

Multi-Objective Monte Carlo Tree Search for Real-Time Games

Diego Perez, *Student Member, IEEE*, Sanaz Mostaghim, Spyridon Samothrakis, *Student Member, IEEE*, Simon M. Lucas, *Senior Member, IEEE*

Abstract—Abstract...

I. INTRODUCTION

Here it goes, the introduction.

II. MONTE CARLO TREE SEARCH

Monte Carlo Tree Search (MCTS) is a tree search algorithm that was originally applied to board games, concretely to the two-players game of Go. This game is played in a square grid board, with a size of 19×19 in the original game, and 9×9 in its reduced version. The game is played in turns, and the objective is to surround the opponent's stones by placing stones in any available position in the board. Due to the very large branching factor of Go, this game is considered the drosophila of Game AI, and MCTS players have reached professional level play in the reduced board size version [1]. After its success in Go, MCTS has been used extensively by many researchers in this and different domains. An extensive survey of MCTS methods, variations and applications, has been written by Browne et al. [2].

MCTS is considered to be an *anytime* algorithm, as it is able to provide a valid next move to choose at any moment in time. This is true independently from how many iterations the algorithm was able to make (although, in general, more iterations usually produce better results). This differs from other algorithms (such as A^*) that normally provide the next ply only after they have finished. This makes MCTS a suitable candidate for real-time domains, where the decision time budget is limited, affecting the number of iterations that can be performed.

MCTS is an algorithm that builds a tree in memory. Each node in the tree maintains statistics that indicate how often a move is played from a given state ($N(s, a)$), how many times each move is played from there ($N(s)$) and the average reward ($Q(s, a)$) obtained after applying move a in state s . The tree is built iteratively by simulating actions in the game, making move choices based on the statistics store in the nodes.

Each iteration of MCTS can be divided into several steps [3]: *Tree selection*, *Expansion*, *Monte Carlo simulation* and *Back-propagation* (all summarized in Figure 1). When the algorithm starts, the tree is formed only by the root node, which holds the current state of the game. During the *selection*

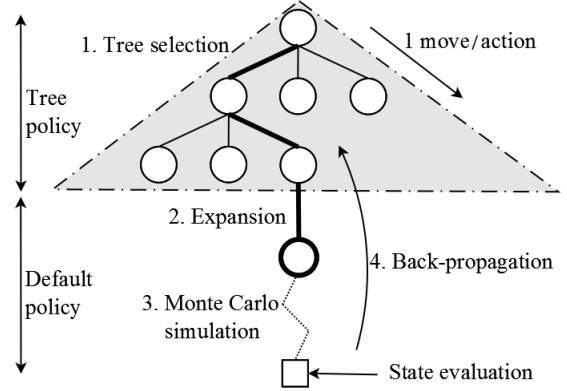


Fig. 1: MCTS algorithm steps.

step, the tree is navigated from the root until a maximum depth or the end of the game has been reached.

In every one of this action decisions, MCTS balances between exploitation and exploration. In other words, this chooses between taking an action that leads to states with the best outcome found so far, and performing a move to go to less explored game states, respectively. In order to achieve this, MCTS uses Upper Confidence Bound (UCB1, see Equation 1) as a *Tree Policy*.

$$a^* = \arg \max_{a \in A(s)} \left\{ Q(s, a) + C \sqrt{\frac{\ln N(s)}{N(s, a)}} \right\} \quad (1)$$

The balance between exploration and exploitation is achieved by setting the value of C . Higher values of C weight more the second term of the UCB1 Equation 1, giving preference to those actions that have been explored less, at the expense of taking actions with the highest average reward $Q(s, a)$. A commonly used value is $\sqrt{2}$, as it balances both facets of the search when the rewards are normalized between 0 and 1. It is worth noting that MCTS, when combined with UCB1 reaches asymptotically logarithmic regret [4].

If, during the *tree selection* phase, a node has less children than the available number of actions from a given position, a new node is added as a child of the current one (*expansion* phase) and the *simulation* step starts. At this point, MCTS executes a Monte Carlo simulation (or roll-out; *default policy*) from the expanded node. This is performed by choosing random (either uniformly random, or biased) actions until the game end or a pre-defined depth is reached, where the state of the game is evaluated.

Finally, during the *back-propagation* step, the statistics $N(s)$, $N(s, a)$ and $Q(s, a)$ are updated for each node visited,

using the reward obtained in the evaluation of the state. These steps are executed in a loop until a termination criteria is met (such as number of iterations).

MCTS has been employed extensively in real-time games in the literature. A clear example of this is the popular real-time game *Ms. PacMan*. The objective of this game is to control Ms. PacMan to clear the maze by eating all pills, without being captured by the ghosts. An important feature of this game is that it is *open-ended*, as an end game situation is, most of the time, far ahead in the future and can not be devised by the algorithm during its iterations. The consequence of this is that MCTS, in its vanilla form, it is not able to know if a given ply will lead to a win or a loss game end state. Robles et al [5] solved this problem by including hand-coded heuristics that guided MCTS simulations towards more promising portions of the search space. Other authors also included domain knowledge to bias the search in MCTS, such as in [6], [7].

MCTS has also been applied to single-player games, like SameGame [8], where the player's goal is to destroy contiguous tiles of the same colour, distributed in a rectangular grid. Another use of MCTS is in the popular puzzle Morpion Solitaire [9], a connection game where the goal is to link nodes of a graph with straight lines that must contain at least five vertices. Finally, the PTSP has also been addressed with MCTS, both in the single-objective [10], [11] and the multi-objective versions [12]. These papers describe the entries that won both editions of the PTSP Competition.

It is worthwhile mentioning that in most cases found in the literature, MCTS techniques have been used with some kind of heuristic that guides the Monte Carlo simulations or the tree selection policy. In the algorithm proposed in this paper, simulations are purely random, as the objective is to compare the search abilities of the different algorithms. The intention is therefore to keep the heuristics to a minimum, and the existing pieces of domain knowledge are shared by all the algorithms presented (as in the case of the score function for MO-PTSP, described later).

