

# AUPO - ABSTRACTED UNTIL PROVEN OTHERWISE: A REWARD DISTRIBUTION BASED ABSTRACTION ALGORITHM

**Robin Schmöcker**

Institute for Information Processing  
Leibniz University Hannover  
Hannover, Germany  
schmoecker@tnt.uni-hannover.de

**Alexander Dockhorn**

SDU Metaverse Lab  
University of Southern Denmark  
Odense, Denmark  
adoc@mmmi.sdu.dk

**Bodo Rosenhahn**

Institute for Information Processing  
Leibniz University Hannover  
Hannover, Germany  
rosenhahn@tnt.uni-hannover.de

## ABSTRACT

We introduce a novel, drop-in modification to Monte Carlo Tree Search’s (MCTS) decision policy that we call *AUPO*. Comparisons based on a range of IPPC benchmark problems show that AUPO clearly outperforms MCTS. AUPO is an automatic action abstraction algorithm that solely relies on reward distribution statistics acquired during the MCTS. Thus, unlike other automatic abstraction algorithms, AUPO requires neither access to transition probabilities nor does AUPO require a directed acyclic search graph to build its abstraction, allowing AUPO to detect symmetric actions that state-of-the-art frameworks like ASAP struggle with when the resulting symmetric states are far apart in state space. Furthermore, as AUPO only affects the decision policy, it is not mutually exclusive with other abstraction techniques that only affect the tree search.

## 1 INTRODUCTION

A plethora of important problems can be viewed as sequential decision-making tasks such as autonomous driving (Liu et al., 2021), energy grid optimization (Sogabe et al., 2018), financial portfolio management (Birge, 2007), or playing video games (Silver et al., 2016). Though arguably state-of-the-art on such decision-making tasks is achieved using machine learning (ML) as demonstrated by DeepMind with their AlphaGo agent for Go (Silver et al., 2016) or OpenAI Five for Dota 2 (Berner et al., 2019), there is still a demand for general domain-knowledge independent, on-the-go-applicable planning methods, properties which ML-based approaches usually lack but which are satisfied by Monte Carlo Tree Search (Browne et al., 2012) (MCTS), the method of interest for this paper. For example, Game Studios rarely implement ML agents as they have to be costly retrained whenever the game and its rules are updated. Though not within the scope of this paper, improvements to MCTS might also potentially translate to ML-based methods that use MCTS as their underlying search.

One family of approaches to improve the performance of MCTS is using abstractions that usually group similarly behaving nodes or actions of the search tree. State-of-the-art abstraction tree searches such as OGA-UCT (Anand et al., 2016) all rely on the reward function being deterministic, having full access to the transition probability of any sampled state-action pair, and on the search graph being a directed acyclic graph. While these methods could in principle still be applied if the first two conditions aren’t met (i.e. by approximating the reward and transition probabilities), they fundamentally rely on doing search on a DAG, which requires being able to check state equalities which is not always guaranteed (e.g., in memory-constrained settings where states may

only be represented as action, stochastic-outcome sequences, in partially observable domains, in continuous-state settings, or in blackbox simulation settings). Until now, no domain-independent, non-learning-based MCTS abstraction algorithms for discrete, fully-observable settings exist that have no additional constraints than MCTS, exist. This is a gap that this paper closes.

Concretely, we introduce **Abstracted Until Proven Otherwise** (AUPO), the first MCTS-based abstraction algorithm that can significantly outperform MCTS in a discrete, fully-observable, non-learning-based setting whilst requiring neither access to transition probabilities nor a directed acyclic search graph, nor a deterministic reward setting. AUPO only affects the decision policy and can thus even be combined and enhanced with other abstraction algorithms during the tree policy. Furthermore, in practice, AUPO can detect symmetric actions that the ASAP (Anand et al., 2015) framework cannot when the resulting symmetric states are far apart in state space, as ASAP needs the search graph to converge. As only the decision policy is affected, AUPO’s runtime overhead vanishes with an increase in the iteration count (see Tab. 3).

The key idea of AUPO is to consider all actions at the root node initially as equivalent, only separating them if the layerwise reward distributions, which were tracked during the MCTS search phase, differ significantly. To our knowledge, AUPO is the first abstraction algorithm to build abstractions based solely on reward distribution statistics.

The paper is structured as follows. First, in **Section 2**, we give an overview of domain knowledge independent abstraction tree searches. Next, in **Section 3** we formalize our problem setting, and lay the theoretical groundwork for understanding AUPO. This is followed by **Section 4** where we formalize AUPO and **Section 5** where we experimentally verify AUPO and discuss the experimental results. Lastly, in **Section 6** we summarise our findings and show avenues for future work.

## 2 RELATED WORK

The literature on abstraction-using planners is vast and ranges from abstractions for strategy games (Moraes & Lelis, 2018; Xu et al., 2023), card games such as Poker (Billings et al., 2003) to board games such as Go (Childs et al., 2008) or even hospital scheduling planners (Friha et al., 1997). Aside from such domain-specific abstractions, general abstraction methods are developed for continuous and/or partially observable domains (Hoerger et al., 2024) or learning-based abstractions such as learning and planning on abstract models (Ozair et al., 2021; Kwak et al., 2024; Chitnis et al., 2020), or building abstractions that rely on learned functions (e.g. a value function) (Fu et al., 2023). A comprehensive survey that focuses on non-learning-based abstraction techniques is provided by Schmöcker & Dockhorn (2025a).

The literature on non-learning-based, fully domain-independent abstraction (the scope of this paper) techniques is small. For MCTS-based planners research has heavily focused on grouping states or state-action pairs of the current MCTS search graph such that their aggregate statistics can be used to enhance the UCB formula. This is achieved by bootstrapping of state-action pair with common successors and inferring which other states or state-action pairs must consequently be equivalent (Jiang et al., 2014; Anand et al., 2015; 2016). In contrast to AUPO, however, these methods can only be applied if one has access to an equality check operator and if the state space graph is not a tree.

All of the above-mentioned techniques can be thought of as pessimistic in that they only abstract actions or states when precise conditions are met. However, in environments where equivalences are the norm and not the exception, optimistic approaches, such as AUPO, can thrive. For example, PARSS (Hostetler et al., 2015) modifies Sparse Sampling by initially grouping all successors of each state-action pair. These groups are refined by repeatedly splitting them in half as the search progresses. Like PARSS, AUPO can also be viewed as a refining and optimistic abstraction algorithm, but whereas PARSS randomly refines its abstractions when it does not have access to additional state information, AUPO does so using statistical evidence. Like PARSS and AUPO, the method of fully abandoning an abstraction mid-search (Xu et al., 2023) or dynamically dropping abstractions on a per-node level (Schmöcker et al., 2025a) can also be seen as an abstraction refinement technique. Another refining approach that is not fully domain-independent is given by Sokota et al. (2021),

who group states based on a domain-specific distance function, and the maximal grouping distance shrinks as the search progresses.

### 3 FOUNDATIONS

We use finite Markov Decision Processes (MDP) (Sutton & Barto, 2018) to formalize the sequential decision-making tasks AUPO attempts to solve.

*Definition:* An MDP is a 6-tuple  $(S, \mu_0, \mathbb{A}, \mathbb{P}, R, T)$  where the components are as follows:

- $S \neq \emptyset$  is the finite set of states,  $T \subseteq S$  is the (possibly empty) set of terminal states, and  $\mu_0 \in \Delta(S)$  is the probability distribution for the initial state where  $\Delta(X)$  denotes the probability simplex of a finite, non-empty set  $X$ .
- $\mathbb{A}: S \mapsto A$  maps each state  $s$  to the available actions  $\emptyset \neq \mathbb{A}(s) \subseteq A$  at state  $s$  where  $|A| < \infty$ .
- $\mathbb{P}: S \times A \mapsto \Delta(S)$  is the stochastic transition function where we use  $\mathbb{P}(s' | s, a)$  to denote the probability of transitioning from  $s \in S$  to  $s' \in S$  after taking action  $a \in \mathbb{A}(s)$  in  $s$ .
- $R: S \times A \mapsto \mathbb{R}$  is the reward function that maps to the set of real-valued random variables.

Starting in  $s_0 \sim \mu_0$ , an MDP progresses from state  $s_t$  to  $s_{t+1}$  by first sampling an action  $a_t \sim \pi(s_t)$  and then sampling  $s_{t+1} \sim \mathbb{P}(\cdot | s_t, a_t)$  where  $\pi$  is any agent for  $M$ . An agent  $\pi: S \mapsto \Delta(A)$  for an MDP  $M$  is a mapping from states to action distributions with  $\pi(s)(a) = 0$  for any  $a \notin \mathbb{A}(s)$ . Crucially, an agent's output depends only on a single state. At each transition  $M$  samples the reward  $r_t \sim R(s_t, a_t)$ .

In this paper, we consider only the finite horizon setting where the game ends after at most  $h \in \mathbb{N}$  steps or earlier when a terminal state is reached. We call  $h$  the horizon and any state-action-reward sequence  $(s_0, a_0, r_0), \dots, (s_n, a_n, r_n), s_{n+1}$  that can be reached a *trajectory*. If additionally  $n + 1 = h$  or  $s_{n+1} \in T$  we call this trajectory an *episode*.

The goal of any agent is to maximize its expected *return*. The return of an episode  $\tau = (s_0, a_0, r_0), \dots, (s_n, a_n, r_n), s_{n+1}$  is defined as the (possibly discounted) sum of rewards, i.e.  $R(\tau) := \sum_{i=0}^n r_i \cdot \gamma^i$  where  $0 < \gamma \leq 1$  is the *discount factor*. For any given state  $s$ , action  $a \in \mathbb{A}(s)$ , and maximum remaining steps  $k \leq h$  we call  $Q^*(s, a, k)$  the Q-value of  $(s, a, k)$  and  $V^*(s, k)$  the state value of  $s$  (given  $k$  remaining steps) which are defined as

$$Q^*(s, a, k) := \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau)], \quad (1)$$

$$V^*(s, k) := \max_{a \in \mathbb{A}(s)} Q^*(s, a, k) \quad (2)$$

where  $\tau(\pi, s, a, k)$  denotes the trajectory distribution of an agent  $\pi$  induced by starting at state  $s$ , directly applying  $a$  and then playing according to  $\pi$  for at most  $k - 1$  steps or until a terminal state is reached. We write  $Q^*(s, a) := Q^*(s, a, h)$  and  $V^*(s) := V^*(s, h)$  and  $V^* := \mathbb{E}_{s_0 \sim \mu_0} [V^*(s_0)]$ .

Our AUPO method will heavily rely on and be compared to MCTS (for a detailed description, see Section A.10). The MCTS version we employ here, uses a greedy decision policy as well as the UCB tree policy.

### 4 METHOD

**The benefit of finding abstractions to the decision policy:** One component of MCTS is its decision policy which decides which action to take given the previously obtained search tree statistics. A common decision policy and the one we employ here is the greedy strategy where one picks the root node action with the highest Q-value (sum of all returns divided by visits).

The Q-value of each (root node) action-visits pair can be viewed as a real-valued random variable. Furthermore, these random variables are independent iff their corresponding actions are different. Let us assume that there are  $n \in \mathbb{N}$  root actions in total. We denote the respective Q-value random

variables by  $Q_1, \dots, Q_n$ . For simplicity, let us assume that each root action has the same number of visits and that the optimal action is the same as  $\arg \max_{1 \leq a \leq n} \mathbb{E}[Q_a]$ .

