

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

Multi-objective optimization has been a field of study of manufacturing, engineering [1] and finance [2], while having little impact in games research. This paper shows that multi-objective optimization has much to offer in developing game strategies that allow for a fine-grain control of alternative policies. The application of such approaches to this field can provide interesting results, especially in games that are long or complex enough where long-term planning is not trivial and achieving a good level of play requires balancing strategies.

From a simple point of view, many competitive games can be seen as a scenario where two or more opponents have the single objective of winning. Achieving victory is usually complex in interesting games, and successful approaches normally provide a state value by combining features of the game state as a weighted sum. An example is a chess heuristic that assigns different weights to each piece according to their value.

Multi-objective approaches can be applied in these scenarios, where several factors must be taken into account to achieve the victory. The algorithms can balance between different objectives, in order to provide a wide range of strategies well suited to the different stages of the game being played, or to face the existing opponents. Application of such approaches could be real-time strategy games, where long term planning must be carried out in order to balance aspects such as attack units, defensive structures and resource gathering.

This paper proposes a multi-objective real-time algorithm version of Monte Carlo Tree Search (MCTS), a popular reinforcement learning approach within the last decade. The proposed algorithm is tested in two different games: a real-time version of the Deep Sea Treasure (DST; a classical multi-objective problem), and the Multi-Objective Physical Travelling Salesman Problem (MO-PTSP). The algorithm is also compared to two other approaches: a single-objective MCTS (that uses a weighted sum of features to value the state) and a rolling horizon Non-dominated Sorting Evolutionary Algorithm II (NSGA-II). This paper extends and formalizes our previous work described in [3].

Two main goals can be identified in this paper: first, the proposed algorithm must be applicable to real-time domains (those where the next move to make must be decided within

a small time budget) and it should obtain better or at least the same performance than the other state of the art algorithms. It is important to highlight that all three algorithms tested employ the same heuristic functions to evaluate the features of a given game state. Thus, the focus of this research is set on how the algorithm explores the search space, instead of providing the best possible solution to each given problem.

Secondly, the algorithm must be able to provide solutions across the multi-objective spectrum: by parametrizing the algorithm, it must be possible to prioritize one objective over the others and therefore converge to solutions according to these preferences.

The paper is structured as follows. First, Sections II and III provide the necessary background for MCTS and multi-objective optimization, respectively. The algorithm proposed in this research is described in detail in Section IV. Then, Section V defines the games used to test the algorithms, with the results discussed in Section VI. Finally, some conclusions and possible extensions of this work are drawn in Section VII.

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 [4]. 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. [5].

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

Diego Perez, Spyridon Samothrakis, Simon M. Lucas (School of Computer Science and Electronic Engineering, University of Essex, Colchester CO4 3SQ, UK; email: {dperez, ssamot, sml}@essex.ac.uk); Sanaz Mostaghim (Department of Knowledge and Language Engineering, Otto-von-Guericke-Universitt Magdeburg, Magdeburg, Germany; email: sanaz.mostaghim@ovgu.de)

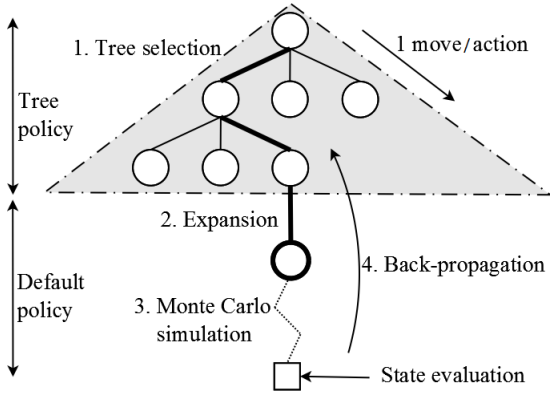


Fig. 1: MCTS algorithm steps.

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 [6]: *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* 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 [7].

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 [8] 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 [9], [10].

MCTS has also been applied to single-player games, like SameGame [11], 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 [12], 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 [13], [14] and the multi-objective versions [15]. 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 $i = 1, 2, \dots, m$.

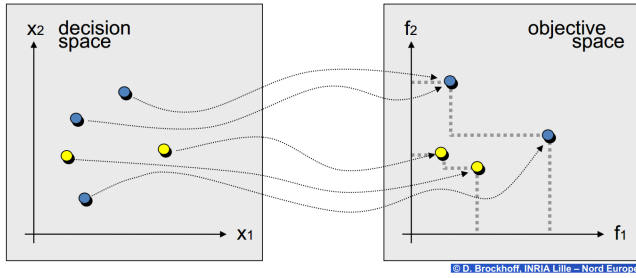


Fig. 2: Decision and Solution spaces in a MOOP with two variables (x_1 and x_2) and two objectives (f_1 and f_2). In the objective space, yellow dots are non-optimal objective vectors, while blue dots form the Pareto-optimal front.