III. MULTI-OBJECTIVE OPTIMIZATION

A multi-objective optimization problem (MOOP) represents a scenario where two or more objective functions are to be optimized (either maximized or minimized). The general form of a MOOP is formally described as a maximization function $F_m(x)$, that transforms points in the decision space (X) to points in the solution space (F). The elements of the decision space are vectors of n variables of the form $x = (x_1, x_2, \dots, x_n)$, while elements in the solution space are vectors with a dimension m : $F_m(x) = (f_1(x), f_2(x), \dots, f_m(x))$. Therefore, each solution provides m different scores (or rewards, or fitness) that are meant to be optimized. Without loss of generality, it is assumed from now on that all objectives must be maximized.

It is said that a solution $F_m(x)$ *dominates* another solution $F_m(y)$ if:

- 1) $F_m(x)$ is not worse than $F_m(y)$ in all objectives for all $i = 1, 2, \dots, m$.

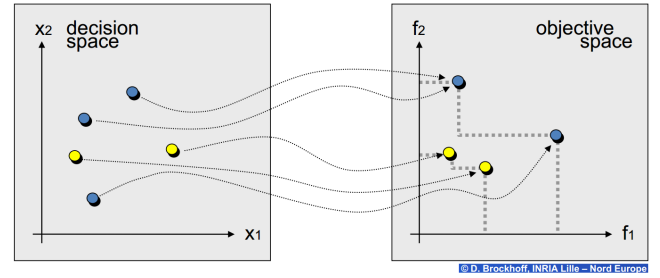


Fig. 2: Decision and Solution spaces in a MOOP with two variables (x_1 and x_2) and two objectives (f_1 and f_2). In the objective space, yellow dots are non-optimal objective vectors, while blue dots form the Pareto-optimal front.

- 2) At least one objective of $F_m(x)$ is better than its analogous counterpart in $F_m(y)$.

When this two conditions apply, it is said that $F_m(x) \preceq F_m(y)$ ($F_m(x)$ dominates $F_m(y)$), and $F_m(x)$ is non-dominated by $F_m(y)$. The *dominance* condition provides a partial ordering between points in the solution space: if $F_m(x) \preceq F_m(y)$, then $F_m(x)$ is considered to be better than $F_m(y)$.

However, there are some cases where it cannot be said that $F_m(x) \preceq F_m(y)$ or $F_m(y) \preceq F_m(x)$. This situation occurs when one of the objectives is better in $F_m(x)$ but a different objective is better in $F_m(y)$ (for instance, when $F_1(x) < F_1(y)$ but $F_2(x) > F_2(y)$). In this case, it is said that these solutions are non-dominated with respect to each other. Solutions that are not dominated by each other are grouped in a *non-dominated set*. Given a non-dominated set P , it is said that P is the *optimal Pareto front* if there is no other solution in the solution space that dominates any member of P . The relation between decision and objective space, dominance and Pareto fronts is depicted in Figure 2.

As many Pareto front can be formed by several points in the solution space, it is important to devise a mechanism to assess the quality of a given front. A possibility is to use the Hypervolume Indicator (HV): given a Pareto front P , $HV(P)$ is defined as the volume of the objective space dominated by P . More formally, $HV(P) = \mu(\{x \in \mathbb{R}^d : \exists r \in P \text{ s.t. } r \preceq x\})$, where μ is the de Lebesgue measure on \mathbb{R}^d . If the objectives are to be maximized, the higher the $HV(P)$, the better the front calculated. Figure 3 shows an example of $HV(P)$ where the objective dimension space is 2.

For more extensive descriptions, definitions, properties and multi-objective optimization in general, the reader can consult the work by K. Deb [13].

Many different algorithms have been proposed to tackle multi-objective optimization problems in the literature. One of the most widely known and used methods is the weighted-sum approach. The procedure consists of giving a weight to each one of the objective and produce a single result as the linear combination of objectives and weights. By varying the weights provided, it is possible to converge to different solutions of the optimal Pareto front, if this is convex. However, K. Deb [13] explains how linear scalarization approaches fail in

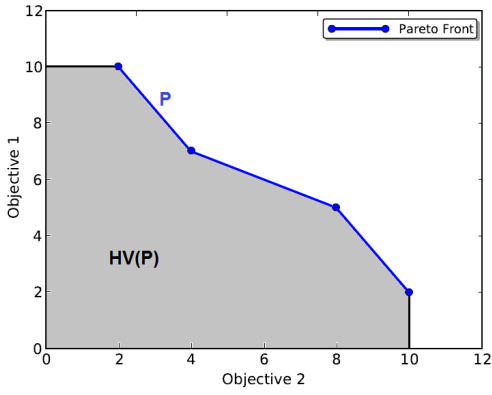


Fig. 3: $HV(P)$ of a given pareto front P .

Algorithm 1 NSGA-2 Algorithm.

```

1: function NSGA-2
2:    $P = \text{NewRandomPopulation}$ 
3:   while Termination criteria not met do
4:      $R = P \cup Q$ 
5:      $F = \text{FASTNONDOMINATEDSORT}(R)$ 
6:     while  $|P| < N$  do
7:        $\text{CROWDINGDISTANCEASSIGNMENT}(F_i)$ 
8:        $P = P \cup F_i$ 
9:      $\text{SORT}(P)$ 
10:     $P = P[0 : N]$ 
11:     $Q = \text{breed}(P)$ 

```

those scenarios where the optimal Pareto front is non-convex.

A popular choice for multi-objective optimization problems are evolutionary multi-objective optimization (EMOA) algorithms [14], [15]. One of the most well known algorithms in the literature is the Non-dominated Sorting Evolutionary Algorithm 2 (NSGA-2), which pseudocode is shown in Algorithm 1. As in any evolutionary algorithm, NSGA-2 evolves a set of individuals or solutions to the problem, with the difference that here they are ranked according to dominance criteria and crowding distance (distances between members of the pareto fronts). A full description of the algorithm can be found in [16].