Furthermore, let us assume that  $\mathbb{E}[Q_1] = \dots = \mathbb{E}[Q_k], k < n$ . Though consequently the actions  $a_1, \dots, a_k$  are value-equivalent they suffer from an overestimation bias in the decision policy that worsens exponentially with increasing  $k$ . The decision policy is invariant under replacing  $Q_1, \dots, Q_k$  by the random variable  $Q^m := \max(Q_1, \dots, Q_k)$ . Trivially,  $\mathbb{E}[Q^m] \geq \mathbb{E}[Q_1]$  and more concretely, it holds that for any constant  $c \in \mathbb{R}$

$$\mathbb{P}(Q^m \geq c) = 1 - \mathbb{P}(Q^m < c) = 1 - \prod_{i=1}^k \mathbb{P}(Q_i < c). \quad (3)$$

If we managed to detect that  $\mathbb{E}[Q_1] = \dots = \mathbb{E}[Q_k]$  and abstract them into a single random variable  $\bar{Q} := \frac{Q_1 + \dots + Q_k}{k}$ , then not only can the previously mentioned overestimation bias be fully mitigated but we can even decrease the variance since

$$\text{Var}(\bar{Q}) = \frac{1}{k^2} \sum_{i=1}^k \text{Var}(Q_i). \quad (4)$$

**Finding abstractions by distribution comparisons:** The main idea of AUPO is to find and utilize action abstractions at the root node during the decision policy by comparing the reward distributions at depths  $1, \dots, D$  of the game tree. Initially, AUPO assumes all actions to be equivalent, however, if the reward distributions of two actions differ significantly at any depth, the two actions are separated.

**Building the abstraction:** Let us assume we are in a state  $s \in S$  with actions  $a_1, \dots, a_n$ . After running standard MCTS for  $m$  iterations, we have sampled  $m$  trajectories where we denote the trajectories that started with action  $a_j$  by  $\tau_{i,j} = (a_{w_1}, r_1, s_1), (a_{w_2}, r_2, s_2), \dots, (a_{w_{D_{i,j}}}, r_{D_{i,j}}, (s_{D_{i,j}}))$ ,  $a_{w_1} = a_j$ ,  $1 \leq i \leq m_j$ ,  $m_1 + \dots + m_n = m$ . Consider the reward sequence  $R_{d,j}$  obtained at depth  $d$  after playing action  $a_j$  at the root node i.e.

$$((R_{d,j})_i)_{1 \leq i \leq m_j} := r_d \text{ with } (a_{w_d}, r_d, s_d) = (\tau_{i,j})_d \quad (5)$$

where we define  $r_d := 0$  in case  $D_{i,j} < d < D$ .

Though this is a heuristic assumption, we assume that all  $R_{d,j}$  are samples from a stationary distribution  $\mathcal{R}_{d,i}$  (this assumption would only hold if we performed a pure Monte Carlo search). Next, we compute the empirical mean and standard deviation for all  $\mathcal{R}_{d,j}$ ,  $d \leq D$ , along with their confidence intervals for a fixed confidence level  $q \in [0, 1]$ . Any pair of actions  $a_j, a_k$  has  $2 \cdot D$  reward distributions associated with them which are  $\mathcal{R}_{1,a_j}, \dots, \mathcal{R}_{D,a_j}$  for  $a_j$  and  $\mathcal{R}_{1,a_k}, \dots, \mathcal{R}_{D,a_k}$  for  $a_k$ . AUPO then groups  $a_j, a_k$  if and only if all confidence intervals (both the mean and std intervals) up to depth  $d \leq D$  of the pairs  $(\mathcal{R}_{d,a_j}, \mathcal{R}_{d,a_k})$  overlap. AUPO builds the standard Gaussian confidence intervals for  $q \in (0, 1)$ . For  $q = 1$ , the interval is defined as  $(-\infty, \infty)$  and for  $q = 0$  the interval is equal to the singleton set containing only the empirical mean of the quantity of interest. If any confidence interval pair does not overlap, then  $a_j, a_k$  are separated. Note that this induces a symmetric and reflexive but not necessarily transitive relation over the root actions.

Optionally, to ensure that in the limit, AUPO does not group non-value-equivalent actions, we may additionally separate two actions, if the distribution of their returns differs significantly (in the sense that their mean and standard deviation confidence intervals do not overlap). The return of a trajectory is the (possibly discounted) sum of all its rewards. We call this option the return filter  $RF \in \{0, 1\}$ . Analogously, whether AUPO uses the standard deviation confidence intervals for distribution separation is denoted by the std filter  $SF \in \{0, 1\}$ .

**Using the abstraction:** We use the abstractions during the decision policy only. AUPO transforms the decision policy into a two-step process. In the first step, we assign each action  $a_j$  its abstract Q-value which is the sum of the returns divided by the sum of the visits of all actions  $a_j$  is grouped with. We select the action  $a^*$  that maximizes the abstract Q-value. Ties are broken randomly. In the second step, we select the action inside the abstraction of  $a^*$  with the highest unabstracted/ground Q-value. This decision policy makes AUPO a generalization of the greedy decision policy as for both  $q \in \{0, 1\}$  AUPO's decision policy degenerates to the greedy

policy. While for  $q = 0$  step two becomes redundant, for  $q = 1$  step one becomes redundant. We summarize AUPO in the Appendix in Alg. 1.

**Theoretical guarantees:** First, without any further assumptions, it can be shown that the root state abstraction AUPO is sound in the iteration limit, i.e., it only groups state-action pairs with the same  $Q^*$  value. The proof of the following is found in Section A.1 of the appendix.

*Theorem:* Assume that MCTS has been run on a state  $s$  for  $m$  iterations and let  $a^{\text{left}}, a^{\text{right}} \in \mathbb{A}(s)$  be two legal actions at  $s$  with  $Q^*(s, a^{\text{left}}) \neq Q^*(s, a^{\text{right}})$ . It then holds that

$$\lim_{m \rightarrow \infty} \mathbb{P}[(s, a^{\text{left}}), (s, a^{\text{right}}) \text{ is abstracted by AUPO with RF} = 1] = 0. \quad (6)$$

The key innovation that makes AUPO work in practice (this will be shown empirically later) is that one does not only compare a single pair of distributions to differentiate a single action pair but rather one compares a number of distributions induced by that action pair. Using some simplifying assumptions, one can show that with an increase in  $D$ , the order of the number of samples required to differentiate two non-equivalent actions changes. The following theorem formalizes this and its proof is found in the appendix Section A.2.

*Theorem:* Again, assume that MCTS has been run on a state  $s$  for and let  $a^{\text{left}}, a^{\text{right}} \in \mathbb{A}(s)$  be two legal actions at  $s$  that both have been played  $n$  times. The following assumptions are made:

1. It is assumed that all MCTS trajectories prior to performing the AUPO abstraction have been generated with a uniformly random tree policy. This ensures that for a fixed depth and root action, the obtained rewards are independent samples from a stationary distribution, a necessary requirement for the following result to hold.
2. All layerwise reward distributions  $\mathcal{R}_{d, a^{\text{left}}}, \mathcal{R}_{d, a^{\text{right}}}, d \in \{1, \dots, D\}$  are assumed to be independently distributed and Gaussians with means  $m^{\text{left}} = (m_1^{\text{left}}, \dots, m_D^{\text{left}})$  and  $m^{\text{right}} = (m_1^{\text{right}}, \dots, m_D^{\text{right}})$  and standard deviations  $\sigma^{\text{left}} = (\sigma_1^{\text{left}}, \dots, \sigma_D^{\text{left}})$ ,  $\sigma^{\text{right}} = (\sigma_1^{\text{right}}, \dots, \sigma_D^{\text{right}})$ .
3. Lastly, it is assumed that AUPO has oracle access to these standard deviations and uses them instead of the empirical standard deviation when constructing the confidence intervals for the means.

Under these assumptions, if AUPO uses neither the return, nor the std-filter then

$$\forall \varepsilon > 0 : \mathbb{P}[\text{AUPO abstracts } a^{\text{left}} \text{ and } a^{\text{right}}] \in \mathcal{O}(f(n)), f(n) = e^{-n \cdot (\varepsilon + \sum_{k=1}^D w_i)} \quad (7)$$

where for  $1 \leq i \leq D$ :  $w_i = \begin{cases} \frac{(\mu_i^{\text{left}} - \mu_i^{\text{right}})^2}{2(\sigma_i^{\text{left}} + \sigma_i^{\text{right}})^2}, & |\mu_i^{\text{left}} - \mu_i^{\text{right}}| \geq \frac{z^*}{\sqrt{n}}(\sigma_i^{\text{left}} + \sigma_i^{\text{right}}) \\ 1, & \text{otherwise} \end{cases}$ , and  $z^*$  is the

critical value of the standard normal distribution for  $q$  (e.g.  $z^* \approx 1.96$  for  $q = 0.95$ ).

**AUPO example:** Next, we illustrate on an instance of the IPPC problem SysAdmin how AUPO detects abstractions. A detailed explanation of this problem is given in the experiment Appendix A.8. Assume we are in a state where all computers, except one outer computer, are online. This is visualized in Fig. 1a. This state features exactly four value-equivalent action types. Idling, rebooting the offline computer (machine 3), rebooting (even though it is still online) the hub computer (machine 0), or rebooting any outer running computer (machines 1-2,5-9). Given enough trajectory samples, AUPO separates and subsequently detects these equivalences as follows.

*Idle action:* All actions except idling have the same immediate reward, the reboot cost. Therefore, the idle action is easily separated by considering only the mean of the 1-step reward distribution.

*Rebooting the offline computer:* This action can be separated from the others by the 2-step reward distribution, as it takes one step for the computer to be rebooted and then another step to receive the reward from the additional running computer. Though a little noisy, the 2-step reward will be on average 1 higher than that of the other actions. We quantitatively verified this in the Appendix in Tab. 2.

*Rebooting the hub computer:* This action can be separated from rebooting any of the outer running computers by the standard deviation of the 3-step reward. If we reboot the hub computer, we safeguard it from randomly crashing in the next step, which prevents the catastrophe where numerous other computers fail in the next step as they are connected to the then-broken hub computer. This scenario happens only rarely but when it does happen it is catastrophic, thus causing the 3-step reward of not rebooting the hub to have a relatively high variance compared to rebooting and thus protecting it. We quantitatively verified this in the Appendix in Tab. 1.

*Rebooting an outer running computer:* Since they are symmetric and thus have identical reward distributions at all downstream steps, AUPO optimistically assumes that are equivalent and thus abstract into a single action.

**Relation to other abstraction frameworks** In practice, AUPO is able to detect abstractions that ASAP could not because the latter requires the state graph to converge on states from which the abstraction building can be bootstrapped. Hence it is practically impossible for ASAP to detect equivalences that arise due to symmetry. For example, while it would be no problem for AUPO to detect that saving any of the four corner cells is equivalent in the Game of Life state visualized in Fig. 1b, ASAP would not be able to detect this with feasible computational resources. Game of Life is defined in the Appendix A.8.