- 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 [16].

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 [16] explains how linear scalarization approaches fail in

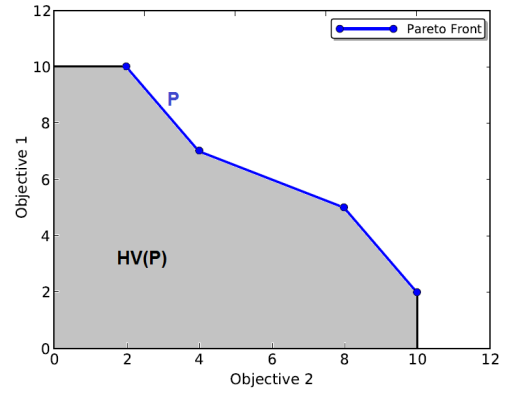


Fig. 3: $HV(P)$ of a given Pareto front P .

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 [17], [18]. 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 [19].

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) [20], 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 [21] 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) [22] 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. [22] 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 [23], 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. [24].

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 [25], 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

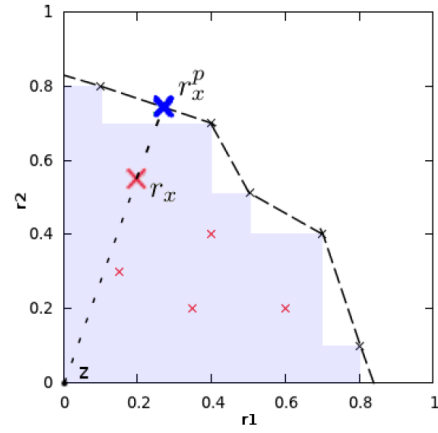


Fig. 4: r_x^p is the projection of the r_x value on the piecewise linear surface (discontinuous line). The shadowed surface represents $HV(P)$. From [26].

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 [26], [27]. 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. $P = \text{node}.P \cup \bar{r}$ 
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

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).

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

This algorithm, same as NSGA-II due to their multi-objective nature, provide a non-dominated front as a solution. However, in planning and control scenarios like the games analyzed in this research, an action must be provided to perform a move in the next step. The question that arises is how to choose the move to make with the information available.

As shown before, it is straightforward to obtain which actions lead to what points in the front given as a solution: first gene in an NSGA-II individual, a root's child in MO-MCTS. Hence, by identifying the point in the Pareto front that the algorithm should converge to, it is possible to execute the action that leads to that point.

In order to select the point in the Pareto front, a weight vector $W(\forall w_i \in W, w_i \in [0, 1])$ can be provided with the number of objectives as its dimension. Two different mechanisms are proposed, for reasons that will be explained in the experiments section (VI):

- **Weighted Sum:** the action chosen is the one that maximizes the weighted sum of the solution vector multiplied by W , for each point in the front.
- **Euclidean distance:** the euclidean distance from each point in the Pareto front (normalized in $[0, 1]$) to the vector W is calculated. The action to choose would be the one that leads to the point in the Pareto front with the smaller distance to W .

Note that, in the vanilla MCTS, there is no Pareto front obtained as a solution. Typically, in this case, rewards are calculated as a weighted sum of the objectives and a weight vector W , and then the action is chosen following any of the mechanisms usually employed in the literature: the action taken more often from the root, the one that leads to the best reward found, the move with the highest expected reward or the action that maximizes Equation 1 in the root.

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. [22]. 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

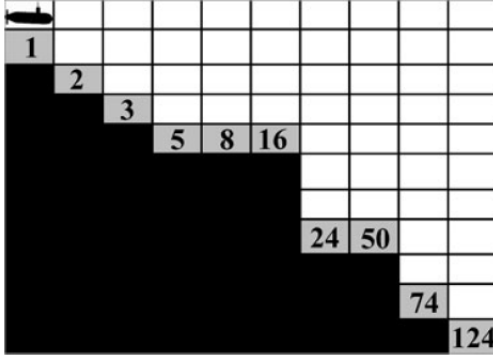


Fig. 5: Environment of the Deep Sea Treasure (from [22]): grey squares represent the treasure (with their different values) available in the map. The black cells are the sea floor and the white ones are the positions that the vessel can occupy freely. The game ends when the submarine picks one of the treasures.