The three main pillars of the NSGA-2 algorithm are:

- A *fast non-dominated sorting* algorithm, that ranks the individuals of the population and groups them in Pareto fronts.
- A *crowding distance*, assigned to each one of the individuals, that measures how close it is to its neighbours. The selection genetic operator chooses individuals based on the ranks of the individuals and their crowding distance.
- *Elitism*, implemented so the algorithm automatically promotes the best N individuals to the next generation.

A more recent approach, developed by Q. Zhang, is the Multi-Objective Evolutionary Algorithm based on Decomposition (MOEA/D) [17], that decomposes the problem into several single optimization sub-problems and an evolutionary algorithm optimizes them all simultaneously. Information is shared between neighbouring sub-problems in order to guide

evolution. The authors show that MOEA/D performs similarly than, and sometimes even outperforms, NSGA-II in the scenarios tested.

Reinforcement Learning (RL) algorithms have also tackled Multi-objective optimization in some scenarios. RL [18] is a broad field in Machine Learning that studies real-time planning and control situations where an agent has to find out the actions (or sequences of actions) that should be applied in order to maximize the reward from the environment.

An RL problem can be defined as a tuple (S, A, T, R, π) . S is the set of possible states in the problem (or game), and s_0 is the initial state. A is the set of available actions the agent can make at any given time, and the transition model $T(s_i, a_i, s_{i+1})$ determines the probability of reaching the state s_{i+1} when action a_i is applied in state s_i . The reward function $R(s_i)$ provides a single value (*reward*) that the agent must optimize, representing the desirability of the state s_i reached. Finally, a decision policy $\pi(s_i) = a_i$ determines which actions a_i must be chosen from each state $s_i \in S$. One of the most important challenges in RL, as shown in Section II, is the trade-off between exploration and exploitation. The decision policy can choose between following actions that provided good rewards in the past and exploring new parts of the search space by selecting new actions.

Multi-objective Reinforcement Learning (MORL) [19] changes this definition by using instead a vector $R = r_0, r_1, \dots, r_n$ as rewards of the problem. Thus, MORL problems differ from RL in having more than one objective objectives (n) that must be maximized. If the objectives are independent or they do not oppose each other, scalarization technique approaches, as described above, could be suitable to tackle the problem. Essentially, this would mean to use a conventional RL algorithm on a single objective obtained by a weighted-sum of the multiple rewards. However, this is not always the case, as it is usual that the objectives are in conflict and the policy (π) must balance among them.

According to how π approaches this problem, Vamplew et al. [19] proposed the following distinction for MORL: *single-policy* algorithms are those that provide a preference order in the objectives available (given by the user or by the nature of the problem). An example of this type of algorithm can be found at [20], where the authors introduce an order of preference in the objectives treated and constraint the value of the rewards desired. Scalarization approaches would also fit in this category, as the work performed by S. Natarajan et al. [21].

The second type of algorithms, *multiple-policy*, target to approximate the optimal Pareto front of the problem. An example of this type of algorithm is the one given by L. Barrett [22], who propose the Convex Hull Iteration Algorithm. This algorithm provides the optimal policy for any linear preference function, by learning all policies that define the convex hull of the Pareto front.

IV. MULTI-OBJECTIVE MONTE CARLO TREE SEARCH

Adapting MCTS into Multi-Objective Monte Carlo Tree Search (MO-MCTS) requires the obvious modification of

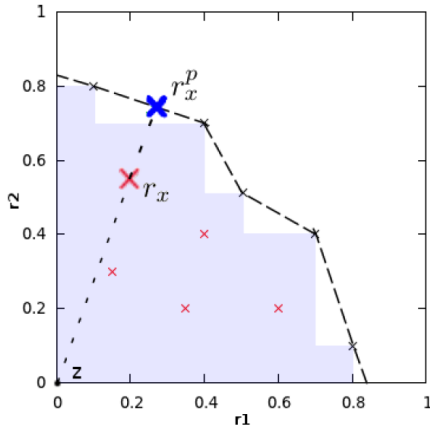


Fig. 4: r_x^p is the projection of the r_x value on the piecewise linear surface (discontinuous line). The shadowed surface represents $HV(P)$. From [23].

dealing with multiple rewards instead of just one. As these are collected at the end of a Monte Carlo simulation, the reward value r now becomes a vector $R = r_0, r_1, \dots, r_n$, where n is the number of objectives to optimize. Derived from this change, the average value $Q(s, a)$ becomes a vector that stores the average reward of each objective. Note that the other statistics ($N(s, a)$ and $N(s)$) do not need to change, as this are just node and action counters. The important question to answer next is how to adapt the vector $Q(s, a)$ to use it in the UCB1 formula (Equation 1).

An initial attempt on Multi-Objective MCTS was addressed by W. Wang and Michele Sebag [23], [24]. In their work, the authors employ a mechanism, based on the HV calculation, to replace the UCB1 equation. The algorithm keeps a Pareto archive (P) with the best solutions found in game end states. Every node in the tree defines \bar{r}_{sa} as a vector of UCB1 values, in which each $\bar{r}_{sa,i}$ is the result of calculating UCB1 for each particular objective i .

The next step is to define the value for each pair of state and action, $W(s, a)$, as in Equation 2. \bar{r}_{sa}^p is the projection of \bar{r}_{sa} into the piecewise linear surface defined by the Pareto archive P (see Figure 4). Then, $HV(P \cup \bar{r}_{sa})$ is declared as the HV of P plus the point \bar{r}_{sa} . If \bar{r}_{sa} is dominated by P , the distance between \bar{r}_{sa} and \bar{r}_{sa}^p is subtracted from the HV calculation. The tree policy selects actions based on a maximization of the value of $W(s, a)$.

$$W(s, a) = \begin{cases} HV(P \cup \bar{r}_{sa}) - dist(\bar{r}_{sa}^p, \bar{r}_{sa}) & \text{Otherwise} \\ HV(P \cup \bar{r}_{sa}) & \text{if } \bar{r}_{sa} \preceq P \end{cases} \quad (2)$$

The proposed algorithm was employed successfully in two domains: the DST and the Grid Scheduling problem, matching state of the art results in both domains, at the expense of a high computational cost.