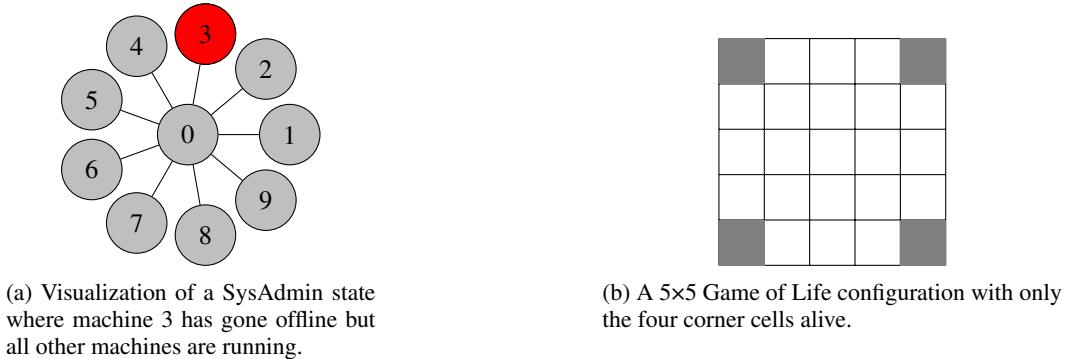


Figure 1: Visualization of two environments considered in this paper.

Furthermore, ASAP struggles with a high stochastic branching factor. While AUPO is able to detect that rebooting any of the outer machines from the SysAdmin example in Section 4 is equivalent, ASAP is not able to detect these equivalences if two equivalent actions have not sampled the exact same set of successors from which there are  $33554432 = 2^{25}$ .

## 5 EXPERIMENTS

In this section, we will present the setup and results for the comparison of AUPO with MCTS showing that AUPO is the first and currently only tree-search abstraction algorithm that does neither require access to the transition probabilities, nor the model having deterministic rewards, nor requires a directed acyclic search graph but can outperform MCTS.

**Problem models:** The problem models that we tested AUPO on are either problems from the International Conference on Probabilistic Planning (IPPC) (Grzes et al., 2014) or appear throughout the literature. For the readers not familiar with these problem models, we give a high-level overview in the appendix in Section A.8. For details and the concrete instances, i.e. model parameter choices, we refer to our publicly available implementation (Schmöcker & Dockhorn, 2025b), which is the translation into C++ of the Relational Dynamic Influence Diagram (RDDI) (Sanner, 2011) descriptions of these environments found at the RDDL repository of Taitler et al. (2022). These environments were deliberately chosen as they appear throughout the abstraction literature (Anand et al., 2015; 2016; Hostetler et al., 2015; Yoon et al., 2008; Jiang et al., 2014) or have been used for planning competitions (Grzes et al., 2014), feature value-equivalent sibling actions, dense rewards, two theoretically necessary requirements for AUPO to yield any

performance increase.

**Experiment setup and reproducibility:** For every experiment, we used a horizon of 50 episode steps. Since we are in the finite-horizon setting, we used a discount of  $\gamma = 1$ . We ran every experiment for at least 2000 episodes, and whenever we denote the mean return of this experiment we additionally provide a 99% confidence interval. We denote the confidence interval of any quantity by its mean and the half of the interval size, e.g. we would denote a return confidence interval  $(1, 3)$  by  $2 \pm 1$ . For both MCTS and AUPO, we performed random playouts until the episode terminates. Additionally, as the problem models vary in their reward scale, we used the dynamic exploration factor Global Std (Schmöcker et al., 2025b) that is given by  $C \cdot \sigma$  where  $\sigma$  is the empirical standard deviation of all Q values of the current search tree and  $C \in \mathbb{R}^+$  is a parameter. For reproducibility, we released our implementation (Schmöcker & Dockhorn, 2025b). Our code was compiled with g++ version 13.1.0 using the -O3 flag (i.e. aggressive optimization).

**Parameter-optimized performances:** First, we tested whether and in which environments AUPO can increase the parameter-optimized performance over MCTS. To do this, we considered the best AUPO performance when varying the parameters exploration constant  $C \in \{0.5, 1, 2, 4, 8, 16\}$ , distribution tracking depth  $D \in \{1, 2, 3, 4\}$ , using the return filter SF  $\in \{0, 1\}$ , using the return filter RF  $\in \{0, 1\}$ , and varying the confidence level  $q \in \{0.8, 0.9, 0.95, 0.99\}$ . Furthermore, since the standard UCB tree policy results in non-uniformly distributed visits, we also considered AUPO’s performance when using a uniform root policy (denoted as U-AUPO) which has two main effects. Firstly, each action, even those that UCB would not exploit, receive visits, thus shrinking their confidence intervals, making them easier to separate from other actions. And secondly, we reduce the risk of separating reward distribution equivalent actions because in MCTS the distributions shift with an increasing visit count as MCTS starts to exploit.

We compare AUPO and U-AUPO to the performance of MCTS and MCTS with a uniform root policy U-MCTS, as well as RANDOM-ABS that is the same as AUPO except that for each action pair they are randomly abstracted at the decision policy with the probability  $p_{\text{random}} \in \{0.1, 0.2, \dots, 0.9\}$ . Hence, RANDOM-ABS is equivalent to MCTS in the cases  $p_{\text{random}} \in \{0, 1\}$ . RANDOM-ABS verifies that the abstractions found by AUPO outperform randomly formed abstractions. Do reduce the amount of visuals; any RANDOM-ABS data points are simply the maximum of both RANDOM-ABS with a uniform root policy and standard root policy. The parameter-optimized performances in dependence of the iteration number are visualized in Fig. 2. The following key observations can be made:

- 1) AUPO can gain a clear performance **advantage** over MCTS (and RANDOM-ABS) in **11 out of the 14 here-considered environments**, in at least one iteration budget. In the environments, Academic Advising, Game of Life, Multi-armed bandit, Push Your Luck, Cooperative Recon, SysAdmin, and Traffic, AUPO maintains a clear performance edge for the majority of iteration budgets.
- 2) Expectedly, U-MCTS mostly performs worse than MCTS, however, the performance improvements between U-MCTS and U-AUPO is mostly significantly greater than the gap between MCTS and AUPO, showing the AUPO as suggested benefits from uniformly distributed visits. Notably, there is an environment, namely Cooperative Recon in which MCTS and U-MCTS perform evenly, where however, U-AUPO clearly outperforms AUPO. Also, in Saving both U-MCTS and U-AUPO outperform their non-uniform counterparts. Hence, using a uniform root policy can be a tool to improve the peak performance.

**Generalization capabilities:** Next, we test AUPO’s generalization capabilities. For this, we computed the pairings and relative improvement scores for all AUPO, U-AUPO, MCTS, U-MCTS, and RANDOM-ABS parameter combinations. These scores are Borda-like rankings of individual parameter-combinations and both lie in the interval  $[-1, 1]$  (1 is the best value and -1 the worst) and they are formalized in the Appendix Section A.11. The results for all iteration budgets and environments combined are visualized in Fig. 3 and show that the best performances with respect to both scores with large margins, are reached with AUPO. These results are qualitatively identical for each iteration budget which is presented in the Appendix Section A.12.

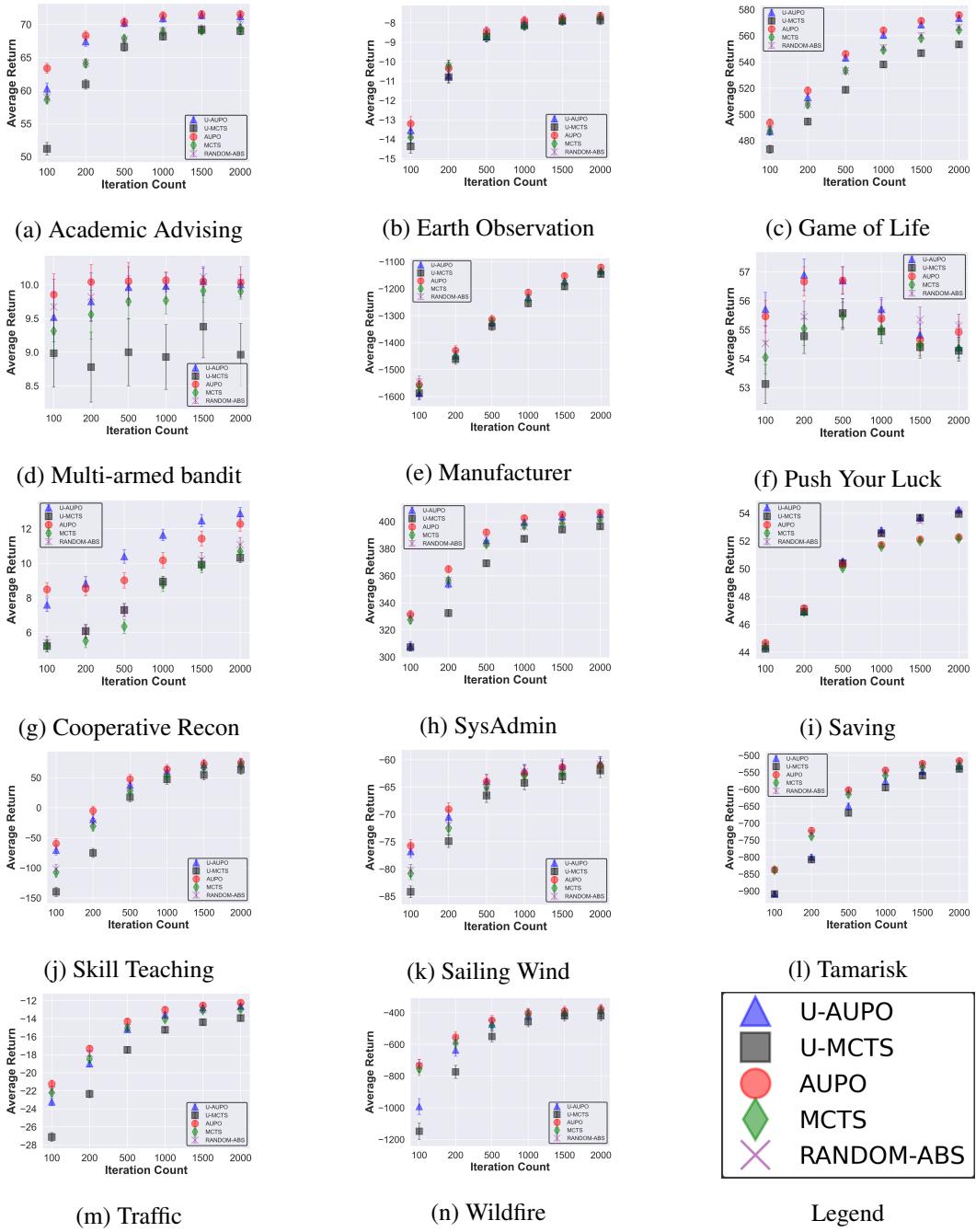
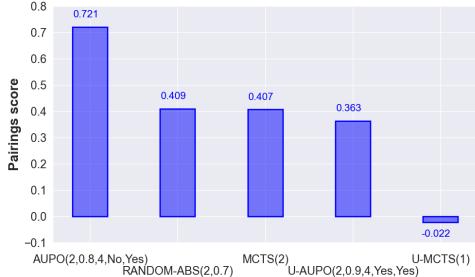
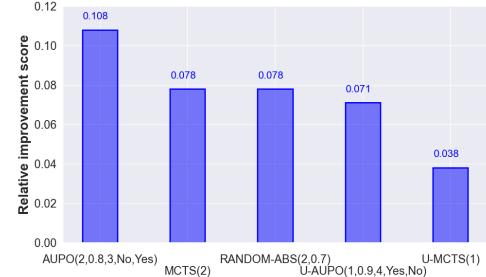


Figure 2: The performance graphs of in dependence of the MCTS iteration count of the parameter optimized versions of AUPO, MCTS, and RANDOM-ABS. The prefix U- denotes AUPO and MCTS using a uniform root policy.