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 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) [28] within MCTS. TT is a technique used to optimize three search

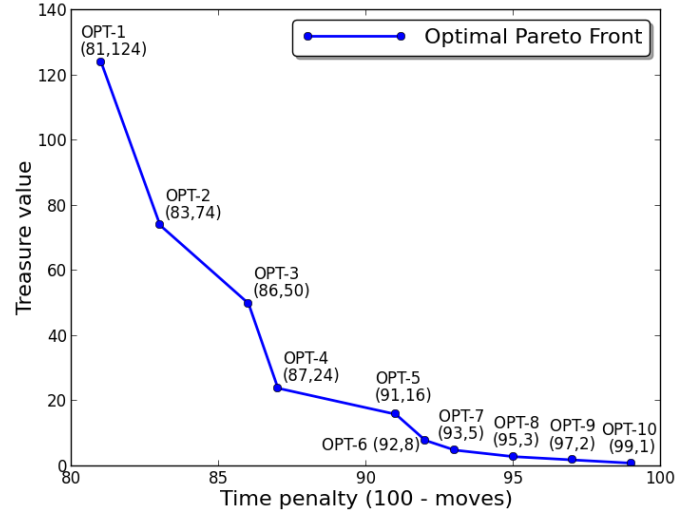


Fig. 6: Optimal Pareto Front of the Deep Sea Treasure, with both objectives to be maximized.

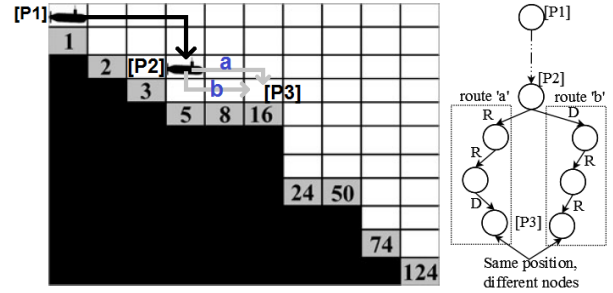


Fig. 7: Example of two different sequences of actions (R: Right, D: Down) that lie in the same position in the map, but different node in the tree.

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 $P2$ 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 [29].

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. [30]. 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 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

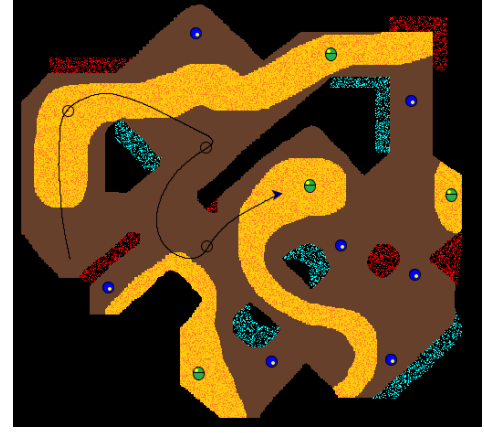


Fig. 8: Sample MO-PTSP map.

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 [13], [31] 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 [13].

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 refeq:fuel: 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, the average of 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, using the same server where the PTSP and MO-PTSP competitions were run¹.

A. Results in DST

As the optimal Pareto front of the DST is known, a measure of performance can be obtained by observing the percentage of times these solutions are found by the agents. As the solution that the algorithms converge to depends on the weights vector employed during the search (see end of Section IV), the approach taken here is to provide different weight vectors W and analyze them separately.

The weight vector for DST has two dimensions. This vector is here referred to as $W = (w_m, w_v)$, where w_m weights moves and w_v the treasure value, and $w_v = 1 - w_m$. w_m takes values between 0 and 1, with a precision of 0.01, and 100 experiments have been run for each pair (w_m, w_v) . Hence, the game has been played a total of 10000 times, for each algorithm.

Figure 9 shows the results obtained after these experiments were performed. The first point to notice is that MCTS only converges (mostly) to the two optimal points located at the edges of the optimal Pareto front (*OPT-1* and *OPT-10*, see also Figure 6). This is an expected result, as K. Deb. suggested in [16] and was also discussed before in Section III: linear scalarization approaches only converge to the edges of the optimal Pareto front if its shape is non-convex.

The results show clearly how approximating the optimal Pareto front allows for finding all possible solutions. Both NSGA-II and MO-MCTS approaches are able to converge to any solution in the front given the appropriate weight vector. It is important to highlight that these two algorithms employed the euclidean distance action selection (see Section IV). Experiments shown that weighted sum action selection provides similar results as in MCTS (this is, convergence to the edges of the front). The crucial distinction to make is that both algorithms NSGA-II and MO-MCTS allow for different action selection mechanisms, that are able to overcome this problem, by approximating a global Pareto front.