As mentioned before, the objective of this paper is to propose an MO-MCTS algorithm that is suitable for real-time domains. Although the work discussed here is influenced by Wang's approach, some modifications need to be done in order

Algorithm 2 Pareto MO-MCTS node update.

```

1: function UPDATE( $node, \bar{r}, dominated = false$ )
2:    $node.Visits = node.Visits + 1$ 
3:    $node.\bar{R} = node.\bar{R} + \bar{r}$ 
4:   if ! $dominated$  then
5:     if  $node.P \preceq \bar{r}$  then
6:        $dominated = true$ 
7:     else
8:        $node.P = node.P \cup \bar{r}$ 
9:   UPDATE( $node.parent, \bar{r}, dominated$ )

```

to overcome the high computational cost involved by their approach.

In the algorithm proposed in this paper, the vector \bar{r} of rewards that was obtained at the end of an Monte Carlo simulation is back-propagated through the nodes visited in the last iteration until the root is reached. In the vanilla algorithm, each node would use this vector \bar{r} to update its own accumulated reward vector \bar{R} . Instead of doing this, each node in the MO-MCTS algorithm keeps a local Pareto front (P), updated at each iteration with the reward vector \bar{r} obtained at the end of the simulation. Algorithm 2 describes how the node statistics are updated.

In this algorithm, if \bar{r} is not dominated by the local Pareto front, it is added to the front and \bar{r} is propagated to its parent. In case \bar{r} is dominated by the node's Pareto front, the local Pareto front and there is no need to keep this propagation up the tree.

Three observations can be made about the mechanism described here:

- Each node in the tree has an estimate of the quality of the solutions reachable from there, both as an average (as in the baseline MCTS) and as the best case scenario (by keeping the non-dominated front P).
- By construction, if a reward \bar{r} is dominated by the local front of a node, it is a given that it will be dominated by the nodes above in the tree, so there is no need to update the fronts of the upper nodes, producing little cost in the efficiency of the algorithm.
- It is easy to infer, from the last point, that the Pareto front of a node cannot be worse than the front of its children (in other words, the front of a child will never dominate that of its parent). Therefore, the root of the tree contains the best non-dominated front ever found during the search.

This last detail is important for two main reasons. First of all, the algorithm allows the root to store information indicating which is the action to take in order to converge to any specific solution in the front discovered. This information can be used, when all iterations have been performed, to select the move to perform next. If weights for the different objectives are provided, this weights can be used to select the desired solution in the Pareto front of the root node, and hence select the action that leads to that point. Secondly, the root's Pareto front can be used to measure the global quality of the search using the hypervolume calculation.

Finally, the information stored at each node regarding the

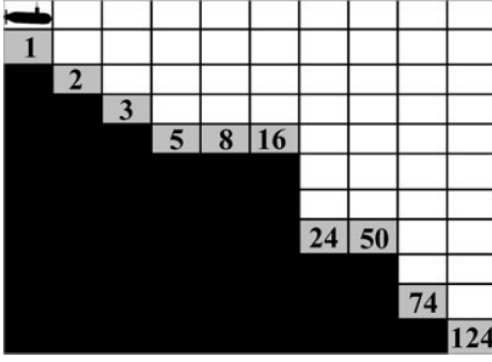


Fig. 5: Environment of the Deep Sea Treasure (from [19]): grey squares represent the treasure (with their different values) available in the map. The black cells are the sea floor and the white ones are the positions that the vessel can occupy freely. The game ends when the submarine picks one of the treasures.

local Pareto front can be used to substitute $Q(s, a)$ in the UCB1 equation. The quality of a given pair (s, a) can be obtained by measuring the HV of the Pareto front stored in the node reached from state s after applying action a . This can be defined as $Q(s, a) = HV(P)/N(s)$, and the Upper Confidence Bound equation, referred to here as *MO-UCB*, is described as in Equation 3).

$$a^* = \arg \max_{a \in A(s)} \left\{ HV(P)/N(s) + C \sqrt{\frac{\ln N(s)}{N(s, a)}} \right\} \quad (3)$$

V. BENCHMARKS

Two different are used in this research to analyze the performance of the algorithm proposed.

A. Deep Sea Treasure

The Deep Sea Treasure (DST) is a well known multi-objective problem introduced by Vamplew et al. [19]. In this single-player puzzle, the agent moves a submarine with the objective of finding a treasure located at the bottom of the ship. The world is divided into a grid of 10 rows and 11 columns, and the vessel starts at the top left board position. There are three types of cells: empty cells (or water), that the submarine can traverse; ground cells that, as the edges of the grid, cannot be traversed; and treasure cells, that provide different rewards and finish the game. Figure 5 (on the left) shows the environment of the DST.

The ship can perform four different moves: *up*, *down*, *right* and *left*. If the action applied takes the ship off the grid or into the sea floor, the vessel's position does not change. There are two objectives in this game: the number of moves performed by the ship, that must be minimized, and the value of the treasure found, that should be maximized. As can be seen in Figure 5, the most valuable treasures are at a greater distance from the initial position, making these objectives orthogonal.

Additionally, the agent can only make up to 100 moves. This allows the problem to be defined as the maximization

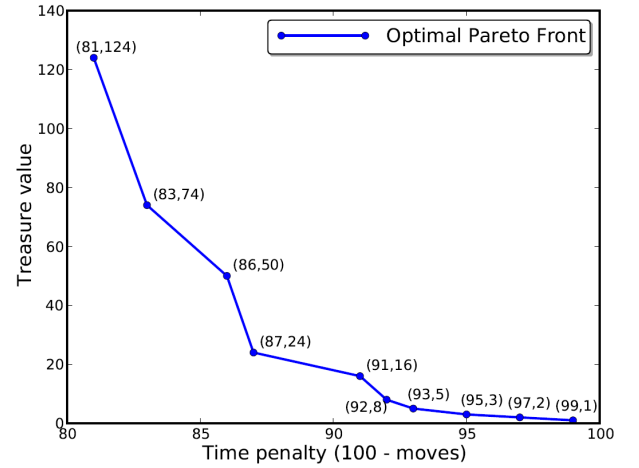


Fig. 6: Optimal Pareto Front of the Deep Sea Treasure, with both objectives to be maximized. From [23].

of two rewards: $(M, T) = (100 - moves, treasureValue)$. Should the ship perform all moves without reaching a treasure, the result would be $(0, 0)$. At each step, the score of a location with no treasure is $(-1, 0)$.