(a) Pairings score



(b) Relative improvement score

Figure 3: The pairings and relative improvement scores across all environments and iteration budgets for different AUPO, U-AUPO (parameter format  $(C, q, D, RF, SF)$ ), MCTS, U-MCTS (parameter format  $(C)$ ), and RANDOM-ABS (parameter format  $(C, p_{random})$ ) agents. The bar charts show the top score reached by each agent type as well as the parameter combination to reach that score. In the case of RANDOM-ABS the score was reached with the standard root policy.

**Ablations:** Lastly, we are going to study the impact of the individual parameters. Instead of displaying the best-performing parameter set, we fix the parameter in question and max over the remaining parameters. With only a few exception such as Multi-armed bandit, the confidence level does only have a significant impact for low iteration counts if it has any impact at all. In the low iteration count regime, lower confidences generally outperform higher confidence levels. Depending on the environment, the impact of the distribution tracking depth can be significant to non-existent. In cases where it does matter, high depths are always preferred with the only exception being Push Your Luck. Filters can be extremely beneficial to some environments, such as Multi-armed Bandit or Wildfire Whilst not causing any harm to the environment where it has little impact. We visualize the concrete performance values for this ablation in the Appendix Fig. 5 which shows the results when varying the confidence level, in Fig. 4, which shows the results when varying the distribution tracking depth, and in Fig. 6 that shows the results when varying either the std filter or the return filter.

## 6 LIMITATIONS AND FUTURE WORK

In this paper, we introduced a novel action abstraction algorithm that we call AUPO which only affects the decision policy of MCTS. We could experimentally show that AUPO outperforms MCTS in a wide range of environments that contain states with value-equivalent sibling actions. Though AUPO introduces four new parameters, their choice mostly has only a minor impact on performance.

First and foremost, for AUPO to achieve any performance gain, the environment must contain state-action pairs with the same parent that have similar  $Q^*$  values, i.e. there need to be abstractions to be detected in the first place. Another key limitation of AUPO is that it is reliant on dense-rewards. For example, in binary-outcome zero-sum two-player games AUPO would have a hard time distinguishing actions, as only the return distribution can be used for differentiation. How this limitation can be overcome, is left as future work. Another weakness of AUPO is that it requires many visits for the distributions to be distinguishable; hence it cannot be used in low iteration settings and therefore not during the tree policy. Therefore, another area for future work is how to make AUPO much more sensitive to be able to deal with low iterations. Furthermore, for future work, as mentioned in the introduction, it could also be of interest to combine AUPO with other abstraction algorithms. For example, one may use state-of-the-art such as OGA-UCT (Anand et al., 2016) during the search phase, replacing only the decision policy with AUPO. In its current form, AUPO uses the same confidence level for each layer. However, it might be worth investigating if additional performance can be achieved by making this parameter layer-dependent.

## REFERENCES

- Giuseppe Abreu. Very simple tight bounds on the  $q$ -function. *IEEE Trans. Commun.*, 60(9):2415–2420, 2012. doi: 10.1109/TCOMM.2012.080612.110075. URL <https://doi.org/10.1109/TCOMM.2012.080612.110075>.
- Ankit Anand, Aditya Grover, Mausam, and Parag Singla. ASAP-UCT: Abstraction of State-Action Pairs in UCT. In Qiang Yang and Michael J. Wooldridge (eds.), *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pp. 1509–1515. AAAI Press, 2015. URL <http://ijcai.org/Abstract/15/216>.
- Ankit Anand, Ritesh Noothigattu, Mausam, and Parag Singla. OGA-UCT: on-the-go abstractions in UCT. In *Proceedings of the Twenty-Sixth International Conference on International Conference on Automated Planning and Scheduling*, ICAPS’16, pp. 29–37. AAAI Press, 2016. ISBN 1577357574.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Christopher Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with Large Scale Deep Reinforcement Learning. *CoRR*, abs/1912.06680, 2019. URL <http://arxiv.org/abs/1912.06680>.
- Darse Billings, Neil Burch, Aaron Davidson, Robert C. Holte, Jonathan Schaeffer, Terence Schauenberg, and Duane Szafron. Approximating Game-Theoretic Optimal Strategies for Full-scale Poker. In Georg Gottlob and Toby Walsh (eds.), *IJCAI-03, Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, Acapulco, Mexico, August 9-15, 2003*, pp. 661–668. Morgan Kaufmann, 2003. URL <http://ijcai.org/Proceedings/03/Papers/097.pdf>.
- John R. Birge. Chapter 20 Optimization Methods in Dynamic Portfolio Management. In John R. Birge and Vadim Linetsky (eds.), *Financial Engineering*, volume 15 of *Handbooks in Operations Research and Management Science*, pp. 845–865. Elsevier, 2007. doi: [https://doi.org/10.1016/S0927-0507\(07\)15020-9](https://doi.org/10.1016/S0927-0507(07)15020-9). URL <https://www.sciencedirect.com/science/article/pii/S0927050707150209>.
- Cameron Browne, Edward Jack Powley, Daniel Whitehouse, Simon M. Lucas, Peter I. Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez Liebana, Spyridon Samothrakis, and Simon Colton. A Survey of Monte Carlo Tree Search Methods. *IEEE Trans. Comput. Intell. AI Games*, 4(1):1–43, 2012. doi: 10.1109/TCIAIG.2012.2186810. URL <https://doi.org/10.1109/TCIAIG.2012.2186810>.
- Benjamin E. Childs, James H. Brodeur, and Levente Kocsis. Transpositions and move groups in Monte Carlo tree search. In Philip Hingston and Luigi Barone (eds.), *Proceedings of the 2008 IEEE Symposium on Computational Intelligence and Games, CIG 2009, Perth, Australia, 15-18 December, 2008*, pp. 389–395. IEEE, 2008. doi: 10.1109/CIG.2008.5035667. URL <https://doi.org/10.1109/CIG.2008.5035667>.
- Rohan Chitnis, Tom Silver, Beomjoon Kim, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Camps: Learning context-specific abstractions for efficient planning in factored mdps. In Jens Kober, Fabio Ramos, and Claire J. Tomlin (eds.), *4th Conference on Robot Learning, CoRL 2020, 16-18 November 2020, Virtual Event / Cambridge, MA, USA*, volume 155 of *Proceedings of Machine Learning Research*, pp. 64–79. PMLR, 2020. URL <https://proceedings.mlr.press/v155/chitnis21a.html>.
- Lamia Friha, P. Berry, and Berthe Y Choueiry. DISA: A Distributed scheduler using abstractions. *Revue d’Intelligence Artificielle*, 11, 01 1997.
- Yangqing Fu, Ming Sun, Bujing Nie, and Yue Gao. Accelerating monte carlo tree search with probability tree state abstraction. In Alice Oh, Tristan Naumann, Amir

Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023*. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/bf89c9fc0ef605571a03666f6a6a44d-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/bf89c9fc0ef605571a03666f6a6a44d-Abstract-Conference.html).

Marek Grzes, Jesse Hoey, and Scott Sanner. International Probabilistic Planning Competition (IPPC) 2014. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, 2014.

Marcus Hoerger, Hanna Kurniawati, Dirk P. Kroese, and Nan Ye. Adaptive discretization using voronoi trees for continuous pomdps. *Int. J. Robotics Res.*, 43(9):1283–1298, 2024. doi: 10.1177/02783649231188984. URL <https://doi.org/10.1177/02783649231188984>.

Jesse Hostetler, Alan Fern, and Thomas G. Dietterich. Progressive Abstraction Refinement for Sparse Sampling. In Marina Meila and Tom Heskes (eds.), *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence, UAI 2015, July 12-16, 2015, Amsterdam, The Netherlands*, pp. 365–374. AUAI Press, 2015. URL <http://auai.org/uai2015/proceedings/papers/81.pdf>.

Nan Jiang, Satinder Singh, and Richard L. Lewis. Improving UCT planning via approximate homomorphisms. In Ana L. C. Bazzan, Michael N. Huhns, Alessio Lomuscio, and Paul Scerri (eds.), *International conference on Autonomous Agents and Multi-Agent Systems, AAMAS '14, Paris, France, May 5-9, 2014*, pp. 1289–1296. IFAAMAS/ACM, 2014. URL <http://dl.acm.org/citation.cfm?id=2617453>.

Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In Johannes Fürnkranz, Tobias Scheffer, and Myra Spiliopoulou (eds.), *Machine Learning: ECML 2006, 17th European Conference on Machine Learning, Berlin, Germany, September 18-22, 2006, Proceedings*, volume 4212 of *Lecture Notes in Computer Science*, pp. 282–293. Springer, 2006. doi: 10.1007/11871842\\_29. URL [https://doi.org/10.1007/11871842\\_29](https://doi.org/10.1007/11871842_29).

Volodymyr Kuleshov and Doina Precup. Algorithms for multi-armed bandit problems. *CoRR*, abs/1402.6028, 2014. URL <http://arxiv.org/abs/1402.6028>.

Yunhyeok Kwak, Inwoo Hwang, Dooyoung Kim, Sanghack Lee, and Byoung-Tak Zhang. Efficient monte carlo tree search via on-the-fly state-conditioned action abstraction. In Negar Kiyavash and Joris M. Mooij (eds.), *Uncertainty in Artificial Intelligence, 15-19 July 2024, Universitat Pompeu Fabra, Barcelona, Spain*, volume 244 of *Proceedings of Machine Learning Research*, pp. 2076–2093. PMLR, 2024. URL <https://proceedings.mlr.press/v244/kwak24a.html>.

Qi Liu, Xueyuan Li, Shihua Yuan, and Zirui Li. Decision-Making Technology for Autonomous Vehicles: Learning-Based Methods, Applications and Future Outlook. In *24th IEEE International Intelligent Transportation Systems Conference, ITSC 2021, Indianapolis, IN, USA, September 19-22, 2021*, pp. 30–37. IEEE, 2021. doi: 10.1109/ITSC48978.2021.9564580. URL <https://doi.org/10.1109/ITSC48978.2021.9564580>.

Rubens O. Moraes and Levi H. S. Lelis. Asymmetric Action Abstractions for Multi-Unit Control in Adversarial Real-Time Games. In Sheila A. McIlraith and Kilian Q. Weinberger (eds.), *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pp. 876–883. AAAI Press, 2018. doi: 10.1609/AAAI.V32I1.11432. URL <https://doi.org/10.1609/aaai.v32i1.11432>.

Sherjil Ozair, Yazhe Li, Ali Razavi, Ioannis Antonoglou, Aäron van den Oord, and Oriol Vinyals. Vector quantized models for planning. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8302–8313. PMLR, 2021. URL <https://proceedings.mlr.press/v139/ozair21a.html>.

