

Optimized look-ahead tree policies: a bridge between look-ahead tree policies and direct policy search

Tobias Jung^{1,*,\dagger}, Louis Wehenkel¹, Damien Ernst¹ and Francis Maes^{1,2}

¹Montefiore Institute, Department of Electrical Engineering and Computer Science, University of Liège, Liège, Belgium
²Department of Computer Science, Katholieke Universiteit Leuven, Leuven, Belgium

SUMMARY

Direct policy search (DPS) and look-ahead tree (LT) policies are two popular techniques for solving difficult sequential decision-making problems. They both are simple to implement, widely applicable without making strong assumptions on the structure of the problem, and capable of producing high-performance control policies. However, computationally, both of them are, each in their own way, very expensive. DPS can require huge offline resources (effort required to obtain the policy) to first select an appropriate space of parameterized policies that works well for the targeted problem and then to determine the best values of the parameters via global optimization. LT policies do not require any offline resources; however, they typically require huge online resources (effort required to calculate the best decision at each step) in order to grow trees of sufficient depth. In this paper, we propose optimized LTs (OLTs), a model-based policy learning scheme that lies at the intersection of DPS and LT. In OLT, the control policy is represented indirectly through an algorithm that at each decision step develops, as in LT by using a model of the dynamics, a small LT until a prespecified online budget is exhausted. Unlike LT, the development of the tree is not driven by a generic heuristic; rather, the heuristic is optimized for the target problem and implemented as a parameterized node scoring function learned offline via DPS. We experimentally compare OLT with pure DPS and pure LT variants on optimal control benchmark domains. The results show that the LT-based representation is a versatile way of compactly representing policies in a DPS scheme (which results in OLT being easier to tune and having lower offline complexity than pure DPS) and at the same time DPS helps to significantly reduce the size of the LTs that are required to take high-quality decisions (which results in OLT having lower online complexity than pure LT). Moreover, OLT produces overall better performing policies than pure DPS and pure LT, and also results in policies that are robust with respect to perturbations of the initial conditions. Copyright © 2013 John Wiley & Sons, Ltd.

Received 9 January 2012; Revised 23 December 2012; Accepted 9 January 2013

KEY WORDS: reinforcement learning; optimal control; direct policy search; look-ahead tree search

1. INTRODUCTION

In the last decades, a wide range of techniques have been developed for solving optimal sequential decision-making problems that arise in the ideologically separated fields of industrial control, reinforcement learning (RL), and AI. Two simple but powerful techniques that are known to work well particularly on difficult problems characterized by large state spaces are *direct policy search* (DPS) [1] and *look-ahead tree* (LT) policies [2].

Direct policy search techniques primarily originate from the field of RL and partly also evolutionary computation. In DPS, optimization is carried out directly in policy space. Control policies are represented as parameterized functions mapping states to actions. The best policy, that is, the best setting of parameters, and sometimes also the functional form itself, is found by local or

*Correspondence to: Tobias Jung, Montefiore Institute, Department of Electrical Engineering and Computer Science, University of Liège, Liège, Belgium.

[†]E-mail: tjung@ulg.ac.be

global optimization over the parameter space, where the objective function—the performance of the induced policy—is estimated through simulation or real-world experiments.

Look-ahead tree policies originate from the field of planning and search, but are increasingly also gaining popularity as a model-based technique in the field of optimal control. Unlike problems studied under the ‘classic’ RL setting, problems in optimal control are typically deterministic, and it is assumed that a model is available that enables to simulate transitions and rewards. LT policies rely on this model during decision making and determine control actions in the following way. First, they develop an LT by simulating the evolution of the system subject to multiple possible future action sequences; then they select an action on the basis of the information collected in this tree.

For practitioners, both DPS and LT methods are very appealing. Both methods are general and do not make any strong assumptions about the nature of the problem (e.g., real-valued vs. unstructured discrete spaces) or its specific properties (regularity assumptions on transitions and rewards), which makes them applicable to a wide range of real-world problems. Both methods are, from the standpoint of technical sophistication required, fairly simple and can be implemented rapidly and easily. Finally, both methods are surprisingly effective and can produce good policies for difficult control problems where more sophisticated methods (such as value function or Pontryagin’s maximum principle based techniques) would otherwise fail.

However, DPS and LT have one major weakness: they both are, each in their own way, computationally very expensive. With DPS, the computational cost arises from the initial effort required to learn the policy. A policy in DPS can be parametrized in many ways (e.g., linear parametrizations, neural networks, or radial basis functions) and will typically also depend on additional hyperparameters that determine their complexity (e.g., the number of hidden neurons and the number of basis functions). There is no simple recipe to determine which particular class of parametrization works best for a given problem, and the final performance of the policy crucially depends on making the right choice. Finding an appropriate representation, tuning its hyperparameters, and then determining the values of the best parameters can thus be a challenging and time consuming process and typically involves a large amount of trial-and-error experimentation. However, once a suitable parameterization is found and its parameters are learned via global optimization, executing the policy, that is, using it to calculate decisions online, is very fast and only requires negligible computational resources. LT policies on the other hand are the exact opposite. LT policies do not require any offline efforts; they can be applied to almost any optimal control problem without any prior knowledge—this without even depending on the dimensionality of the state space. The main issue with LT policies is that they typically require a huge amount of online resources to grow a tree of sufficient depth to derive high-quality decisions from. In practice, these online resources are often limited, because of either real-time constraints (e.g., 100 ms per decision) or human-scale constraints (e.g., 1 h or 1 day per decision). Given such constraints, the actions suggested by LT policies may be arbitrarily sub-optimal.