Finally, in the comparison between NSGA-II and MO-MCTS, the latter algorithm obtains higher percentages for each one of the points in the front. This result suggests that, with a limited number of iterations/evaluations, the proposed algorithm is able to explore the search space more efficiently than a state of the art algorithm such NSGA-II.

¹Intel Core i5, 230GHz, 4GB of RAM

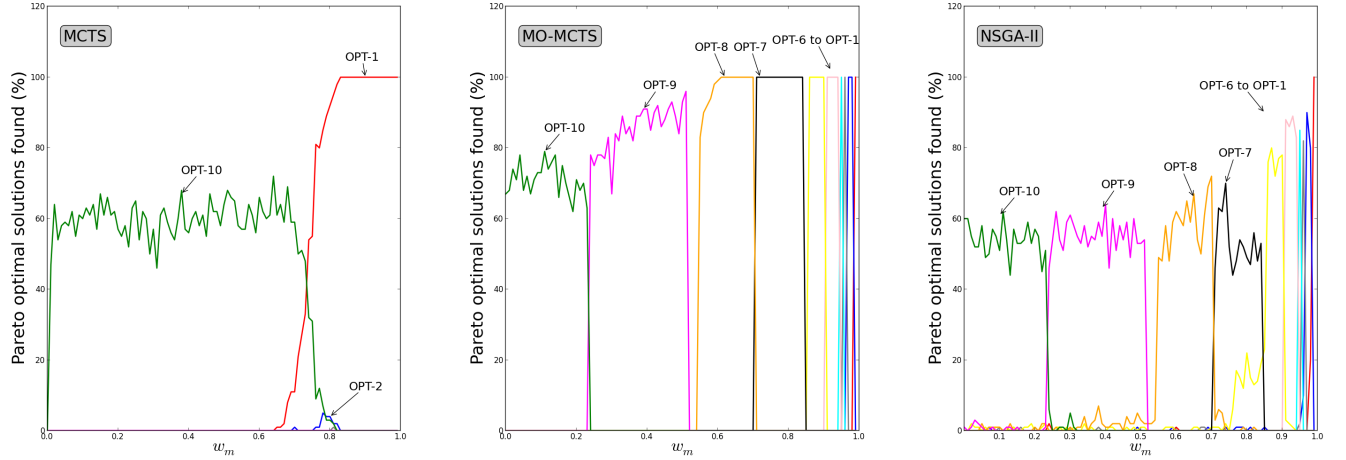


Fig. 9: Results in DST: percentages of each optima found during 100 games played with different weight vectors. Scalarization approaches converge to the edges of the optimal front, whereas Pareto approaches are able to find all optimal solutions. The proposed algorithm, MO-MCTS, finds these solutions significantly more often than NSGA-II.

B. Results in MO-PTSP

The same three algorithms have been tested in the MO-PTSP domain. The framework² comes with 10 maps, used in this research to run the algorithms in. In this case, the optimal Pareto front is not known in advance, and normally it will be different from one map to another. Hence, the mechanisms to compare the performance of the algorithms tested need to be different than the one used for DST.

The idea is as follows. First of all, 4 different weight vectors are tested: $W_1 = (0.33, 0.33, 0.33)$, $W_2 = (0.1, 0.3, 0.6)$, $W_3 = (0.1, 0.6, 0.3)$ and $W_4 = (0.6, 0.1, 0.3)$, where each w_i corresponds to the weight for an objective (w_t for time, w_f for fuel and w_d for damage). W_1 treats all objectives with the same weight, while the other three give more relevance to different objectives. These vectors provide then a wide spectrum for weights in this benchmark. In this case, MO-MCTS uses this weights to select an action based on a weighted sum, which has shown in preliminary experiments better performance than the Euclidean distance mechanism.

The first point to check is if the different weight vectors affect the solutions obtained by MO-MCTS. Table I shows the results of executing the MO-MCTS controller during 30 games in the 10 maps of the MO-PTSP, for every weight vector. It can be seen that the highest weight in W leads, in most of the cases, to the best solution in that objective in the map.