- Scott Sanner. Relational dynamic influence diagram language (rddl): Language description. 01 2011. [https://users.cecs.anu.edu.au/~ssanner/IPPC\\_2011/RDDL.pdf](https://users.cecs.anu.edu.au/~ssanner/IPPC_2011/RDDL.pdf) Accessed: 22-01-2025.
- Robin Schmöcker and Alexander Dockhorn. A survey of non-learning-based abstractions for sequential decision-making. *IEEE Access*, 13:100808–100830, 2025a. doi: 10.1109/ACCESS.2025.3572830.
- Robin Schmöcker and Alexander Dockhorn. Aupo, 2025b. Repository available at: <https://github.com/codebro634/Aupo>.
- Robin Schmöcker, Lennart Kampmann, and Alexander Dockhorn. Time-Critical and Confidence-Based Abstraction Dropping Methods. In *2025 IEEE Conference on Games (CoG)*, 2025a. doi: 10.1109/CoG64752.2025.11114261.
- Robin Schmöcker, Christoph Schnell, and Alexander Dockhorn. Investigating scale independent uct exploration factor strategies, 2025b. URL <https://arxiv.org/abs/2510.21275>.
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P. Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nat.*, 529(7587):484–489, 2016. doi: 10.1038/NATURE16961. URL <https://doi.org/10.1038/nature16961>.
- Tomoh Sogabe, Dinesh Bahadur Malla, Shota Takayama, Seiichi Shin, Katsuyoshi Sakamoto, Koichi Yamaguchi, Thakur Praveen Singh, Masaru Sogabe, Tomohiro Hirata, and Yoshitaka Okada. Smart Grid Optimization by Deep Reinforcement Learning over Discrete and Continuous Action Space. In *2018 IEEE 7th World Conference on Photovoltaic Energy Conversion (WCPEC) (A Joint Conference of 45th IEEE PVSC, 28th PVSEC and 34th EU PVSEC)*, pp. 3794–3796, 2018. doi: 10.1109/PVSC.2018.8547862.
- Samuel Sokota, Caleb Ho, Zaheen Farraz Ahmad, and J. Zico Kolter. Monte Carlo Tree Search With Iteratively Refining State Abstractions. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 18698–18709, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/9b0ead00a217ea2c12e06a72eec4923f-Abstract.html>.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, 2nd edition, 2018.
- Ayal Taitler, Michael Gimelfarb, Sriram Gopalakrishnan, Xiaotian Liu, and Scott Sanner. pyrddl-gym: From RDDL to gym environments. *CoRR*, abs/2211.05939, 2022. doi: 10.48550/ARXIV.2211.05939. URL <https://doi.org/10.48550/arXiv.2211.05939>.
- Linjie Xu, Alexander Dockhorn, and Diego Perez-Liebana. Elastic Monte Carlo Tree Search. *IEEE Transactions on Games*, 15(4):527–537, 2023. doi: 10.1109/TG.2023.3282351.
- Sung Wook Yoon, Alan Fern, Robert Givan, and Subbarao Kambhampati. Probabilistic Planning via Determinization in Hindsight. In Dieter Fox and Carla P. Gomes (eds.), *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence, AAAI 2008, Chicago, Illinois, USA, July 13-17, 2008*, pp. 1010–1016. AAAI Press, 2008. URL <http://www.aaai.org/Library/AAAI/2008/aaai08-160.php>.

## A APPENDIX

### A.1 PROOF OF AUPO SOUNDNESS

In this section, Equation 6 from Section 4 is proven.

*Theorem:* Assume that MCTS has been run on a state  $s$  for  $m$  iterations and let  $a^{\text{left}}, a^{\text{right}} \in \mathbb{A}(s)$  be two legal actions at  $s$  with  $Q^*(s, a^{\text{left}}) \neq Q^*(s, a^{\text{right}})$ . It then holds that

$$\lim_{m \rightarrow \infty} \mathbb{P}[(s, a^{\text{left}}), (s, a^{\text{right}}) \text{ is abstracted by AUPO with RF} = 1] = 0. \quad (8)$$

*Proof:* In the iteration limit, the Q values of  $(s, a^{\text{right}})$  and  $(s, a^{\text{left}})$  converge in probability to their corresponding  $Q^*$  values. Furthermore, since our MDP model definition encompasses only finite state-action pair sets, the reward function and consequently the set of possible episode returns are bounded and thus the empirical standard deviation is also bounded. Therefore, the lengths of the Gaussian Q value confidence intervals converge to 0. And since the confidence intervals are centered at their respective Q values, the probability of them overlapping converges to zero. Consequently, the probability of the state-action pairs being abstracted by AUPO also converges to zero.  $\square$

### A.2 PROOF OF ABSTRACTION PROBABILITY THEOREM

In this section, Equation 7 from Section 4 is proven which is the following.

*Theorem:* Again, assume that MCTS has been run on a state  $s$  for and let  $a^{\text{left}}, a^{\text{right}} \in \mathbb{A}(s)$  be two legal actions at  $s$  that both have been played  $n$  times. The following assumptions are made:

1. It is assumed that all MCTS trajectories prior to performing the AUPO abstraction have been generated with a uniformly random tree policy. This ensures that for a fixed depth and root action, the obtained rewards are independent samples from a stationary distribution, a necessary requirement for the following result to hold.
2. All layerwise reward distributions  $\mathcal{R}_{d, a^{\text{left}}}, \mathcal{R}_{d, a^{\text{right}}}, d \in \{1, \dots, D\}$  are assumed to be independently distributed and Gaussians with means  $m^{\text{left}} = (m_1^{\text{left}}, \dots, m_D^{\text{left}})$  and  $m^{\text{right}} = (m_1^{\text{right}}, \dots, m_D^{\text{right}})$  and standard deviations  $\sigma^{\text{left}} = (\sigma_1^{\text{left}}, \dots, \sigma_D^{\text{left}})$ ,  $\sigma^{\text{right}} = (\sigma_1^{\text{right}}, \dots, \sigma_D^{\text{right}})$ .
3. Lastly, it is assumed that AUPO has oracle access to these standard deviations and uses them instead of the empirical standard deviation when constructing the confidence intervals for the means.

Under these assumptions, if AUPO uses neither the return, nor the std-filter then

$$\forall \varepsilon > 0 : \mathbb{P}[\text{AUPO abstracts } a^{\text{left}} \text{ and } a^{\text{right}}] \in \mathcal{O}(f(n)), f(n) = e^{-n \cdot (\varepsilon + \sum_{k=1}^D w_i)} \quad (9)$$

where for  $1 \leq i \leq D$ :  $w_i = \begin{cases} \frac{(\mu_i^{\text{left}} - \mu_i^{\text{right}})^2}{2(\sigma_i^{\text{left}} + \sigma_i^{\text{right}})^2}, & |\mu_i^{\text{left}} - \mu_i^{\text{right}}| \geq \frac{z^*}{\sqrt{n}}(\sigma_i^{\text{left}} + \sigma_i^{\text{right}}) \\ 1, & \text{otherwise} \end{cases}$ , and  $z^*$  is the

critical value of the standard normal distribution for  $q$  (e.g.  $z^* \approx 1.96$  for  $q = 0.95$ ).

*Proof:* Firstly, we will derive a general upper bound for the probability of confidence intervals overlapping and then use this result in the context of AUPO's abstraction mechanism. **1)** Let  $n \in \mathbb{N}$  and  $X_1, \dots, X_n, Y_1, \dots, Y_n$  be i.i.d. Gaussian random variables with respective means and stds of  $\mu_X \geq \mu_Y$  and  $\sigma_X, \sigma_Y$ . For any confidence level  $q \in [0, 1]$ , the confidence interval for  $\mu_X$  (analogously  $\mu_Y$ ) is of the form

$$[\bar{X} \pm \frac{z^* \cdot \sigma_X}{\sqrt{n}}] \quad (10)$$

where  $z^* \in \mathbb{R}$  is the z-score for the given confidence level  $q$  and  $\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k$  ( $\bar{Y}$  is defined analogously). The probability that the confidence intervals for  $\mu_X$  and  $\mu_Y$  overlap is thus given by

$$\mathbb{P}[\underbrace{|\bar{X} - \bar{Y}|}_{Z :=} \leq \underbrace{\frac{z^*}{\sqrt{n}}(\sigma_X + \sigma_Y)}_{T :=}]. \quad (11)$$

Since  $Z$  is Gaussian and the mean of  $Z$  is  $\mu_Z := \mu_X - \mu_Y$  and the std is  $\sigma_Z := \frac{\sigma_X + \sigma_Y}{\sqrt{n}}$ , and since  $\mathbb{P}[|Z| \leq T] = \mathbb{P}[Z \leq T] - \mathbb{P}[Z \leq -T]$  one obtains

$$\mathbb{P}[|Z| \leq T] = \frac{1}{2} \left[ \operatorname{erf} \left( \frac{T + \mu_Z}{\sqrt{2}\sigma_Z} \right) + \operatorname{erf} \left( \frac{T - \mu_Z}{\sqrt{2}\sigma_Z} \right) \right] \quad (12)$$

using the identity  $\Phi(\frac{x-\mu}{\sigma}) = \frac{1}{2}(1 + \operatorname{erf}(\frac{x-\mu}{\sigma\sqrt{2}}))$  that holds for any Gaussian with mean  $\mu$  and std  $\sigma$  where  $\Phi$  is the CDF for the standard Gaussian distribution and  $\operatorname{erf}$  is the Gauss error function. Next, using that  $\operatorname{erf}$  is an odd function with range  $(-1, 1)$ , yields

$$\mathbb{P}[|Z| \leq T] = \frac{1}{2} \left[ -\operatorname{erfc} \left( \frac{\mu_Z + T}{\sqrt{2}\sigma_Z} \right) + \operatorname{erfc} \left( \frac{\mu_Z - T}{\sqrt{2}\sigma_Z} \right) \right] \leq \frac{1}{2} \operatorname{erfc} \left( \frac{\mu_Z - T}{\sqrt{2}\sigma_Z} \right) \text{ with } \operatorname{erfc} := 1 - \operatorname{erf}. \quad (13)$$

Next, two cases are differentiated. If  $\mu_Z - T < 0$ , we simply bound  $\mathbb{P}[|Z| \leq T]$  by 1. In the other case,  $\mu_Z - T \geq 0$ , one can use an upper bound derived by Giuseppe Abreu (Abreu, 2012) to further estimate this expression in terms of the exponential function. Concretely this yields,

$$\frac{1}{2} \operatorname{erfc} \left( \frac{\mu_Z - T}{\sqrt{2}\sigma_Z} \right) \leq \frac{1}{50} e^{-x^2} + \frac{1}{2(x+1)} e^{-x^2/2} \leq e^{-x^2/2}, x = \frac{\mu_Z - T}{\sigma_Z}, \quad (14)$$

which is a function of the form

$$e^{-\tilde{\lambda}_1 + \tilde{\lambda}_2 \sqrt{n} - w \cdot n}, \text{ with } w = \frac{(\mu_X - \mu_Y)^2}{2(\sigma_X + \sigma_Y)^2}; \tilde{\lambda}_1, \tilde{\lambda}_2 \in \mathbb{R}^+. \quad (15)$$