The optimal Pareto front of the DST is shown in Figure 6. There are 10 non-dominated solutions in this front, one per each treasure in the board. The front is globally concave, with local concavities at the second $(83, 74)$, fourth $(87, 24)$ and sixth $(92, 8)$ points from the left. The HV value of the optimal front is 10455.

Section III introduced the problems of linear scalarization approaches when facing non-convex optimal Pareto front. The concave shape of the DST's optimal front makes those approximations to converge to the non dominated solutions located at the edges of the Pareto front $((-19, 1), (-1, 124))$. Note that this happens independently from the weights chosen for the linear approximation: some solutions of the front just can't be reached with these approaches. Thus, successful approaches should be able to find all elements of the optimal Pareto front and converge to any of the non dominated solutions.

1) *Transposition Tables in DST*: The DST is a problem specially suited for the use of Transposition Tables (TT) [25] within MCTS. TT is a technique used to optimize three search based algorithms when two or more states can be found at different locations in the tree. It consists of sharing information between these equivalent states in a centralized manner, in order to avoid managing these positions as completely different states. Figure 7 shows an example of this situation in the DST.

In this example, the submarine starts in the initial position ($P1$, root of the tree) and makes a sequence of moves that places it in $P2$. From this location, if the vessel wishes to move to $P3$, two optimal trajectories would be route a and b . In a tree that does not use TT, there would be two different nodes to represent $P3$, although the location and number of moves performed up to this point are the same (thus, the states are equivalent). It is worthwhile to highlight that the coordinates of the vessel are not enough to identify two equivalent states, but also the number of moves is needed: imagine a third route from

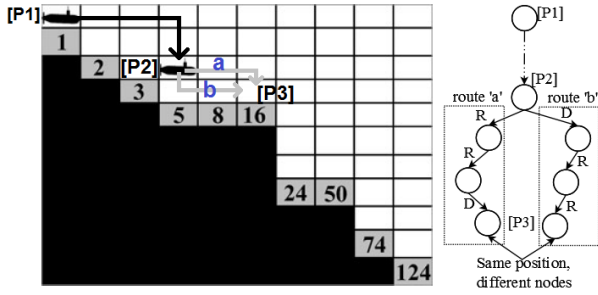


Fig. 7: Example of two different sequences of actions (R: Right, D: Down) that lie in the same position in the map, but different node in the tree.

$P2$ to $P3$ with the moves: *Up, Right, Right, Down, Down*. As the submarine performs 5 moves, the states are not the same, and the node where the ship is in $P3$ now would be two levels deeper in the tree.

TT tables are implemented in the MCTS algorithms tested in this benchmark by using hash tables that stores a *representative* node for each pair (*position, moves*) found. The key of the hash map needs then to be obtained from three values: position coordinates x and y of the ship in the board, and number of moves, indicated by the depth of the node in the tree. Hence, transpositions can only happen at the same depth within the tree, a feature that has been successfully tried before in the literature [26].

2) *Heuristics for DST*: As DST has two different objectives, number of moves and value of the treasure, the quality of a state can be assessed by two rewards, r_m and r_v , for each objective respectively. r_m is adjusted to be maximized using the maximum number of moves in the game, while r_v is just the value in the cell with the treasure. Equation 4 summarizes these rewards:

$$\begin{aligned} r_m &= 100 - \text{moves} \\ r_v &= \text{treasureValue} \end{aligned} \quad (4)$$

B. Multi-Objective PTSP

The Multi-Objective Physical Travelling Salesman Problem (MO-PTSP) is a game employed in a competition held in the IEEE Conference on Computational Intelligence in Games (CIG) in 2013, a modification of the Physical Travelling Salesman Problem (PTSP), previously introduced by Perez et al. [27]. The MO-PTSP is a real-time game where the agent navigates a ship and must visit 10 waypoints scattered around the maze. All waypoints must be visited to consider a game as complete, and a game ticks timer is reset every time a waypoint is collected, finishing the game prematurely (and unsuccessfully) if it reaches 800 game steps before visiting another waypoint.

This needs to be done while optimizing three different objectives: 1) **Time**: the player must collect all waypoints scattered around the maze in as few time steps as possible; 2) **Fuel**: the fuel consumption by the end of the game must

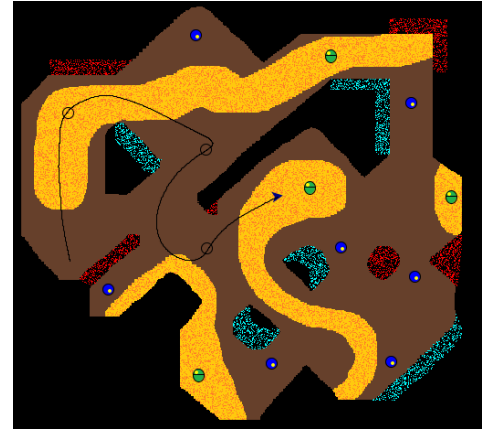


Fig. 8: Sample MO-PTSP map.

be minimized; and 3) **Damage**: the ship should end the game with as little damage as possible.

In the game, the agent must provide an action every 40 milliseconds. The available actions are combinations of two different inputs: *throttle* (that could be *on* or *off*) and *steering* (that could be *straight*, *left* or *right*). This allows for 6 different actions that modify the ship's position, velocity and orientation. These vectors are kept from one step to the next, keeping the inertia of the ship, and making the navigation task not trivial.

The ship starts with 5000 units of fuel, and one unit is spent every time an action supplied has the throttle input *on*. There are, however, two ways of collecting fuel: each waypoint visited grants the ship with 50 more units of fuel. Also, fuel canisters are scattered around the maze and their collection provides with 250 more units.