Some exceptions can actually be explained analyzing the map: in map 1, a fuel canister is always collected near the end, restoring the fuel level to its maximum. This leaves too few cycles to make a difference in the controller (note also that this is the map with the fewest fuel consumption). Also, in map 10, there is no difference in the damage objective: but precisely map 10 is a map with no obstacles to damage the ship (thus only lava lakes deals damage). This also results in this map being the one with lowest damage overall. It can also be seen that, in those maps where time has priority, the results in this objective are not as dominant as the other two. This

Map	$W : (w_t, w_f, w_d)$	Time	Fuel	Damage
Map 1	(0.33, 0.33, 0.33)	1654 ± 7	131 ± 2	846 ± 13
	(0.1, 0.3, 0.6)	1657 ± 8	130 ± 2	773 ± 11
	(0.1, 0.6, 0.3)	1681 ± 11	131 ± 2	837 ± 15
	(0.6, 0.1, 0.3)	1649 ± 8	132 ± 2	833 ± 13
Map 2	(0.33, 0.33, 0.33)	1409 ± 7	235 ± 4	364 ± 3
	(0.1, 0.3, 0.6)	1402 ± 6	236 ± 5	354 ± 2
	(0.1, 0.6, 0.3)	1416 ± 8	219 ± 4	360 ± 3
	(0.6, 0.1, 0.3)	1396 ± 8	245 ± 5	361 ± 2
Map 3	(0.33, 0.33, 0.33)	1373 ± 6	221 ± 3	301 ± 7
	(0.1, 0.3, 0.6)	1378 ± 5	211 ± 3	268 ± 5
	(0.1, 0.6, 0.3)	1385 ± 5	203 ± 4	291 ± 7
	(0.6, 0.1, 0.3)	1363 ± 4	229 ± 4	285 ± 7
Map 4	(0.33, 0.33, 0.33)	1383 ± 6	291 ± 5	565 ± 5
	(0.1, 0.3, 0.6)	1385 ± 7	304 ± 4	542 ± 4
	(0.1, 0.6, 0.3)	1423 ± 6	273 ± 4	583 ± 5
	(0.6, 0.1, 0.3)	1388 ± 6	309 ± 4	559 ± 5
Map 5	(0.33, 0.33, 0.33)	1405 ± 7	467 ± 4	559 ± 4
	(0.1, 0.3, 0.6)	1431 ± 9	447 ± 4	541 ± 4
	(0.1, 0.6, 0.3)	1467 ± 9	411 ± 5	567 ± 5
	(0.6, 0.1, 0.3)	1399 ± 9	469 ± 4	547 ± 3
Map 6	(0.33, 0.33, 0.33)	1575 ± 7	549 ± 5	303 ± 4
	(0.1, 0.3, 0.6)	1626 ± 9	540 ± 6	286 ± 5
	(0.1, 0.6, 0.3)	1703 ± 11	499 ± 4	316 ± 7
	(0.6, 0.1, 0.3)	1571 ± 7	559 ± 5	294 ± 4
Map 7	(0.33, 0.33, 0.33)	1434 ± 5	599 ± 6	284 ± 6
	(0.1, 0.3, 0.6)	1475 ± 10	602 ± 5	243 ± 6
	(0.1, 0.6, 0.3)	1489 ± 12	549 ± 3	264 ± 6
	(0.6, 0.1, 0.3)	1407 ± 8	618 ± 5	270 ± 6
Map 8	(0.33, 0.33, 0.33)	1761 ± 9	254 ± 5	382 ± 3
	(0.1, 0.3, 0.6)	1804 ± 10	269 ± 4	357 ± 4
	(0.1, 0.6, 0.3)	1826 ± 10	230 ± 3	392 ± 7
	(0.6, 0.1, 0.3)	1732 ± 9	311 ± 8	379 ± 6
Map 9	(0.33, 0.33, 0.33)	2501 ± 14	926 ± 6	574 ± 9
	(0.1, 0.3, 0.6)	2503 ± 10	921 ± 10	524 ± 8
	(0.1, 0.6, 0.3)	2641 ± 14	833 ± 5	574 ± 14
	(0.6, 0.1, 0.3)	2470 ± 9	956 ± 5	573 ± 8
Map 10	(0.33, 0.33, 0.33)	1430 ± 8	630 ± 4	205 ± 2
	(0.1, 0.3, 0.6)	1493 ± 13	615 ± 4	209 ± 2
	(0.1, 0.6, 0.3)	1542 ± 10	554 ± 4	229 ± 5
	(0.6, 0.1, 0.3)	1378 ± 5	663 ± 6	202 ± 4

TABLE I: MO-MCTS results in MO-PTSP with different weight vectors. Significantly better results per map are in bold.

²www.ptsp-game.net

can be explained by the fact that the time heuristic is actually

part of the fuel and damage heuristics (Equation 5 in Eqs. 6 and 7).