**2)** By definition, AUPO using no return or std filter with a distribution tracking depth  $D$  only abstracts  $a^{\text{left}}$  and  $a^{\text{right}}$  iff their mean confidence intervals up to depth  $D$  all overlap. Since all reward distributions are independent by assumption, the probability of all confidence intervals overlapping, is given as the product of the individual ones overlapping. Hence, we can use the previously obtained results about a single pair of confidence intervals to obtain the following for every  $\varepsilon > 0$

$$\mathbb{P}[\text{AUPO abstracts } a^{\text{left}} \text{ and } a^{\text{right}}] \leq e^{-\lambda_1 + \sqrt{n} \cdot \lambda_2 - n \cdot \sum_{k=1}^D w_i} \in \mathcal{O}(f(n)), \lambda_1, \lambda_2 \in \mathbb{R}^+, \quad (16)$$

where  $f(n) = e^{-n \cdot (\varepsilon + \sum_{k=1}^D w_i)}$  and for  $1 \leq i \leq D$ :

$$w_i = \begin{cases} \frac{(\mu_i^{\text{left}} - \mu_i^{\text{right}})^2}{2(\sigma_i^{\text{left}} + \sigma_i^{\text{right}})^2}, & |\mu_i^{\text{left}} - \mu_i^{\text{right}}| \geq \frac{z^*}{\sqrt{n}} (\sigma_i^{\text{left}} + \sigma_i^{\text{right}}) \\ 1, & \text{otherwise} \end{cases} \quad (17)$$

This proves the original statement.  $\square$

### A.3 AUPO PSEUDOCODE

---

**Algorithm 1:** AUPO

---

**Parameters:**  $q$ ,  $D$ ,  $filter\_std$ ,  $filter\_return$ ,  $mcts\_args$

**Input:**  $state$

// Run MCTS and collect reward distribution data

1  $n = \text{num\_actions}(state)$ ,  $R[d, j] = [] \forall d, j$

2 **for**  $i = 1 \dots mcts\_iterations$  **do**

3   sample MCTS trajectory with rewards  $r_1, \dots, r_{D^*}$  and first action  $a_j$

4   **for**  $d = 1 \dots D$  **do**

5     |  $R[d, j].append(r_d \text{ if } d \leq D^* \text{ else } 0)$

6   **end**

7   |  $R^*[j].append(r_1 + \dots + r_{\min(D, D^*)})$

8 **end**

// Compute confidence intervals

9 **for**  $j = 1 \dots n$  **do**

10   **for**  $d = 1 \dots D$  **do**

11     |  $mean\_interval[d, j] = \text{mean\_conf\_interval}(R[d, j], q)$

12     |  $std\_interval[d, j] = \text{std\_conf\_interval}(R[d, j], q)$

13   **end**

14   |  $return\_mean\_interval[j] = \text{mean\_conf\_interval}(R^*[j], q)$

15   |  $return\_std\_interval[j] = \text{std\_conf\_interval}(R^*[j], q)$

16 **end**

// Compute abstractions

17 **for**  $i = 1 \dots n$  **do**

18   |  $abstract\_visits = 0, abstract\_value = 0$

19   |  $abstraction[i] = \{\}$

20   **for**  $j = 1 \dots n$  **do**

21     |  $abstracted = \text{true}$

22     **for**  $d = 1 \dots D$  **do**

23       |  $\text{if } mean\_interval[d, j] \cap mean\_interval[d, i] == \emptyset \text{ or } filter\_std \text{ and }$

24       |  $std\_interval[d, j] \cap std\_interval[d, i] == \emptyset \text{ then}$

25       |   |  $abstracted = \text{false}$

26       **end**

27       **if**  $filter\_return$  **and** (

28       |  $return\_mean\_interval[d, j] \cap return\_mean\_interval[d, i] == \emptyset \text{ or } filter\_std \text{ and }$

29       |  $return\_std\_interval[d, j] \cap return\_std\_interval[d, i] == \emptyset \text{ then}$

30       |   |  $abstracted = \text{false}$

31       **end**

32       **if**  $abstracted$  **then**

33         |  $abstract\_visits += action\_visits(j)$

34         |  $abstract\_value += action\_returns(j)$

35         |  $abstraction[i].insert(j)$

36       **end**

37       |  $abstract\_Q[i] = \frac{abstract\_value}{abstract\_visits}$

38 **end**

// Action selection

39  $abs\_action = \arg \max_{i=1 \dots n} abstract\_Q[i]$

40  $ground\_action = \arg \max_{i \in abstraction[abs\_action]} Q[i]$

41 **return**  $ground\_action$ ;

---

#### A.4 ABLATION: DISTRIBUTION TRACKING DEPTH $D$

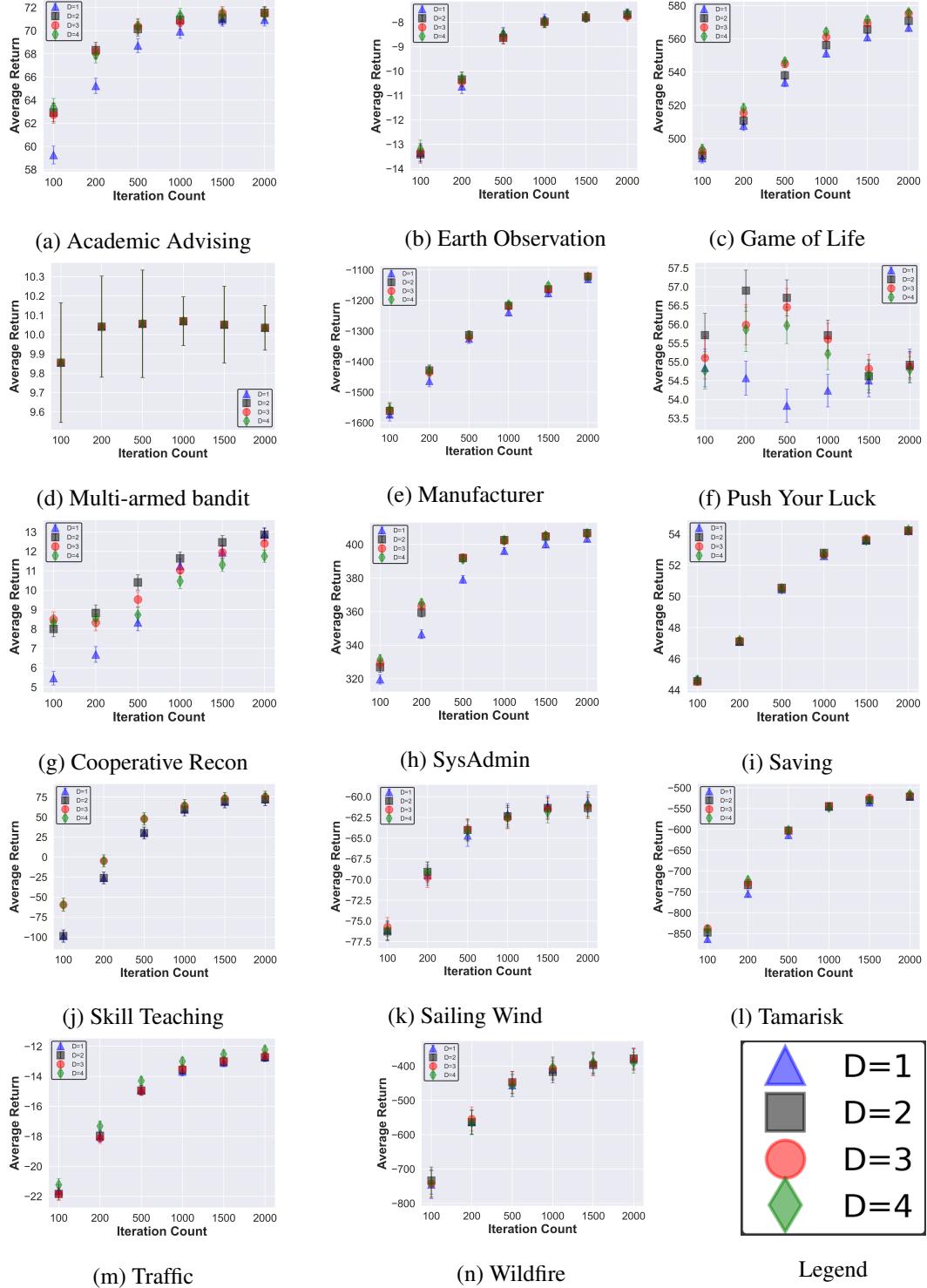


Figure 4: The performance graphs of in dependence of the MCTS iteration count of the parameter optimized versions of AUPO using different fixed values for the distribution tracking depth  $D$ .

### A.5 PERFORMANCES IN DEPENDENCE OF THE CONFIDENCE LEVEL $q$

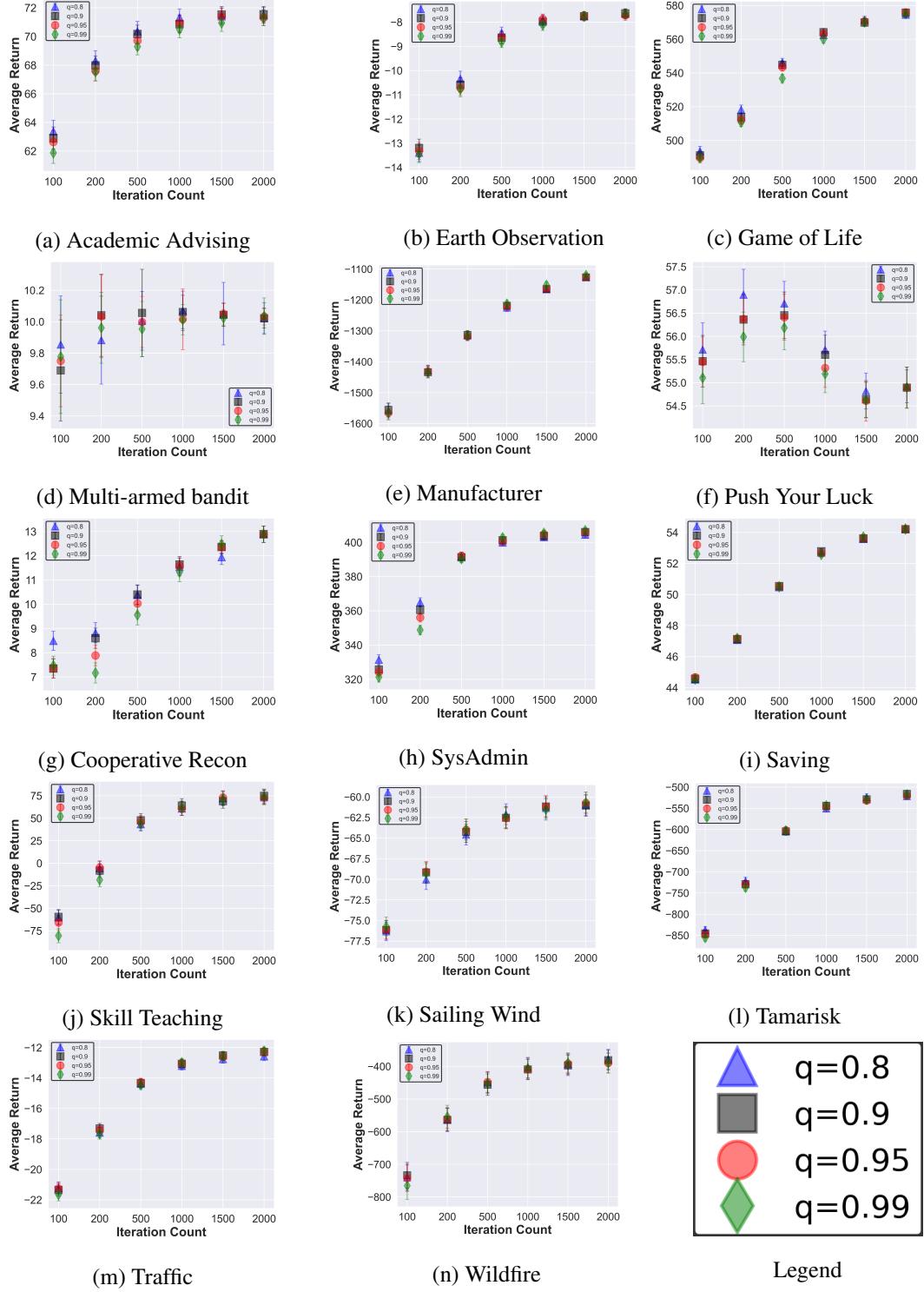


Figure 5: The performance graphs of in dependence of the MCTS iteration count of the parameter optimized versions of AUPO using different fixed values for the confidence  $q$ .