Regarding the third objective, the ship can suffer damage in two different ways: by colliding with obstacles and driving through lava. In the former case, the ship can collide with normal obstacles (that subtract 10 units of damage) and specially damaging obstacles (30 units). In the latter, lava lakes are abundant in the MO-PTSP levels and, in contrast with normal surfaces, they deal one unit of damage for each time step the ship spends over this type of surface. All this subtractions are deduced from an initial counter of 5000 points of damage.

Figure 8 shows an example of an MO-PTSP map, as drawn by the framework. Waypoints yet not visited are painted as blue circles, while those already collected are depicted as empty circles. The green ellipses represent fuel canisters, normal surfaces are drawn in brown and lava lakes are printed as red-dotted yellow surfaces. Normal obstacles black, damaging obstacles are drawn in red. Blue obstacles are elastic walls that produce no damage to the ship. The vessel is drawn as a blue polygon and its trajectory is traced with a black line.

1) *Macro-actions for MO-PTSP*: All algorithms tested in the MO-PTSP in this research employ macro-actions. Macro-actions is a concept that can be used for coarsening the action space by applying the action chosen by the control algorithm in several consecutive cycles, instead of just in the next one.

Previous research in PTSP [10], [28] suggests that using

macro-actions for real-time navigation domains increases the performance of the algorithms, and it has been used in the previous PTSP competitions by the winner and other entries.

A macro-action of length L is defined as a repetition of a given action during L consecutive time steps. The main advantage of macro-actions is that the control algorithms can see further in the future, and then reduce the implication of open-endedness of real-time games (see Section II). As this reduces the search space significantly, better algorithm performance is allowed. A consequence of using macro-actions is also that the algorithm can employ L consecutive steps to plan the next move (the next macro-action) to make. Therefore, instead of spending 40 milliseconds, as in MO-PTSP, to define the next action, the algorithm can employ $40 \times L$ milliseconds for this task, and hence perform a more extensive search.

In MO-PTSP, the macro-action size is $L = 15$, a value that has shown its proficiency before in PTSP. For more information about macro-actions and their application to real-time navigation problems such as the PTSP, the reader is referred to [10].

2) *Heuristics for MO-PTSP*: In order to evaluate a game state in MO-PTSP, three different measures or rewards are taken, r_t, r_f, r_d , one for each one of the objectives time, fuel and damage, respectively. All these rewards are defined so they have to be maximized. The first reward uses a measure of distance for the time objective, as indicated in Equation 5:

$$r_t = 1 - d/D \quad (5)$$

The MO-PTSP framework includes a path-finding library that allows controllers to query the shortest distance of a route of waypoints. d indicates the distance from the current position until the last waypoint, following the desired route, and D is the distance of the whole route from the starting position. Minimizing the distance to the last waypoint will lead to the end of the game.

Equation 6 shows how the value of the fuel objective, r_f , is obtained:

$$r_f = (1 - (f_c/f_0)) \times \alpha + r_t \times (1 - \alpha) \quad (6)$$

f_c is the fuel consumed so far, and f_0 is the initial fuel at the start of the game. α is a value that balances between the fuel component and the time objective from Equation 5. Note that, as waypoints needs to keep being visited, it is necessary to include a distance measure for this reward. Otherwise, an approach that prioritizes this objective would not minimize distance to waypoints at all (the ship could just stand still: no fuel consumed is optimal), and therefore could not complete the game. The value of α has been determined empirically, in order to provide a good balance between these two components, and it is set to 0.66.

Finally, Equation 7 gives the method to calculate the damage objective, r_d :

$$r_d = \begin{cases} (1 - (d_s/d_m)) \times \beta_1 + r_t \times (1 - \beta_1), & s > \gamma \\ (1 - (d_s/d_m)) \times \beta_2 + r_t \times (1 - \beta_2), & s \leq \gamma \end{cases} \quad (7)$$

d_s is the damage suffered so far in the game, and d_m is the maximum damage the ship can take. In this case, three different variables are used to regulate the behaviour of this objective: γ, β_1 and β_2 . Both β_1 and β_2 have the same role as α in Equation 5: they balance between the time objective and the damage measure. The difference is that β_1 is used in high speeds, while β_2 is employed with low velocities. This is distinguished by the parameter γ , that can be seen as a threshold value for the ship's speed (s). This differentiation is made in order to avoid low speeds in lava lakes, as this increase the damage suffered importantly. The values for these variables have been determined also empirically, and they are set to $\gamma = 0.8, \beta_1 = 0.75$ and $\beta_2 = 0.25$.

VI. EXPERIMENTATION

The experiments performed in this research compare three different algorithms in the two benchmarks presented in Section V: a single objective MCTS (referred to here simply as *MCTS*), the Multi-Objective MCTS (*MO-MCTS*) and a rolling horizon version of the NSGA-II algorithm described in Section III (*NSGA-II*). This NSGA-II version evolves a population where the individuals are sequences of actions (macro-actions in the MO-PTSP case), obtaining the fitness from the state of the game after applying such sequence.

All algorithms have a limited number of evaluations before providing an action to perform. In order for these games to be real-time, the time budget allowed is close to 40 milliseconds. With the objective of avoiding congestion peaks at the machine where the experiments are run, an average of the number of evaluations doable in 40 ms. is calculated and employed in the tests. This led to 4500 evaluations in the DST and 500 evaluations for MO-PTSP, in the same server where the PTSP and MO-PTSP competitions were run¹.

A. Results in DST

Figure 9

B. Results in MO-PTSP

Table I

C. A step further in MO-PTSP: dynamically varying weights

VII. CONCLUSIONS AND FUTURE WORK

ACKNOWLEDGMENTS

This work was supported by EPSRC grants EP/H048588/1, under the project entitled "UCT for Games and Beyond".