These results suggest that the weights effectively work on making the algorithm converge to different points in the Pareto front discovered by the algorithm. Now, it is time to compare these results with the other algorithms. In order to do this, the same number of experiments is run for NSGA-II and MCTS in the 10 maps of the benchmark.

In order to compare several algorithms between each other, the results are compared in pairs, in terms of dominance. The procedure is as follows: once all games on a single map have been run, the MannWhitneyWilcoxon non-parametric test with 95% confidence is calculated on the three objectives. If the measures on all objectives are assumed to be drawn from different distributions, their averages are compared for a dominance test. If one result dominates another, then the algorithm dominates the other in that particular map.

Extending this comparison to all maps, each pair of algorithms ends with a triplet (D, \emptyset, d) , where D is the number of maps where the first algorithm dominates the second, \emptyset is the amount of maps where no dominance can be established, and d states the number of maps where the first is dominated by the second. Table II summarizes these results for all algorithms tested.

One of the first things to notice is that MO-MCTS dominates MCTS and NSGA-II in most of the maps, and it is never dominated by these. Particularly, the dominance of MO-MCTS over MCTS is outstanding, even dominating in all 10 maps for two of the weight vectors. MO-MCTS also dominates NSGA-II in more maps than in those where there is no dominance, and it is as well never dominated by NSGA-II in any map.

It is also interesting to see that NSGA-II dominates also the weighted sum version of MCTS, although for the vector W_2 there is a technical draw, as they dominate each other in 4 different maps each and there is no dominance in the other 2.

Additionally, Table II contains a fourth entry, *PurofMovio*, that the algorithms are compared against. *PurofMovio* is the winner entry of the 2013 MO-PTSP competition, a controller based in a weighted-sum MCTS approach (see [15] for details of its implementation). As can be seen, *PurofMovio* obtains better results than the algorithm proposed in this paper.

However, it is very important to highlight that *PurofMovio* is not using the same heuristics as the ones presented in this research. Hence, nothing can be concluded from making pairwise comparisons directly with the winning entry of the competition. It is likely that *PurofMovio*'s heuristics are more efficient than the ones presented here, but the goal of this paper is not to develop the best heuristics for the MO-PTSP, but to provide an insight of how a multi-objective version of MCTS compares to other algorithms using the same heuristics.

Anyway, the inclusion of *PurofMovio* in this comparison is not pointless: it is possible to assess the quality of the three algorithms tested here by comparing their performance relatively, against this high quality entry. Attending to this criteria, it can be seen how MO-MCTS is the algorithm that is dominated less often by him, producing similar results on an

average of 4.75 out of the 10 maps, and being dominated in 5.25 maps. MCTS and NSGA-II are dominated more often than MO-MCTS, being dominated on averages of 7.5 and 5.75 of the maps, respectively. This comparative result shows again that MO-MCTS is achieving the best results among the algorithms compared.

C. A step further in MO-PTSP: segments and weights

There is another aspect that can still be improved in the MO-PTSP benchmark, but that it is also applicable to other domains. It is naive to think that an unique weight vector will be the ideal one for the whole game. Specifically in the MO-PTSP, there are regions of the map where there are more obstacles on lava lakes, hence most likely to damage the ship. Also, the route followed during the game affects the relative ideal speed between waypoints, or sometimes a fuel canister is going to be picked up, which affects how the fuel objective can be managed. In general, many real-time games go through different phases, with different objectives and priorities.

A way to provide different weights at different times in MO-PTSP is straightforward. Given the route of waypoints (or fuel canisters) that is being followed, one can divide it into *segments*, where a segment starts and ends with a waypoint (or fuel canister). Then, each segment can be assigned a particular weight vector W .

The question is then how to assign these weight vectors. Three different ways can be devised:

- Manually set the weight vectors. This was attempted and it proved to be a non trivial task.
- Setting the appropriate weight for each segment dynamically, based on the segment's characteristics. This involves the creation or discovery of features and some kind of function approximation to assign the values.
- Learn, for each specific map, the combination of weight vectors that produces better results.

This section details the efforts made in order to test the third variant, by using a stochastic hill climbing algorithm on each one of the maps. The goal is to check if by varying the weight vectors between segments, better solutions can be achieved.

A solution is identified by a string of integers, where each integer refers to one of the weight vectors utilized in the previous sections. The solution is evaluated playing a particular map 10 times, and its fitness is obtained by calculating the average of those runs. A population of 10 individuals is kept, and the solutions of the initial population are created either randomly, or mutated from base individuals. These base individuals have all segments with the same weight vector W_1 , W_2 , or W_3 (W_4 is out of this experiment, as it shown not to be that influential).