## A.6 PERFORMANCES WHEN USING DIFFERENT FILTER COMBINATIONS

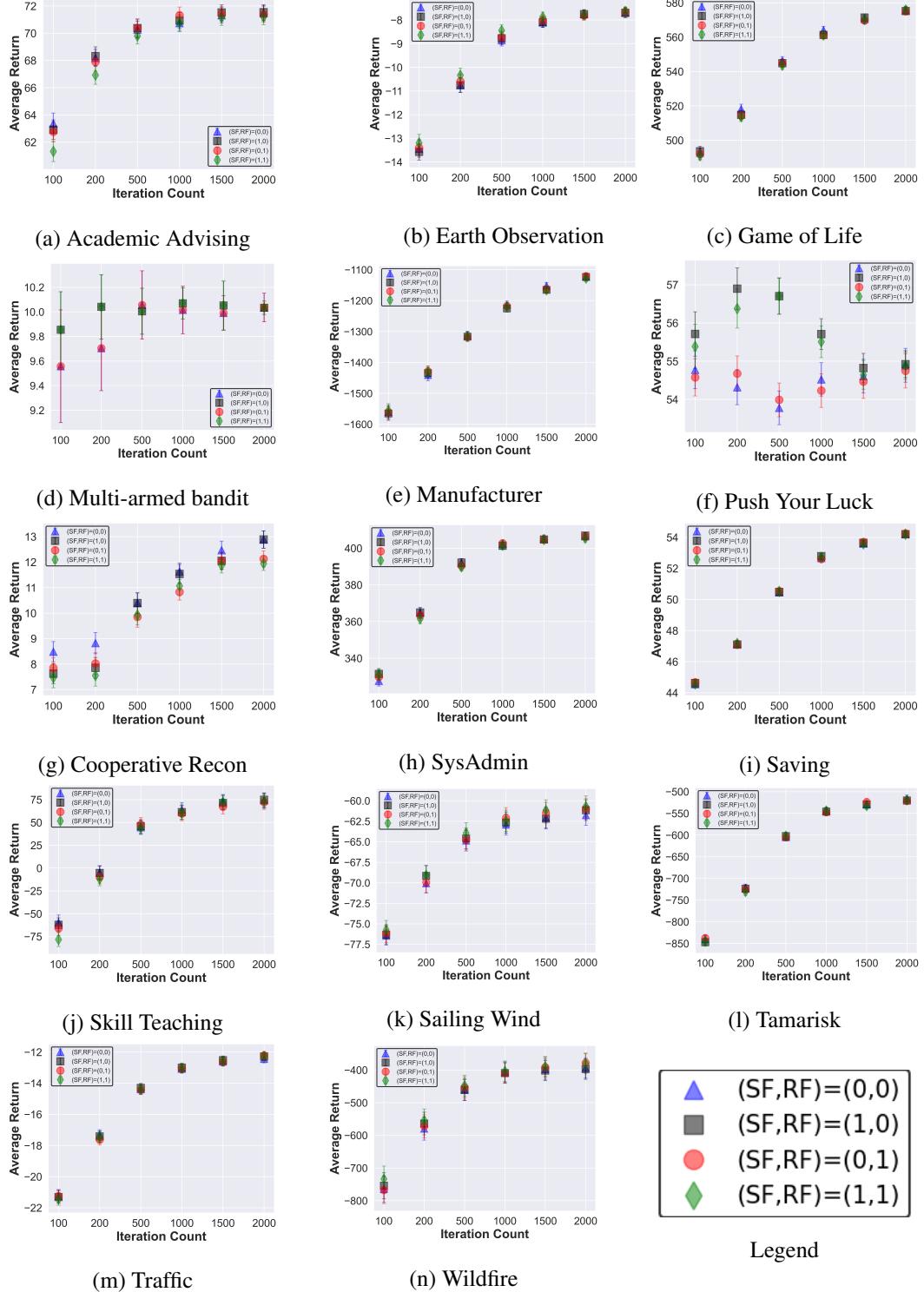


Figure 6: The performance graphs of in dependence on the MCTS iteration count of the parameter optimized versions of AUPO using different fixed filter settings. Both the return (RF) and std filter (SF) are varied.

### A.7 REWARD DISTRIBUTION CONFIDENCE INTERVALS FOR SYSADMIN

Table 1: 3-step standard deviation 95% confidence intervals for the reward distribution after a different number of MCTS iterations on the SysAdmin state of Fig. 1a. Note that with higher iteration counts, rebooting the hub can be separated from the remaining actions

Action	1000 iterations	2000 iterations	3000 iterations	4000 iterations
Hub (0)	(0.89, 1.08)	(0.91, 1.04)	(0.92, 1.03)	(0.92, 1.02)
1	(1.14, 1.38)	(1.15, 1.32)	(1.18, 1.31)	(1.21, 1.33)
2	(1.13, 1.37)	(1.16, 1.32)	(1.19, 1.32)	(1.19, 1.31)
3	(1.16, 1.40)	(1.16, 1.33)	(1.17, 1.31)	(1.21, 1.33)
4	(1.11, 1.34)	(1.15, 1.32)	(1.18, 1.32)	(1.19, 1.30)
5	(1.14, 1.38)	(1.17, 1.34)	(1.17, 1.31)	(1.19, 1.31)
6	(1.15, 1.39)	(1.17, 1.34)	(1.19, 1.32)	(1.18, 1.30)
7	(1.12, 1.35)	(1.16, 1.33)	(1.19, 1.33)	(1.19, 1.31)
8	(1.12, 1.36)	(1.16, 1.33)	(1.19, 1.33)	(1.20, 1.32)
9	(1.11, 1.35)	(1.17, 1.34)	(1.18, 1.32)	(1.19, 1.31)
Idle	(1.12, 1.36)	(1.19, 1.36)	(1.19, 1.33)	(1.21, 1.33)

Table 2: 2-step mean 95% confidence intervals for the reward distribution after different numbers of MCTS iterations on the SysAdmin state of Fig. 1a. Note that even with very low iteration counts, rebooting machine 3 can easily be separated from the other actions.

Action	200 iterations	500 iterations	1000 iterations	2000 iterations
Hub (0)	(7.72, 8.16)	(7.83, 8.10)	(7.88, 8.06)	(7.90, 8.03)
1	(7.67, 8.11)	(7.76, 8.03)	(7.81, 8.00)	(7.85, 7.98)
2	(7.71, 8.14)	(7.76, 8.04)	(7.81, 8.00)	(7.84, 7.97)
3	(8.62, 9.04)	(8.70, 8.97)	(8.73, 8.92)	(8.74, 8.88)
4	(7.68, 8.13)	(7.77, 8.04)	(7.82, 8.01)	(7.85, 7.98)
5	(7.72, 8.16)	(7.78, 8.05)	(7.83, 8.02)	(7.85, 7.99)
6	(7.67, 8.12)	(7.76, 8.04)	(7.81, 8.00)	(7.84, 7.97)
7	(7.69, 8.14)	(7.78, 8.05)	(7.83, 8.02)	(7.84, 7.98)
8	(7.68, 8.13)	(7.78, 8.05)	(7.81, 8.01)	(7.84, 7.98)
9	(7.68, 8.12)	(7.76, 8.04)	(7.82, 8.01)	(7.83, 7.97)
Idle	(7.64, 8.10)	(7.74, 8.02)	(7.77, 7.97)	(7.80, 7.94)

### A.8 PROBLEM MODELS

In the following, we provide a brief description of each domain/environment that was used in this paper. Some of these environments can be parametrized (e.g., choosing a concrete map size for Sailing Wind). The concrete parameter settings can be found in the *ExperimentConfigs* folder in our publicly available GitHub repository (Schmöcker & Dockhorn, 2025b). In the survey paper (Schmöcker & Dockhorn, 2025a) as well as in Schmöcker et al. (2025b), descriptions for most of the environments considered here can be found. For a detailed description of these environments, we refer to our implementation. In the following, descriptions for the environments that are not contained in the previous two papers are given.

**Academic Advising:** Though this problem is described by Schmöcker & Dockhorn (2025a), we use a modified version in this paper. Originally, the agent would also always receive a negative reward as long as there is one mandatory course that has not been passed. We increased the reward density, by letting this negative reward be dependent on the number of missing mandatory courses. Furthermore, we also added a reward for every course passed.

**Multi-armed bandit:** Multi-armed bandits (MAB) (Kuleshov & Precup, 2014) are 1-step MDPs. Each action  $1 \leq a \leq n$  is called an arm, and its execution yields an immediate random reward sampled from the probability distribution associated with the  $a$ -th arm. We use Gaussians as the reward distributions. All actions whose associated arms have the same mean are equivalent. We deliberately chose a MAB instance with a high number of equivalences.

### A.9 RUNTIME MEASUREMENTS

We validate the claim that AUPO adds only a minor runtime overhead over vanilla MCTS for high iteration budgets, the following table, Tab. 3 lists the average decision-making times for each environment of AUPO compared to MCTS for 100 and 2000 iterations on states sampled from a distribution induced by random walks. This shows that while AUPO adds a significant overhead for low iteration budgets, the impact of the decision policy and therefore AUPO’s runtime overhead vanishes. Note, though, that this runtime is both heavily implementation and hardware-dependent, and more efficient implementations might reduce this overhead. In particular, we are using highly optimized environment implementations that could be the runtime bottleneck in more complex environments.

Table 3: Average decision-making times of AUPO and MCTS in milliseconds for 100 and 2000 iterations. For AUPO the most computational heavy version has been used, which uses  $p = 0.8$ ,  $D = 4$ , the return- and std filter. This data was obtained using an Intel(R) Core(TM) i5-9600K CPU @ 3.70GHz. The data shows a median runtime overhead of  $\approx 8\%$  for 100 iterations and  $\approx 4\%$  for 2000 iterations.