We thus make the distinction in our evaluation of methods between *offline computational complexity* (effort required to obtain the policy) and *online computational complexity* (effort required to calculate the best decision at each step) and find that DPS and LT lie at opposite ends of the complexity spectrum: DPS typically requires huge offline resources to first select an appropriate space of parameterized policies that works well for the targeted problem and then to determine the best values of the parameters via global optimization. However, the online complexity of DPS can be very low. LT policies do not require any offline resources; however, they typically require huge online resources in order to grow trees of sufficient depth.

Starting from this observation, it is natural to wonder how we can combine these two ideas in order to create a new hybrid solution for optimal control that leverages the advantages of both.

In this paper, we propose to bridge the gap between LT and DPS policies with a new model-based hybrid technique: *optimized LT (OLT) policies*.[‡] From the point of view of LT policies, our approach

[‡]This technique was initially proposed in [3] and is here extended through a more mature exposition of the method and an extensive experimental study, which covers many more aspects than the initial paper. Furthermore, we introduce the use of Gaussian process optimization (GPO) to improve the sample efficiency of the approach.

starts from the observation that the efficiency of an LT policy, given a finite online computational budget, can be improved by more intelligently directing the way in which the tree is grown, for example, by giving priority to sequences of actions that seem more promising [2]. In fact, a large number of different strategies could be considered to develop LTs, by combining aspects from breadth-first exploration to depth-first exploration and from random exploration to purely greedy exploration. In OLT policies, the strategy that is used to develop the tree is not fixed *a priori*, but is learned in a problem-dependent way through DPS. We use DPS as a principled way to exploit offline computational time in order to optimize the way how the agent will allocate its finite amount of online computational time.

We experimentally compare our hybrid OLT method with pure LT policies by using four different exploration strategies. We also consider two other state-of-the-art DPS techniques, where the policies are represented through neural networks and adaptive radial basis functions, respectively. On the basis of an extensive study on four challenging optimal control problems, we show that our approach indeed leverages the advantages of the two techniques from which it originates. In particular, we show that OLT policies generally achieve a better performance than either pure LT or pure DPS techniques (in one of the benchmark domains, significantly better), while requiring fewer online resources than LT and fewer offline resources than DPS.

The content of the paper is structured as follows. In Section 2, we begin by introducing basic notation and state formally the type of sequential decision-making problems this paper is about. Section 3 presents a detailed description of OLTs. The node scoring function of OLTs can be learned using any derivative-free global optimizer. One such optimizer that was shown to be highly relevant to DPS is GPO [4, 5]. We give a brief (but for all practical purposes fully sufficient) summary of this approach in Section 4. Section 5 presents the results of extensive experimental evaluations, wherein we compare OLTs with different pure DPS and pure LT techniques on a number of well-known benchmark domains. We present related work in Section 6, discuss weaknesses in Section 7, and finally conclude in Section 8.

2. PROBLEM STATEMENT

As a first step towards investigating the combination of DPS with LT policies, we focus in this paper on an elementary setting for optimal control: we consider fully observable deterministic problems with discrete actions. Dealing with continuous actions, stochasticity or partial observability is left for future work, and several directions to extend our ideas to these more complex settings are presented in Section 7.

Formally, we consider a discrete-time system whose dynamics are given by $x_{t+1} = f(x_t, a_t)$, where x is an element of some state space X and a is an element of some action space A . We do not make any assumptions about the form of the state space, as the presented approach can deal with any form of X . We do, however, make strong assumptions about the action space A : we require that it is finite, that is, $|A| = K$, and, for computational reasons, that the number of actions is small. In the examples we will consider later, X will be continuous and vector-valued, that is, $X \subset \mathbb{R}^{n_x}$ and each discrete action will be mapped to a specific continuous control vector (bang-bang control).

The system evolves at discrete-time steps with decision points occurring at $t = 1, 2, \dots$. For every transition we make, we observe a scalar *reward*[§] $q(x_t, a_t)$, which serves as a performance measure of the sum of which we want to optimize over time. We denote by $\underline{B} = \inf\{q(x_t, a_t) : (x, a) \in X \times A\}$ and $\overline{B} = \sup\{q(x_t, a_t) : (x, a) \in X \times A\}$, without however assuming that they are finite. As is common for this class of problems, myopic optimization, that is, choosing the one action a_t that maximizes q at each decision point t , does generally *not* produce optimality with respect to the accumulated sum of rewards. Let $\pi : X \rightarrow A$ denote a stationary *policy*, that is, a mapping from states to actions. For any given policy π and state x_0 , the *infinite horizon discounted sum of rewards*

[§]Note that because the dynamics are deterministic, the reward can implicitly also depend on the successor state $x_{t+1} = f(x_t, a_t)$.

$V^\pi(x_0)$ is defined as

$$V^\pi(x_0) := \lim_{T \rightarrow \infty} \sum_{t=0}^T \gamma^t \varrho(x_t, \pi(x_t)), \quad \text{where } \forall t : x_{t+1} = f(x_t, \pi(x_t)), \quad (1)$$

and where $0 < \gamma < 1$ is a discount factor. The optimal value function is defined as the maximum over all policies, $V^*(x_0) := \max_\pi V^\pi(x_0)$, and satisfies the discrete-time Hamilton–Jacobi–Bellman equation

$$V^*(x) = \max_a [\varrho(x, a) + \gamma V^*(f(x, a))] \quad \forall x \in X. \quad (2)$$

If we are able to determine the solution V^* of Equation (2), the optimal action for every state x can be derived as $\pi^*(x) = \operatorname{argmax}_a [\varrho(x, a) + \gamma V^*(f(x, a))]$. However, solving Equation (2) exactly is possible only in some cases (e.g., linear-quadratic regulator (LQR)); in the general nonlinear case and for higher-dimensional spaces, this presents an open research problem.