REFERENCES

- [1] C.-S. Lee, M.-H. Wang, G. M. J.-B. Chaslot, J.-B. Hoock, A. Rimmel, O. Teytaud, S.-R. Tsai, S.-C. Hsu, and T.-P. Hong, "The Computational Intelligence of MoGo Revealed in Taiwan's Computer Go Tournaments," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 1, no. 1, pp. 73–89, 2009.
- [2] C. Browne, E. Powley, D. Whitehouse, S. Lucas, P. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakakis, and S. Colton, "A Survey of Monte Carlo Tree Search Methods," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 4:1, pp. 1–43, 2012.

¹Intel Core i5, 230GHz, 4GB of RAM

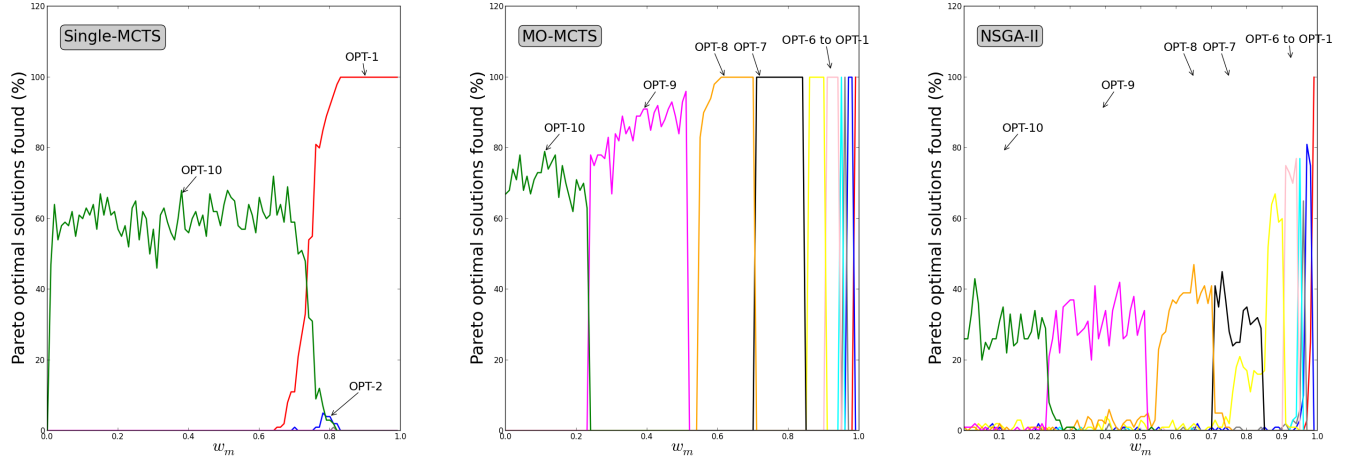


Fig. 9: Results DST

| | $W : (w_t, w_f, w_d)$ | MO-MCTS (D, \emptyset, d) | MCTS (D, \emptyset, d) | NSGA-II (D, \emptyset, d) | PurofMovio (D, \emptyset, d) |
|------------|----------------------------|-----------------------------|--------------------------|-----------------------------|--------------------------------|
| MO-MCTS | $W_0 : (0.33, 0.33, 0.33)$ | — | $(8, 2, 0)$ | $(8, 2, 0)$ | $(0, 5, 5)$ |
| | $W_1 : (0.1, 0.3, 0.6)$ | | $(10, 0, 0)$ | $(4, 6, 0)$ | $(0, 6, 4)$ |
| | $W_2 : (0.1, 0.6, 0.3)$ | | $(8, 2, 0)$ | $(7, 3, 0)$ | $(0, 5, 5)$ |
| MCTS | $W_0 : (0.33, 0.33, 0.33)$ | $(0, 8, 2)$ | — | $(0, 2, 8)$ | $(0, 2, 8)$ |
| | $W_1 : (0.1, 0.3, 0.6)$ | $(0, 0, 10)$ | | $(4, 2, 4)$ | $(0, 3, 7)$ |
| | $W_2 : (0.1, 0.6, 0.3)$ | $(0, 2, 8)$ | | $(0, 1, 9)$ | $(0, 6, 4)$ |
| NSGA-II | $W_0 : (0.33, 0.33, 0.33)$ | $(0, 2, 8)$ | $(8, 2, 0)$ | — | $(0, 4, 6)$ |
| | $W_1 : (0.1, 0.3, 0.6)$ | $(0, 6, 4)$ | $(4, 2, 4)$ | | $(0, 4, 6)$ |
| | $W_2 : (0.1, 0.6, 0.3)$ | $(0, 3, 7)$ | $(9, 1, 0)$ | | $(0, 5, 5)$ |
| PurofMovio | $W_0 : (0.33, 0.33, 0.33)$ | $(5, 5, 0)$ | $(8, 2, 0)$ | $(6, 4, 0)$ | — |
| | $W_1 : (0.1, 0.6, 0.3)$ | $(4, 6, 0)$ | $(7, 3, 0)$ | $(6, 4, 0)$ | |
| | $W_2 : (0.1, 0.3, 0.6)$ | $(5, 5, 0)$ | $(6, 4, 0)$ | $(5, 5, 0)$ | |