Domain	AUPO-100	MCTS-100	AUPO-2000	MCTS-2000
Academic Advising	1.15	1.59	23.71	25.36
Cooperative Recon	2.41	2.57	52.25	54.21
Earth Observation	7.03	6.95	130.57	136.73
Game of Life	3.93	4.31	65.72	65.18
Manufacturer	10.31	9.82	185.76	186.81
Sailing Wind	1.90	2.01	34.34	35.15
Saving	0.82	0.87	17.41	18.24
Skills Teaching	2.42	2.61	52.20	53.91
SysAdmin	1.20	1.58	22.97	24.29
Tamarisk	2.78	2.88	47.84	48.95
Traffic	3.05	3.85	61.23	63.74
Triangle Tireworld	1.25	1.32	25.43	27.22
Push Your Luck	2.26	2.47	43.21	45.38
Multi-armed bandit	0.16	1.16	3.13	4.33
Wildfire	1.48	2.20	34.22	35.73

### A.10 MONTE CARLO TREE SEARCH

AUPO heavily relies on Monte Carlo Tree Search (MCTS) which we are going to describe now. Let  $M$  be a finite horizon MDP. On a high level, MCTS repeatedly samples trajectories starting at some state  $s_0 \in S$  where a decision has to be made until a stopping criterion is met. The final decision is then chosen as the action at  $s_0$  with the highest average return. In contrast to a pure Monte Carlo search, MCTS improves subsequent trajectories by building a tree from a subset of the states encountered in the last iterations which is then exploited. In contrast to pure Monte Carlo search, MCTS is guaranteed to converge to the optimal action.

An MCTS search tree is made of two components. Firstly, the state nodes, that represent states and Q nodes that represent state action pairs. Each state node, saves only its children which are a set of Q nodes. Q nodes save both its children which are state nodes and the number of and the sum of the returns of all trajectories that were sampled starting at the Q node.

Initially, the MCTS search tree consists only of a single state node representing  $s_0$ . Until some stopping criterion is met, the following steps are repeated.

1. **Selection phase:** Starting at the root node, MCTS first selects a Q node according to the so-called *tree policy*, which may use the nodes’ statistics, and then samples one of the Q node’s successor states. If either a terminal state node, a state node with at least one non-visited action (partially expanded), or a new Q node successor state is sampled, the selection phase ends.

A commonly used tree policy (**and the one we used**) that is synonymously used with MCTS is Upper Confidence Trees (UCT) (Kocsis & Szepesvári, 2006) which selects an

action that maximizes the Upper Confidence Bound (UCB) value. Let  $s \in S$  and  $V_a, N_a$  with  $a \in \mathbb{N}$  be the return sum and visits and of the Q nodes of the node representing  $s$ . The UCB value of any action  $a$  is then given by

$$\text{UCB}(a) = \underbrace{\frac{V_a}{N_a}}_{\text{Q term}} + \lambda \underbrace{\sqrt{\frac{\log \left( \sum_{a' \in \mathbb{A}(s)} N_{a'} \right)}{N_a}}}_{\text{Exploration term}}. \quad (18)$$

The exploration term quantifies how much the Q term could be improved if this Q node was fully exploited and is controlled by the exploration constant  $\lambda \in \mathbb{R} \cup \{\infty\}$ . If one chose  $\lambda = 0$ , the UCT selection policy becomes the greedy policy and for  $\lambda = \infty$ , the selection policy becomes a uniform policy over the visits. In case of equality, some tiebreak rule has to be selected, which is typically a random tiebreak. From here, will use MCTS and UCT (MCTS with UCB selection formula) synonymously.

2. **Expansion:** Unless the selection phases ended in a terminal state node, the search tree is expanded by a single node. In case the selection phase ended in a partially expanded state node, then one unexpanded action is selected (e.g. randomly, or according to some rule), the corresponding Q node is created and added as a child and one successor state of that Q node is sampled and added as a child to the new Q node. If the selection phase ended because a new successor of a Q node was sampled, then a state node representing this new state is added as a child to that Q node.
3. **Rollout/Simulation phase:** Starting at the state  $s_{\text{rollout}}$  of the newly added state node of the expansion phase (or at a terminal state node reached by the selection phase), actions according to the *rollout policy* are repeatedly selected and applied to  $s_{\text{rollout}}$  until a terminal state is reached. All states encountered during this phase are not added to the search tree.
4. **Backpropagation:** In this phase, the statistics of all Q nodes that were part of the last sampled trajectory that corresponds to a path in the search tree are updated by incrementing their visit count and adding the trajectory's return (of the trajectory starting at the respective Q node) to their return sum statistic.

Once the MCTS search tree has been built (by reaching an iteration limit in our case) and statistics have been gathered, the final decision is made by the *decision policy* that in our MCTS version simply chooses the action with the highest final Q value.

#### A.11 DEFINITION OF RELATIVE IMPROVEMENT AND PAIRINGS SCORE

In the main experimental section, we evaluated AUPO with respect to the relative improvement and pairings score, which are formalized here. This pairings score was also used in Schmöcker et al. (2025b). While the pairings score is calculated by summing over the number of tasks where some agent performed better than another, the relative improvement score also takes the percentage of the improvement into account; however, it is prone to outliers. Hence, we considered both scores to paint the full picture.

**Definition:** Concretely, let  $\{\pi_1, \dots, \pi_n\}$  be  $n$  agents (e.g., concrete parameter settings for possibly different base algorithms such as AUPO or MCTS) where each agent was evaluated on  $m$  tasks (in this paper, a task will always be a given MCTS iteration budget and an environment) where  $p_{i,k} \in \mathbb{R}$  denotes the performance of agent  $\pi_i$  on the  $k$ -th task. The *pairings score matrix*  $M^{\text{pairings}} \in \mathbb{R}^{n \times n}$  is defined as

$$M_{i,j}^{\text{pairings}} = \frac{1}{m-1} \sum_{1 \leq k \leq m} \text{sgn}(p_{i,k} - p_{j,k}) \quad (19)$$

where  $\text{sgn}$  is the signum function. The *pairings score*  $s_i^{\text{pairings}}, i \leq n$  is given by

$$s_i^{\text{pairings}} = \frac{1}{n-1} \sum_{1 \leq l \leq n, l \neq i} M_{i,l}^{\text{pairings}}. \quad (20)$$

The *relative improvement matrix*  $M^{\text{rel}} \in \mathbb{R}^{n \times n}$  is defined as

$$M_{i,j}^{\text{rel}} = \frac{1}{m-1} \sum_{1 \leq k \leq m} \frac{p_{i,k} - p_{j,k}}{\max(|p_{i,j}|, |p_{j,k}|)} \quad (21)$$

and the *relative improvement score*  $s_i^{\text{rel}}, i \leq n$  is given by

$$s_i^{\text{rel}} = \frac{1}{n-1} \sum_{1 \leq l \leq n, l \neq i} M_{i,l}^{\text{rel}}. \quad (22)$$

## A.12 PAIRINGS AND RELATIVE IMPROVEMENT SCORES

Table 4: The pairings and relative improvement scores for the **100, 200, and 500 iterations** setting for the parameters combination of AUPO, U-AUPO, RANDOM-ABS, MCTS, and U-MCTS with the highest respective scores as well as the concrete parameters used to reach that score. The parameters and environments used to obtain these scores are the same as the experiments of Section 5. The parameter format for AUPO and U-AUPO is  $(C, q, D, RF, SF)$ , the format RANDOM-ABS is  $(C, p_{\text{random}})$ , and for both MCTS and U-MCTS is  $(C)$ . For RANDOM-ABS the best scores are obtained using the standard root policy.

100 iterations relative improvement score.

Parameters	Score
AUPO(2,0.8,4,No,No)	0.120
U-AUPO(16,0.8,4,No,No)	0.068
RANDOM-ABS(2,0.9)	0.055
MCTS(2)	0.049
U-MCTS(4)	-0.006

200 iterations relative improvement score.

Parameters	Score
AUPO(2,0.8,3,No,Yes)	0.137
RANDOM-ABS(1,0.8)	0.073
U-AUPO(0.5,0.9,4,Yes,No)	0.068
MCTS(2)	0.062
U-MCTS(1)	-0.005

500 iterations relative improvement score.

Parameters	Score
AUPO(2,0.9,4,Yes,No)	0.139
RANDOM-ABS(2,0.4)	0.109
U-AUPO(0.5,0.9,3,Yes,No)	0.107
MCTS(2)	0.105
U-MCTS(0.5)	0.069

100 iterations pairings score.

Parameters	Score
AUPO(2,0.8,3,No,Yes)	0.742
RANDOM-ABS(2,0.7)	0.360
U-AUPO(1,0.8,4,Yes,No)	0.319
MCTS(2)	0.301
U-MCTS(0.5)	-0.099

200 iterations pairings score.

Parameters	Score
AUPO(2,0.8,3,Yes,Yes)	0.826
RANDOM-ABS(2,0.9)	0.438
MCTS(2)	0.368
U-AUPO(0.5,0.8,4,Yes,No)	0.330
U-MCTS(0.5)	-0.020

500 iterations pairings score.

Parameters	Score
AUPO(2,0.9,4,Yes,Yes)	0.793
U-AUPO(2,0.9,4,Yes,Yes)	0.459
MCTS(2)	0.431
RANDOM-ABS(1,0.7)	0.417
U-MCTS(0.5)	-0.007

Table 5: The pairings and relative improvement score for the **1000, 1500, and 2000 iterations** setting for the parameters combination of AUPO, U-AUPO, RANDOM-ABS, MCTS, and U-MCTS with the highest respective score as well as the concrete parameters used to reach that score. The parameters and environments used to obtain these scores are the same as the experiments of Section 5. The parameter format for AUPO and U-AUPO is  $(C, q, D, RF, SF)$ , the format RANDOM-ABS is  $(C, p_{\text{random}})$ , and for both MCTS and U-MCTS is  $(C)$ . For RANDOM-ABS the best scores are obtained using the standard root policy.

1000 iterations relative improvement score.

Parameters	Score
AUPO(2,0.9,2,Yes,Yes)	0.108
RANDOM-ABS(2,0.9)	0.089
U-AUPO(1,0.9,4,Yes,Yes)	0.089
MCTS(2)	0.084
U-MCTS(1)	0.057

1500 iterations relative improvement score.

Parameters	Score
AUPO(2,0.99,2,Yes,Yes)	0.099
MCTS(2)	0.084
RANDOM-ABS(2,0.8)	0.083
U-AUPO(1,0.8,4,Yes,Yes)	0.083
U-MCTS(1)	0.061

2000 iterations relative improvement score.

Parameters	Score
AUPO(2,0.95,4,Yes,Yes)	0.098
RANDOM-ABS(2,0.7)	0.087
MCTS(2)	0.086
U-AUPO(1,0.99,4,Yes,Yes)	0.086
U-MCTS(1)	0.059

1000 iterations pairings score.

Parameters	Score
AUPO(2,0.9,4,Yes,Yes)	0.770
U-AUPO(1,0.9,4,Yes,Yes)	0.499
RANDOM-ABS(2,0.9)	0.447
MCTS(2)	0.417
U-MCTS(1)	0.036

1500 iterations pairings score.

Parameters	Score
AUPO(2,0.9,4,Yes,Yes)	0.763
U-AUPO(2,0.9,4,Yes,Yes)	0.532
MCTS(2)	0.438
RANDOM-ABS(2,0.7)	0.418
U-MCTS(1)	0.026

2000 iterations pairings score.

Parameters	Score
AUPO(2,0.95,4,Yes,Yes)	0.753
U-AUPO(2,0.95,4,Yes,Yes)	0.538
RANDOM-ABS(2,0.8)	0.532
MCTS(2)	0.487
U-MCTS(1)	0.028