Direct policy search approaches do not try to estimate V^π or V^* . Instead, they solve the conceptually simpler problem of finding a policy that ‘works well’ for some initial conditions. In our case, we define this objective as follows: given a set of initial states $X_0 \subset X$, our goal is to find a policy that maximizes the performance over the states[¶] in X_0 , that is, we want to find

$$\operatorname{argmax}_\pi \sum_{x_0 \in X_0} V^\pi(x_0) \quad (3)$$

(or $\operatorname{argmax}_\pi \sum_{x_0 \in X_0} p(x_0) V^\pi(x_0)$, where $p(x_0) \geq 0$ are some weights).

In order to solve this problem, one assumes that policies are functions $\pi_\theta := \pi(\cdot; \theta)$ parameterized by some vector $\theta \in \mathbb{R}^d$. The optimization over policies in Equation (3) is thus turned into an optimization over real-valued vectors

$$\operatorname{argmax}_{\theta \in \mathbb{R}^d} V_{X_0}(\theta) := \sum_{x_0 \in X_0} V^{\pi_\theta}(x_0). \quad (4)$$

Given a vector θ , the objective function V_{X_0} can be evaluated by simulating the system under the policy π_θ from all the states in X_0 and summing the rewards as in Equation (1). Note that, to avoid the infinite sum, we have to truncate the infinite horizon to a (typically large) number of finite simulation steps H . The objective function in Equation (4) is thus effectively replaced by

$$\operatorname{argmax}_{\theta \in \mathbb{R}^d} V_{X_0}(\theta) := \sum_{x_0 \in X_0} \sum_{t=0}^H \gamma^t \varrho(x_t, \pi(x_t; \theta)), \quad \text{where } \forall t : x_{t+1} = f(x_t, \pi(x_t; \theta)). \quad (5)$$

The novelty of this paper is that we consider policies $\pi(\cdot; \theta)$, which, instead of being defined through function approximators, are non-trivial algorithms that at some point depend on θ and on the model (f, ϱ) . To emphasize this requirement, we will write $\pi_{f, \varrho}(\cdot; \theta)$ whenever we refer explicitly to an OLT policy.

3. OPTIMIZED LOOK-AHEAD TREES

Given a state x_t , we now describe a way of selecting action $\pi_{f, \varrho}(x_t; \theta)$ on the basis of the construction of an LT. The construction of the tree will, in general, be non-uniform and is controlled by a function that is parameterized by θ (non-uniform meaning that leaf nodes are not necessarily at the same depth).

[¶]Note that in the experiments we report on later, we will only use a single initial state. In the general case, having multiple initial states tend to increase the robustness of the policy.

3.1. Notation

Let K be the number of possible actions and \mathcal{T} denote the LT. \mathcal{T} is composed of nodes $n_{i,h} \in \mathcal{T}$, where h denotes the depth and i denotes the index within the depth h (i.e., $i \in \{1, \dots, K^h\}$). The state x_t for which we want to find action $\pi_{f,\varrho}(x_t; \theta)$ is placed in the root node $n_{1,0}$. Each node $n_{i,h}$ corresponds to a particular sequence of actions and states. The children of $n_{i,h}$ are generated by applying one of the K actions: assume $a^{(h)} \in \{1, \dots, K\}$ is the index of the action taken at depth h , then the child of $n_{i,h}$ is $n_{K(i-1)+a^{(h)},h+1}$.

For each node $n_{i,h}$, there exists one unique path from the root $n_{1,0}$ to $n_{i,h}$:

$$\text{path from } n_{1,0} \text{ to } n_{i,h} : n_{1,0} \longrightarrow n_{i_1,1} \longrightarrow n_{i_2,2} \longrightarrow \dots \longrightarrow n_{i_h,h} = n_{i,h},$$

which corresponds to a particular sequence of actions $a^{(1)}, a^{(2)}, \dots, a^{(h)}$ taken at each intermediate step and which induces a partial trajectory of length h :

$$x_t \xrightarrow[\varrho(x_t, a^{(1)})]{a^{(1)}} x_{t+1}^{a^{(1)}} \xrightarrow[\varrho(x_{t+1}^{a^{(1)}}, a^{(2)})]{a^{(2)}} x_{t+2}^{a^{(1)}a^{(2)}} \dots x_{t+h-1}^{a^{(1)}\dots a^{(h-1)}} \xrightarrow[\varrho(x_{t+h-1}^{a^{(1)}\dots a^{(h-1)}}, a^{(h)})]{a^{(h)}} x_{t+h}^{a^{(1)}\dots a^{(h)}}.$$

The successor states are generated according to the transition function f , for example, $x_{t+i}^{a^{(1)}\dots a^{(i)}} = f(x_{t+i-1}^{a^{(1)}\dots a^{(i-1)}}, a^{(i)})$ and rewards according to the reward function ϱ . For each node $n_{i,h}$, we define $\varrho(n_{i,h})$ to be the reward obtained in the last step: $\varrho(n_{i,h}) := \varrho(x_{t+h-1}^{a^{(1)}\dots a^{(h-1)}}, a^{(h)})$.

Every time we expand a node, we generate all of its children. A node whose children have been generated is called an *inner node*. Otherwise, it is called a *terminal node*.[¶] We denote $\mathcal{T} = \mathcal{T}_{\text{inner}} \cup \mathcal{T}_{\text{leaf}}$.