The best solution (determined by dominance) is promoted to the next generation and it is mutated to generate other individuals. Also, a portion of individuals is generated randomly at every generation, until the end of the algorithm. Table III shows the results obtained on each run, one per map.

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".

	$W : (w_t, w_f, w_d)$	MO-MCTS (D, \emptyset, d)	MCTS (D, \emptyset, d)	NSGA-II (D, \emptyset, d)	PurofMovio (D, \emptyset, d)
MO-MCTS	$W_1 : (0.33, 0.33, 0.33)$	—	(8, 2, 0)	(8, 2, 0)	(0, 5, 5)
	$W_2 : (0.1, 0.3, 0.6)$		(10, 0, 0)	(4, 6, 0)	(0, 6, 4)
	$W_3 : (0.1, 0.6, 0.3)$		(8, 2, 0)	(7, 3, 0)	(0, 5, 5)
	$W_4 : (0.6, 0.1, 0.3)$		(10, 0, 0)	(3, 7, 0)	(0, 3, 7)
MCTS	$W_1 : (0.33, 0.33, 0.33)$	(0, 8, 2)	—	(0, 2, 8)	(0, 2, 8)
	$W_2 : (0.1, 0.3, 0.6)$	(0, 0, 10)		(4, 2, 4)	(0, 3, 7)
	$W_3 : (0.1, 0.6, 0.3)$	(0, 2, 8)		(0, 1, 9)	(0, 6, 4)
	$W_4 : (0.6, 0.1, 0.3)$	(0, 0, 10)		(3, 3, 4)	(0, 1, 9)
NSGA-II	$W_1 : (0.33, 0.33, 0.33)$	(0, 2, 8)	(8, 2, 0)	—	(0, 4, 6)
	$W_2 : (0.1, 0.3, 0.6)$	(0, 6, 4)	(4, 2, 4)		(0, 4, 6)
	$W_3 : (0.1, 0.6, 0.3)$	(0, 3, 7)	(9, 1, 0)		(0, 5, 5)
	$W_4 : (0.6, 0.1, 0.3)$	(0, 7, 3)	(4, 3, 3)		(0, 4, 6)
PurofMovio	$W_1 : (0.33, 0.33, 0.33)$	(5, 5, 0)	(8, 2, 0)	(6, 4, 0)	—
	$W_2 : (0.1, 0.6, 0.3)$	(4, 6, 0)	(7, 3, 0)	(6, 4, 0)	
	$W_3 : (0.1, 0.3, 0.6)$	(5, 5, 0)	(6, 4, 0)	(5, 5, 0)	
	$W_4 : (0.6, 0.1, 0.3)$	(7, 3, 0)	(9, 1, 0)	(6, 4, 0)	

TABLE II: Results in MO-PTSP: Each cell indicates the triplet (D, \emptyset, d) , where D is the number of maps where the row algorithm dominates the column one, \emptyset is the amount of maps where no dominance can be established, and d states the number of maps where the row algorithm is dominated by the column one. All algorithms followed the same route (order of waypoints and fuel canisters) in every map tested.

REFERENCES

- [1] R. T. Marler and J. S. Arora, "Survey of Multi-objective Optimization Methods for Engineering," *Structural and Multidisciplinary Optimization*, vol. 26, pp. 369–395, 2004.
- [2] C. Coello, *Handbook of Research on Nature Inspired Computing for Economy and Management*. Idea Group Publishing, 2006, ch. Evolutionary Multi-Objective Optimization and its Use in Finance.
- [3] D. Perez, S. Samothrakis, and S. Lucas, "Online and offline learning in multi-objective monte carlo tree search," in *Proceedings of the Conference on Computational Intelligence and Games (CIG)*, 2013.
- [4] 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.
- [5] C. Browne, E. Powley, D. Whitehouse, S. Lucas, P. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, 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.
- [6] 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.
- [7] P.-A. Coquelin, R. Munos *et al.*, "Bandit Algorithms for Tree Search," in *Uncertainty in Artificial Intelligence*, 2007.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] 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.
- [14] 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.
- [15] —, "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.
- [16] K. Deb, *Multi-Objective Optimization using Evolutionary Algorithms*. Wiley, 2001.
- [17] 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.
- [18] 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.
- [19] 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.
- [20] 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.
- [21] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction (Adaptive Computation and Machine Learning)*. A Bradford Book, 1998.
- [22] 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.
- [23] Z. Gabor, Z. Kalmar, and C. Szepesvari, "Multi-criteria Reinforcement Learning," in *The fifteenth international conference on machine learning*, 1998, pp. 197–205.
- [24] 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.
- [25] L. Barrett and S. Narayanan, "Learning All Optimal Policies with Multiple Criteria," in *Proceedings of the international conference on machine learning*, 2008.
- [26] W. Weijia and M. Sebag, "Multi-objective Monte Carlo Tree Search," in *Proceedings of the Asian Conference on Machine Learning*, 2012, pp. 507–522.
- [27] —, "Hypervolume indicator and dominance reward based multi-objective Monte-Carlo Tree Search," *Machine Learning*, vol. 92:2–3, pp. 403–429, 2013.
- [28] 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.
- [29] T. Kozelek, "Methods of MCTS and the game Arimaa," M.S. thesis, Charles Univ., Prague, 2009.
- [30] 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.
- [31] —, "Monte Carlo Tree Search: Long Term versus Short Term