TABLE I: Results MO-PTSP

- [3] S. Gelly, Y. Wang, R. Munos, and O. Teytaud, "Modification of UCT with Patterns in Monte-Carlo Go," Inst. Nat. Rech. Inform. Auto. (INRIA), Paris, Tech. Rep., 2006.
- [4] P.-A. Coquelin, R. Munos *et al.*, "Bandit Algorithms for Tree Search," in *Uncertainty in Artificial Intelligence*, 2007.
- [5] D. Robles and S. M. Lucas, "A Simple Tree Search Method for Playing Ms. Pac-Man," in *Proceedings of the IEEE Conference of Computational Intelligence in Games*, Milan, Italy, 2009, pp. 249–255.
- [6] S. Samothrakis, D. Robles, and S. M. Lucas, "Fast Approximate Max-n Monte-Carlo Tree Search for Ms Pac-Man," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 3, no. 2, pp. 142–154, 2011.
- [7] N. Ikehata and T. Ito, "Monte Carlo Tree Search in Ms. Pac-Man," in *Proc. 15th Game Programming Workshop*, Kanagawa, Japan, 2010, pp. 1–8.
- [8] S. Matsumoto, N. Hirose, K. Itonaga, K. Yokoo, and H. Futahashi, "Evaluation of Simulation Strategy on Single-Player Monte-Carlo Tree Search and its Discussion for a Practical Scheduling Problem," in *Proceedings of the International MultiConference of Engineers and Computer Scientists*, vol. 3, Hong Kong, 2010, pp. 2086–2091.
- [9] S. Edelkamp, P. Kissmann, D. Sulewski, and H. Messerschmidt, "Finding the Needle in the Haystack with Heuristically Guided Swarm Tree Search," in *Multikonf. Wirtschaftsinform.*, Göttingen, Germany, 2010, pp. 2295–2308.
- [10] D. Perez, E. J. Powley, D. Whitehouse, P. Rohlfshagen, S. Samothrakis, P. I. Cowling, and S. M. Lucas, "Solving the Physical Travelling Salesman Problem: Tree Search and Macro-Actions," *IEEE Trans. Comp. Intell. AI Games (submitted)*, 2013.
- [11] E. J. Powley, D. Whitehouse, and P. I. Cowling, "Monte Carlo Tree Search with macro-actions and heuristic route planning for the Physical Travelling Salesman Problem," in *Proc. IEEE Conf. Comput. Intell. Games*, 2012, pp. 234–241.
- [12] —, "Monte Carlo Tree Search with Macro-Actions and Heuristic Route Planning for the Multiobjective Physical Travelling Salesman Problem," in *Proc. IEEE Conf. Comput. Intell. Games*, 2013, pp. 73–80.
- [13] K. Deb, *Multi-Objective Optimization using Evolutionary Algorithms*. Wiley, 2001.
- [14] C. Coello, "An Updated Survey of Evolutionary Multiobjective Optimization Techniques: State of the Art and Future Trends," in *Proc. of the Congress on Evolutionary Computation*, 1999, pp. 3–13.
- [15] A. Zhou, B.-Y. Qu, H. Li, S.-Z. Zhao, P. N. Suganthan, and Q. Zhang, "Multiobjective Evolutionary Algorithms: A Survey of the State of the Art," *Swarm and Evolutionary Computation*, vol. 1, pp. 32–49, 2011.
- [16] K. Deb, S. Agrawal, A. Pratap, and T. Meyarivan, "A Fast Elitist Non-dominated Sorting Genetic Algorithm for Multi-objective Optimization: NSGA-II," in *Proceedings of the 6th International Conference on Parallel Problem Solving from Nature*, 2000, pp. 1–11.
- [17] Q. Zhang and H. Li, "MOEA/D: A Multi-objective Evolutionary Algorithm Based on Decomposition," *IEEE Transactions on Evolutionary Computation*, vol. 11, no. 6, pp. 712–731, 2007.
- [18] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction (Adaptive Computation and Machine Learning)*. A Bradford Book, 1998.
- [19] P. Vamplew, R. Dazeley, A. Berry, R. Issabekov, and E. Dekker, "Empirical Evaluation Methods for Multiobjective Reinforcement Learning Algorithms," *Machine Learning*, vol. 84, pp. 51–80, 2010.
- [20] Z. Gabor, Z. Kalmar, and C. Szepesvari, "Multi-criteria Reinforcement Learning," in *The fifteenth international conference on machine learning*, 1998, pp. 197–205.
- [21] S. Natarajan and P. Tadepalli, "Dynamic Preferences in Multi-Criteria Reinforcement Learning," in *In Proceedings of International Conference of Machine Learning*, 2005, pp. 601–608.
- [22] L. Barrett and S. Narayanan, "Learning All Optimal Policies with Multiple Criteria," in *Proceedings of the international conference on machine learning*, 2008.
- [23] W. Weijia and M. Sebag, "Multi-objective Monte Carlo Tree Search," in *Proceedings of the Asian Conference on Machine Learning*, 2012, pp. 507–522.
- [24] —, "Hypervolume indicator and dominance reward based multi-

objective Monte-Carlo Tree Search,” *Machine Learning*, vol. 92:2–3, pp. 403–429, 2013.

- [25] B. E. Childs, J. H. Brodeur, and L. Kocsis, “Transpositions and Move Groups in Monte Carlo Tree Search,” in *Proceedings of IEEE Symposium on Computational Intelligence and Games*, 2008, pp. 389–395.
- [26] T. Kozelek, “Methods of MCTS and the game Arimaa,” M.S. thesis, Charles Univ., Prague, 2009.
- [27] D. Perez, P. Rohlfshagen, and S. Lucas, “The Physical Travelling Salesman Problem: WCCI 2012 Competition,” in *Proceedings of the IEEE Congress on Evolutionary Computation*, 2012.
- [28] —, “Monte Carlo Tree Search: Long Term versus Short Term Planning,” in *Proceedings of the IEEE Conference on Computational Intelligence and Games*, 2012, pp. 219 – 226.



Diego Perez received a B.Sc. and a M.Sc. in Computer Science from University Carlos III, Madrid, in 2007. He is currently pursuing a Ph.D. in Artificial Intelligence applied to games at the University of Essex, Colchester. He has published in the domain of Game AI, participated in several Game AI competitions and organized the Physical Travelling Salesman Problem competition, held in IEEE conferences during 2012. He also has programming experience in the videogames industry with titles published for game consoles and PC.

Sanaz Mostaghim BIO HERE.

**PHOTO
HERE**



Spyridon Samothrakis is currently pursuing a PhD in Computational Intelligence and Games at the University of Essex. His interests include game theory, computational neuroscience, evolutionary algorithms and consciousness.



Simon Lucas (SMIEEE) is a professor of Computer Science at the University of Essex (UK) where he leads the Game Intelligence Group. His main research interests are games, evolutionary computation, and machine learning, and he has published widely in these fields with over 160 peer-reviewed papers. He is the inventor of the scanning n-tuple classifier, and is the founding Editor-in-Chief of the IEEE Transactions on Computational Intelligence and AI in Games.