3.2. Using look-ahead trees to make decisions

We first describe how to use the information in an LT to select an action. For each terminal node $n_{i,h} \in \mathcal{T}_{\text{leaf}}$, we define the ℓ -score as the discounted sum of rewards obtained along the path from the root node $n_{1,0}$ to $n_{i,h}$ plus a lower bound on the cumulated rewards not yet observed:

$$\forall n_{i,h} \in \mathcal{T}_{\text{leaf}} : \ell(n_{i,h}) := \sum_{t=1}^h \varrho(n_{i,t}) \gamma^{t-1} + \sum_{t=h+1}^{\infty} \underline{B} \gamma^{t-1} = \sum_{t=1}^h \varrho(n_{i,t}) \gamma^{t-1} + \frac{\underline{B} \gamma^h}{1-\gamma}. \quad (6)$$

For each non-terminal node $n_{i,h} \in \mathcal{T}_{\text{inner}}$, we define the ℓ -score recursively as the maximum of the ℓ -scores of its children (Figure 1):

$$\forall n_{i,h} \in \mathcal{T}_{\text{inner}} : \ell(n_{i,h}) := \max_{n \in \text{children}(n_{i,h})} \ell(n). \quad (7)$$

The ℓ -score of the root node $n_{1,0}$ (which corresponds to state x_t) is a lower bound on the optimal value $V^*(x_t)$ [2]. Given an LT, we adopt the conservative strategy that consists in selecting the action that leads to the successor state with maximal lower bound. With the naming scheme for nodes introduced earlier, we can write $\ell(n_{1,0}) = \max_a \ell(n_{a,1})$, and thus

$$\pi_{f,\varrho}(x_t) = \underset{a}{\operatorname{argmax}} \ell(n_{a,1}). \quad (8)$$

In the rest of this paper, we generally assume that $\underline{B} = 0$ and that \overline{B} is finite, although neither of these conditions is a requirement for our OLT method.

[¶]Note that the following terminology is equivalent:

non-terminal node = inner node = explored node = closed node
terminal node = leaf node = unexplored node = open node.

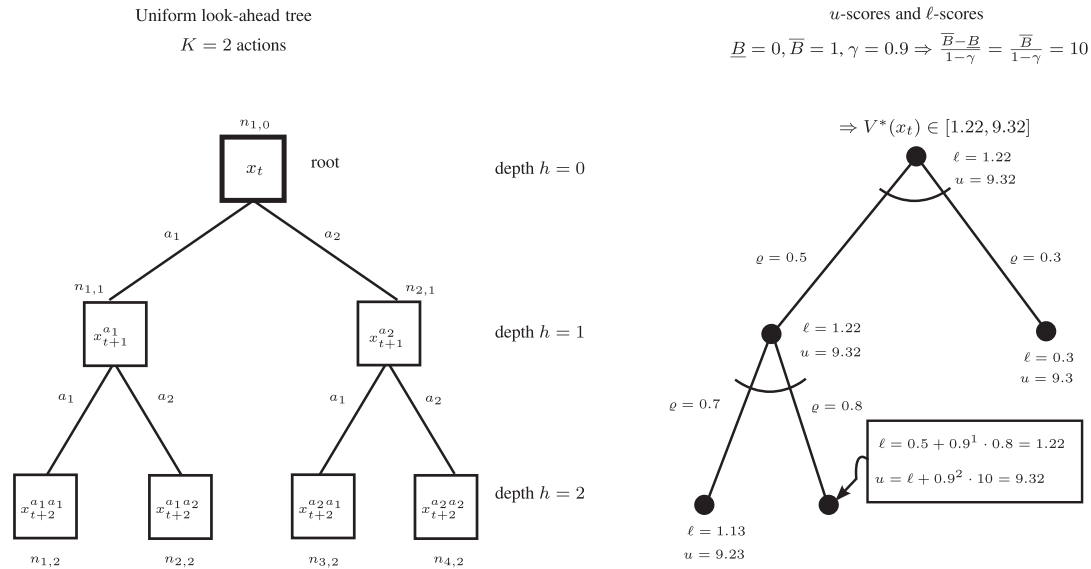


Figure 1. Left side: A uniform look-ahead tree developed from state x_t for $h = 2$ stages and $K = 2$ actions. In it, nodes are labeled by $n_{i,h}$ and correspond to states; edges correspond to choosing an action. Each node encodes the result of taking a particular sequence of actions; for example, the node $n_{2,2}$ corresponds to state $x_{t+2}^{a_1 a_2}$, which is the result of starting from state x_t and choosing action a_1 in the first stage and action a_2 in the second stage. Right side: Illustrates how ℓ -scores and u -scores are calculated (which form an upper and lower bound for the true value $V^*(x_t)$ of the root node). The plot shows a non-uniform tree where the terminal node with the highest u -score is expanded first. Edges are labeled with the rewards ρ associated with the transition they represent. Note that the ℓ -scores and u -scores are calculated only for terminal nodes and propagated back to their parent via the \max operator (denoted by the arc).

3.3. Developing look-ahead trees

We now discuss the construction of LTs. If the rewards are upper bounded (\bar{B} is finite), we can also upper bound the optimal values $V^*(x)$. For each terminal node $n_{i,h} \in \mathcal{T}_{\text{leaf}}$, we define the u -score as the discounted sum of rewards obtained along the path from the root node $n_{1,0}$ to $n_{i,h}$, plus an upper bound on the cumulated rewards not yet observed:

$$\forall n_{i,h} \in \mathcal{T}_{\text{leaf}} : u(n_{i,h}) := \sum_{t=1}^h \rho(n_{i,t}) \gamma^{t-1} + \sum_{t=h+1}^{\infty} \bar{B} \gamma^{t-1} = \ell(n_{i,h}) + \frac{(\bar{B} - \underline{B}) \gamma^h}{1 - \gamma}. \tag{9}$$

For each non-terminal node $n_{i,h} \in \mathcal{T}_{\text{inner}}$, we define the u -score recursively as the maximum of the u -scores of its children (Figure 1):

$$\forall n_{i,h} \in \mathcal{T}_{\text{inner}} : u(n_{i,h}) := \max_{n \in \text{children}(n_{i,h})} u(n). \tag{10}$$

See Figure 1 for a graphical illustration of an LT with associated ℓ -scores and u -scores. The u -score of the root node $n_{1,0}$ is an upper bound on the optimal value $V^*(x_t)$. Furthermore, the ℓ -score of the root increases with the number of expanded nodes in the tree, and at the same time its u -score decreases [2]. In other words, the more nodes the tree has and the farther we develop it into the future, the tighter our bounds on the optimal value will become** and thus the better a decision we can hope to make with Equation (8). (Bear in mind that, unlike search trees for constraint satisfaction problems, which are grown up to a goal state, our LTs are only grown until a predefined computational budget is exhausted.)

