Essex, CO4 3SQ, U.K. (e-mail: sml@essex.ac.uk; prohl@essex.ac.uk; dperez@essex.ac.uk; ssamot@essex.ac.uk).

E. Powley, D. Whitehouse, and P. I. Cowling are with the School of Computing, Informatics and Media, University of Bradford, Bradford, West Yorkshire BD7 1DP, U.K. (e-mail: e.powley@bradford.ac.uk; d.whitehouse1@bradford.ac.uk; p.i.cowling@bradford.ac.uk).

Digital Object Identifier 10.1109/TCAIG.2012.2186810