Map	Weight genome	Time	Fuel	Damage	D
Map 1	11111111111111	1654 ± 7	131 ± 2	846 ± 13	Y
	22222222222222	1657 ± 8	130 ± 2	773 ± 11	Y
	33333333333333	1681 ± 11	131 ± 2	837 ± 15	Y
	21201101220001	1619 ±	130 ±	744 ±	
Map 2	11111111111111	1409 ± 7	235 ± 4	364 ± 3	Y
	22222222222222	1402 ± 6	236 ± 5	354 ± 2	Y
	33333333333333	1416 ± 8	219 ± 4	360 ± 3	Y
	12020201212102	1390 ±	210 ±	353 ±	
Map 3	11111111111111	1373 ± 6	221 ± 3	301 ± 7	Y
	22222222222222	1378 ± 5	211 ± 3	268 ± 5	Ø
	33333333333333	1385 ± 6	203 ± 4	291 ± 7	Ø
	00001112001120	1364 ±	223 ±	269 ±	
Map 4	11111111111111	1383 ± 6	291 ± 5	565 ± 5	Ø
	22222222222222	1385 ± 7	304 ± 4	542 ± 4	Ø
	33333333333333	1423 ± 6	273 ± 4	583 ± 5	Y
	22212221201211	1408 ±	271 ±	584 ±	
Map 5	11111111111111	1405 ± 7	467 ± 4	559 ± 4	Y
	22222222222222	1431 ± 9	447 ± 4	541 ± 4	Y
	33333333333333	1467 ± 9	411 ± 5	567 ± 5	Ø
	10200102000100	1397 ±	448 ±	535 ±	
Map 6	11111111111111	1575 ± 7	549 ± 5	303 ± 4	Y
	22222222222222	1626 ± 9	540 ± 6	286 ± 5	Y
	33333333333333	1703 ± 11	499 ± 4	316 ± 7	Ø
	20010201000000	1570 ±	535 ±	266 ±	
Map 7	11111111111111	1434 ± 5	599 ± 6	284 ± 6	Y
	22222222222222	1475 ± 10	602 ± 5	243 ± 6	Y
	33333333333333	1489 ± 12	549 ± 3	264 ± 6	Ø
	00221100210221	1401 ±	563 ±	230 ±	
Map 8	11111111111111	1761 ± 9	254 ± 5	382 ± 3	Y
	22222222222222	1804 ± 10	269 ± 4	357 ± 4	Ø
	33333333333333	1826 ± 10	230 ± 3	392 ± 7	Ø
	12110020202212	1747 ±	247 ±	363 ±	
Map 9	11111111111111	2501 ± 14	926 ± 6	574 ± 9	Y
	22222222222222	2503 ± 10	921 ± 10	524 ± 8	Y
	33333333333333	2641 ± 14	833 ± 5	574 ± 14	Ø
	10021220222112	2463 ±	891 ±	523 ±	
Map 10	11111111111111	1430 ± 8	630 ± 4	205 ± 2	Ø
	22222222222222	1493 ± 13	615 ± 4	209 ± 2	Ø
	33333333333333	1542 ± 10	554 ± 4	229 ± 5	Ø
	00200000011100	1418 ±	623 ±	197 ±	

TABLE III: MO-PTSP Results with different weights: Each row corresponds to a run in the associated map. It gives results for two weights: the one below is the best after the run, while the one on the top is that from Table II with the most similar genome. Each genome is a string of the form xyz , that represents $W_x W_y W_z$, where each value is a weight used in that particular segment.

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.