**The tightness of the bounds will also depend on γ . The closer γ is to 1, the more nodes we will have to expand. However, there is little we can do about it, as γ is given in advance.

Knowing that, the big question is *How should we develop the tree such that we arrive at a near-optimal decision within the allowed computational budget?*

A typical way of developing an LT is to build a uniform tree, by expanding nodes with a breadth-first strategy. However, to fully develop a tree of depth h , we need $K^0 + K^1 + K^2 + \dots + K^h$ node expansions, which for growing h will rapidly become computationally infeasible.

Best-first search is the common solution to this dilemma; it develops trees non-uniformly and only expands those nodes to a deeper depth that look ‘promising’. In order to do this, the algorithm relies on a *node scoring function* $e : \mathcal{T}_{\text{leaf}} \rightarrow \mathbb{R}$ that is used to assign a certain score to every terminal node $n_{i,h} \in \mathcal{T}_{\text{leaf}}$; whenever we want to decide which terminal node to expand next, we compare all the scores and choose the node with the highest score. Note that the default uniform development strategy can be obtained as a particular case of best-first search with the following node scoring function:

$$e^{\text{uniform}}(n_{i,h}) := -h. \quad (11)$$

The work in [2] suggests to use as scoring function the u -scores of the terminal nodes (which makes the method become closely related to A^*):

$$e^{\text{optimistic}}(n_{i,h}) := u(n_{i,h}). \quad (12)$$

The authors of [2] show theoretically that non-uniform trees developed by the u -score will never perform worse than uniform trees for the same budget of expansions, but certain conditions must be met to make them perform better. The extent to which these conditions are fulfilled is problem specific and depends on how ‘informative’ the reward is (which in turn depends on the coupling of dynamics and rewards). Informally speaking, a problem will be difficult (i.e., the rewards will be non-informative) if there are many nodes with the property that, given the observed rewards from the root to the node in question, one cannot decide whether the node lies on an optimal path or not. This is, for example, the case for optimal control problems with a flat reward structure, such as in the mountain car domain [6] (where every transition has -1 reward and only entering the goal gives 0). On the other hand, having an informative reward is not a totally exotic requirement; for many control problems, an informative reward comes as the natural definition of performance (e.g., the many pole balancing or inverted pendulum domains, where rewards are taken as a quadratic function of angle and angular velocity).

3.4. Parameterizing the development of the look-ahead tree

It turns out that a large number of strategies could be used to develop LTs, and it is probably the case that there exists no single best strategy for all problems. Instead of searching for the best possible generic strategy, we adopt the approach first introduced in [3], which consists in learning a specific look-ahead node development strategy in a problem-driven way.

Optimized LTs rely on a parameterized node scoring function $e(\cdot; \theta) : \mathcal{T}_{\text{leaf}} \rightarrow \mathbb{R}$, where $\theta \in \mathbb{R}^d$ is a parameter vector. This function should be flexible enough to represent a large variety of tree-development strategies. In particular, the parameters should encode important aspects of search, such as the extent to which depth-first search is preferred over breadth-first search, or how short-term rewards should be used to bias search. Furthermore, as different search strategies may be optimal in different regions of the state space, we would like $e(\cdot; \theta)$ to enable state-dependent strategies. To meet these requirements, we take, as in [3], the parameterized node scoring function to be a simple weighted sum of features extracted from the information encoded in the path from the root node to the node in question.

In our investigations, the node scoring function is defined for each terminal node $n_{i,h} \in \mathcal{T}_{\text{leaf}}$ in the following way: let $x_{t+h}^{a^{(1)} \dots a^{(h)}} = (x_{t+h}^{(1)}, \dots, x_{t+h}^{(n_x)}) \in \mathbb{R}^{n_x}$ denote the n_x -dimensional state corresponding to $n_{i,h}$, and let $q_{t+h} = q(n_{i,h})$ denote its reward. We consider three blocks of features: the first n_x features correspond to the components of the state, the next n_x features correspond to components of the state multiplied by the reward (enabling the tree to be grown in a more or less

directed way), and the final set of n_x features corresponds to the components of the state multiplied by the depth h (enabling to control breadth/depth trade-off). Let $\theta \in \mathbb{R}^{3n_x}$ be the vector of parameters. The parameterized scoring function $e(\cdot; \theta)$ can then be written as

$$\forall n_{i,h} \in \mathcal{T}_{\text{leaf}} : e(n_{i,h}; \theta) := \sum_{j=1}^{n_x} x_{t+h}^{(j)} (\theta_j + \theta_{n_x+j} \cdot \varrho_{t+h} + \theta_{2n_x+j} \cdot h). \quad (13)$$

Notice that with this linear parameterization, the outcome of which node is selected is invariant under scaling of the parameter vector by positive scalars. Of course, other kinds of features are also possible and may in fact turn out to be more suitable in some cases (this may be an avenue for future research).

3.5. Summary: the algorithm

Figure 2 presents a simple algorithm based on a sorted list to implement policies as parameterized LTs. The algorithm requires as input a state x_t and the parameter vector θ and returns the action $\pi_{f,\varrho}(x_t; \theta)$. It depends on the allowed number of node expansions h_{\max} , the domain (represented by generative model f, ϱ), and the discount factor γ . The computational complexity for evaluating $\pi_{f,\varrho}(x_t; \theta)$ for a budget of h_{\max} node expansions is obtained as follows. Let D be the cost for one call to the generative model and E be the cost to recursively update the internal scores of a node. We assume that the leaf nodes are stored in an appropriate data structure, which allows incremental insertion in logarithmic time and allows finding the maximum of the scores in constant time (such that for every iteration of the main loop, if the list of terminal nodes contains N elements, we can find the maximum score and update the structure with $\log N$ operations). Note that for every iteration $j = 1, 2, \dots$ of the main loop, our list of terminal nodes contains $(K-1)(j-1) + K$ elements,

Input: state x_t , weights θ	
Output: policy action $\pi_{f,\varrho}(x_t; \theta)$	Node n is a struct consisting of fields:
Depends on:	$n.x$ state
f transition function	$n.h$ depth
ϱ reward function	$n.\ell$ ℓ -score=cumulative reward
γ discount factor	$n.e$ tree development score from $e(\cdot; \theta)$
$e(\cdot; \theta)$ parameterized score function	$n.\pi$ first action on path
h_{\max} maximum number of node expansions	

1. Initialize
list $\leftarrow\emptyset$
For $i = 1 \dots K$ /* Expand root node for all actions */
 $x' := f(x_t, a_i)$; $r' := \varrho(x_t, a_i)$; $h' := 1$; $e := e(x', r', h'; \theta)$
Generate node n :
 $n.x := x'$; $n.\ell := r'$; $n.h := h'$; $n.e := e$; $n.\pi := a_i$
Add n to list

2. Main loop
While number of node expansions $j < h_{\max}$ (the max number of node expansions)
Find $n^* := \operatorname{argmax}_{n \in \text{list}} n.e$ and remove n^* from list
For $i = 1 \dots K$ /* Expand node n^* for all actions */
 $x' := f(n^*.x, a_i)$; $r' := \varrho(n^*.x, a_i)$; $h' := n^*.h + 1$; $e := e(x', r', h'; \theta)$
Generate node n :
 $n.x := x'$; $n.\ell := n^*.\ell + \gamma^{h'-1} r'$; $n.h := h'$; $n.e := e$; $n.\pi := n^*.\pi$
Add n to list

3. Get best action
Return $\pi_{f,\varrho}(x_t; \theta) := n^*.\pi$ where $n^* = \operatorname{argmax}_{n \in \text{list}} n.\ell$ /* ties broken randomly */

Figure 2. Implementing a policy represented by a parameterized look-ahead tree.

which can be shown by simple induction (every time we expand a node, we remove one element from the list and add K new ones). At each iteration j of the main loop, we do the following: we first find the best node to expand by looping over the list of current terminal nodes. We then generate all of its K successors, where each one costs D for having to call the generative model and E to recursively update its internal scores. After the main loop has run for h_{\max} times, we have to find the best ℓ -score from the list. The computational complexity of the algorithm is thus

$$\sum_{j=1}^{h_{\max}} [K(D + E) + \log\{(K - 1)(j - 1) + K\}] + (K - 1)(h_{\max} - 1) + K \quad (14)$$

$$\leq h_{\max}(K(D + E) + \log\{(K - 1)(h_{\max} - 1) + K\}) + K - 1 + 1.$$

In the particular case where we are expanding by the heuristic *uniform* (i.e., breadth-first), we skip the maximization within each iteration of the main loop. The computational complexity then becomes

$$\sum_{j=1}^{h_{\max}} [K(D + E)] + (K - 1)(h_{\max} - 1) + K = h_{\max}(K(D + E) + K - 1) + 1. \quad (15)$$

Note that in order to expand uniformly, h_{\max} must be equal to one of $K^0, K^0 + K^1, K^0 + K^1 + \dots + K^d$, where d is the depth of the tree.

4. GAUSSIAN PROCESS OPTIMIZATION

We now turn to the problem of solving Equation (5), that is, find a vector θ such that the induced policy $\pi(\cdot; \theta)$ (globally) maximizes the score $V_{X_0}(\theta)$ over the set of initial states X_0 . In general, this will be a difficult optimization problem. First, the dependence of V_{X_0} on θ can be a complex one with local extrema occurring frequently; thus we may need to sample the search space exhaustively. Second, there is no closed-form expression for evaluating the objective function V_{X_0} or its gradient; instead we have to simulate the system (or run real-world experiments), which is expensive. Among the many alternatives that have been proposed in the past for this purpose, such as cross-entropy (CE) [7], various stochastic search alternatives [8], or Lipschitzian optimization [9], GPO is considered to be one of the most efficient methods to optimize expensive functions [10].

The purpose of this section is to provide the required background in GPO to the reader who is not aware of this. Note that in our experiments, we will illustrate the optimization of LTs with both GPO and the CE method, the latter being much simpler to understand and to implement than the former. As choosing a policy representation and optimizing its parameters are two separate tasks, we keep the description of GPO general and independent of the specific nature of the OLT policies. Readers less interested in the optimization part may safely skip this section and directly jump to Section 5.

4.1. Notation

To enhance the readability of this section and to conform with the standard notation used in the literature, we will define a local notation for this section, which overlaps with what we use in the remainder of the article. Specifically, we will now write $f(x)$ for the objective function we want to maximize (instead of $V_{X_0}(\theta)$) and x for its input arguments (instead of θ).

4.2. Overview

Gaussian process optimization is an iterative sampling-based search procedure, which constructs a surrogate for the objective function and optimizes that in place of the original one. GPO is able to incorporate prior knowledge about the problem and provides a principled way to build and exploit the surrogate so as to trade-off exploration and exploitation of the search space.

Each iteration of GPO consists of two steps. Assume that at iteration n , we have already sampled the objective function at locations x_1, \dots, x_n with values $f(x_1), \dots, f(x_n)$. The first step is

to fit a regression model by using Gaussian process regression to the samples gathered so far, that is, to fit a regression model to the training data $\mathcal{D}_n := \{(x_i, f(x_i))\}_{i=1}^n$. Let us call the resulting model GP_n . The second step is to use the regression model GP_n as input to a scoring heuristic \mathcal{U} (the so-called acquisition function) to find the most promising point x_{n+1} at which to evaluate the objective function next. Typically, acquisition functions are defined in an optimistic way such that high scores correspond to *potentially* high values of the objective function.

The regression model GP_n acts as a surrogate of the objective function: it can be evaluated at any given point x of the search space to produce an estimate for $f(x)$ (more precisely, GP_n will produce a distribution over $f(x)$). The reason for building the surrogate is that, unlike the true objective function, it is computationally very cheap to evaluate (it can be done analytically and in closed form). Thus we can afford sampling it as exhaustively as necessary to find the maximum. On the other hand, as the values the surrogate produces will only be estimates, we can never be sure that what we get from maximizing the surrogate is indeed the maximum of the objective function or even close to it. Instead, we have to take into account how accurate the estimates produced by GP_n are; this in turn will depend on the general smoothness of the objective function and the number of data points we have collected in the neighborhood of x .

Each time GP_n is evaluated at a location x , it outputs two values $\mu_n(x)$ and $\sigma_n^2(x)$, which together define the Gaussian predictive distribution $N(\mu_n(x), \sigma_n^2(x))$ over values $f(x)$. The mean of this distribution, $\mu_n(x)$, can be directly taken as point estimate for $f(x)$. The variance of this distribution, $\sigma_n^2(x)$, can be taken as a measure of how certain the GP is about this estimate. Both $\mu_n(x)$ and $\sigma_n^2(x)$ will be used by the acquisition function \mathcal{U} to assign a score to x , which we will write as $(\mathcal{U} GP_n)(x)$. To determine the next sample location x_{n+1} , we thus have to determine

$$x_{n+1} := \operatorname{argmax}_x (\mathcal{U} GP_n)(x), \quad (16)$$

which, unlike Equation (5), can be solved efficiently by any black-box global optimization method (in our experimental studies, we will use DIRECT [9]).

In the following two sections, we will describe each of these steps in more detail; a summary of the algorithm is also given as pseudo-code in Figure 3.

Input: objective function $V_{X_0}(\cdot)$ taking arguments $\theta \in \mathbb{R}^d$ (from Eq. (4))
Output: $\theta^* \approx \operatorname{argmax}_{\theta} V_{X_0}(\theta)$

Depends on:
 Acquisition function \mathcal{U} that takes as functional input a GP model and as argument a vector $\theta \in \mathbb{R}^d$

1. Initialize
 Generate initial sample locations $\theta_1, \dots, \theta_n$ by random sampling or space filling methods (e.g., Latin hypercube sampling)
 Evaluate the objective function in $\theta_1, \dots, \theta_n$ to produce training data $\mathcal{D}_n := \{(\theta_i, V_{X_0}(\theta_i))\}_{i=1}^n$

2. Main loop $i = n, n + 1, \dots$ /* loop until we run out of computational resources */
 Fit GP model to the current training data:
 $GP_i =$ Gaussian process regression on training data \mathcal{D}_i
 Get most promising next sample location:
 $\theta_{i+1} := \operatorname{argmax}_{\theta} (\mathcal{U} GP_i)(\theta)$ /* using e.g. DIRECT */
 Evaluate objective function V_{X_0} in θ_{i+1} , add result to training data, and repeat:
 $\mathcal{D}_{i+1} = \mathcal{D}_i \cup \{(\theta_{i+1}, V_{X_0}(\theta_{i+1}))\}$

3. Return
 Training point in \mathcal{D}_n with highest score:
 $\operatorname{argmax}_{\theta \in \mathcal{D}_n} V_{X_0}(\theta)$

Figure 3. Implementing Gaussian process optimization.

4.3. Using Gaussian processes for regression

Here we briefly review how GPs can be used for function estimation [11]. Note that we will present GPs for the more general case of a stochastic objective function; although the problems we consider later are all deterministic, this leaves the door open for stochastic returns in future work.

Suppose we are looking for a function $f : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ from which we have observed noisy samples $(x_1, y_1), \dots, (x_n, y_n)$, where $x_i \in \mathcal{X}$ is the input and $y_i = f(x_i) + \varepsilon_i$ the output corrupted by independent zero-mean Gaussian noise with common variance σ_0^2 , that is, $\varepsilon_i \sim_{\text{iid}} \mathcal{N}(0, \sigma_0^2)$. To estimate the function value $f(x)$ at any given input location x , we proceed as follows. We suppose that the sought function is a realization of a zero-mean^{††} Gaussian process with covariance function $k_\vartheta(x, x')$, which we write as $f \sim \mathcal{GP}(0, k_\vartheta(x, x'))$, where ϑ is a vector of hyperparameters (as explained below). Hence, the vector of function values at the n observed input locations is assumed to be drawn from a joint Gaussian distribution

$$(f(x_1), \dots, f(x_n)) | X, \vartheta \sim \mathcal{N}(0_{n \times 1}, K_\vartheta), \quad (17)$$

where $X := [x_1, \dots, x_n]$, and K_ϑ is the $n \times n$ covariance matrix with entries $[K_\vartheta]_{i,j} = k_\vartheta(x_i, x_j)$. The covariance function $k_\vartheta(\cdot, \cdot)$ can be thought of as a way to encode our prior information about the ‘smoothness’ of the functions f we believe to come up; typically $k_\vartheta(x, x')$ is chosen as a function of the distance $\|x - x'\|$ in which case k_ϑ measures the expected amount of variation of the function values $f(x), f(x')$ in terms of the distance between the locations x and x' .

A typical choice for k is the translation-invariant squared exponential, which is of the form

$$k_\vartheta(x, x') := v_0 \exp\{-0.5(x - x')^T \Omega (x - x')\} \quad (18)$$

and which itself is parameterized by hyperparameters $v_0 > 0$ and $\Omega = \text{diag}(a_1, \dots, a_d)$, $a_i > 0$. All hyperparameters specifying the GP are thus collected in the vector $\vartheta := (v_0, a_1, \dots, a_d)$, which together with σ_0^2 characterizes the prior distribution of our noisy samples. The actual ‘training’ of a GP thus consists of finding a good pair (ϑ, σ_0^2) from the data, for which various alternative procedures exist in the literature: here we stick to the most common one, which is the optimization of the marginal likelihood. See [11] for a detailed description.

Now suppose that we know ϑ and σ_0^2 . As the noise is zero-mean Gaussian, white, and independent of the function f , it follows from Equation (17) that the $n \times 1$ vector of observed outputs $Y := [y_1, \dots, y_n]^T$ will also be jointly Gaussian and distributed as follows:

$$Y | X, \vartheta, \sigma_0^2 \sim \mathcal{N}(0_{n \times 1}, K_\vartheta + \sigma_0^2 I_{n \times n}). \quad (19)$$

Furthermore, the joint distribution of the function value $f(x)$ at query location x and the vector of observed values of Y is also Gaussian and characterized in the following way:

$$\begin{bmatrix} Y \\ f(x) \end{bmatrix} | X, \sigma_0^2, \vartheta, x \sim \mathcal{N}\left(\begin{bmatrix} 0_{n \times 1} \\ 0 \end{bmatrix}, \begin{bmatrix} K_\vartheta + \sigma_0^2 I_{n \times n} & k_\vartheta(x) \\ k_\vartheta(x)^T & \kappa_\vartheta \end{bmatrix}\right), \quad (20)$$

where the $n \times 1$ vector $k_\vartheta(x)$ is defined by $k_\vartheta(x) := [k_\vartheta(x, x_1), \dots, k_\vartheta(x, x_n)]^T$ and scalar κ_ϑ by $\kappa_\vartheta := k_\vartheta(x, x)$. Conditioning $f(x)$ on Y , we thus obtain

$$f(x) | X, \vartheta, \sigma_0^2, x, Y \sim \mathcal{N}(\mu(x), \sigma^2(x)), \quad (21)$$

where

$$\mu(x) := k_\vartheta(x)^T (K_\vartheta + \sigma_0^2 I_{n \times n})^{-1} Y \quad (22)$$

$$\sigma^2(x) := \kappa_\vartheta - k_\vartheta(x)^T (K_\vartheta + \sigma_0^2 I_{n \times n})^{-1} k_\vartheta(x). \quad (23)$$

Thus given (noisy) observations $(x_1, y_1), \dots, (x_n, y_n)$ from an unknown function f , with GP regression, we first infer from the observations hyperparameters ϑ and σ_0^2 , and then obtain for any new point x the distribution over function values $p(f(x)|X, \vartheta, x, y) = \mathcal{N}(\mu(x), \sigma^2(x))$.

^{††}The assumption of zero mean is made here for notational convenience. (Centering the data combined with an ergodicity assumption allows to avoid a non-zero mean).

4.4. Choosing an acquisition function

Early work [12, 13] suggested to take as acquisition function the probability of improving over the current maximum $x^+ := \operatorname{argmax}_{x \in \mathcal{D}_n} f(x)$ within the training data \mathcal{D}_n . The resulting probability of improvement (PI) acquisition function is given by $P(f(x) \geq f(x^+) + \zeta)$, which we write as

$$(PI\ GP_n)(x) := \Phi_{0,1} \left(\frac{\mu_n(x) - f(x^+) - \zeta}{\sigma_n(x)} \right), \quad (24)$$

where $\mu_n(x)$ is the mean and $\sigma_n^2(x)$ the variance of the predictive distribution as output by GP_n for point x (Equations (22) and (23)), and $\Phi_{0,1}$ is the standard normal cumulative distribution. The trade-off parameter $\zeta \geq 0$ controls the compromise between the strength of the potential improvement and its probability to be realized.

An alternative acquisition function is the expected improvement (EI), which can be evaluated analytically [12], giving

$$(EI\ GP_n)(x) := \sigma_n(x)(Z\Phi_{0,1}(Z) + \phi_{0,1}(Z)), \quad (25)$$

where $Z := (\mu_n(x) - f(x^+) - \zeta)/\sigma_n(x)$, and $\phi_{0,1}$ is the probability density of the standard normal distribution.

Figure 4 illustrates how GPO works and how PI and EI give rise to distinct sampling behavior by using the same GP and data. Which of the acquisition functions works best for a given problem is usually difficult to say in general and needs experimentation; in our examples in Section 5, we have found that EI worked best.

5. EXPERIMENTS

This section reports an extensive evaluation of OLTs. The purpose of the experiments is to compare OLT to both pure DPS and pure LT approaches and to demonstrate its strengths across a variety of domains. To this end, we have chosen four challenging benchmark problems, which are described in Section 5.1: *inverted pendulum*, *double inverted pendulum*, *acrobat handstand*, and *HIV drug treatment*. We compare our approach to two of today's standard methods for DPS and to four generic LT policies, all of which are described in Section 5.2. The comparison is structured as follows:

- In Section 5.3, we study the performance of OLT policies in terms of accumulated reward achieved after learning and compare OLT with both state-of-the-art DPS and LT approaches. We show that OLT policies significantly outperform the other approaches in some cases and perform about the same in the remaining ones. We also show that OLT requires significantly fewer online resources than comparable LT policies.
- In Section 5.4, we study the offline complexity of OLT learning and compare its sample efficiency (both in terms of policy evaluations and number of simulated transitions) with that of DPS. We show that most of the time OLT requires fewer samples than its DPS competitors to reach high-performance policies.
- In Section 5.5, we compare how the different policy parameterizations behave with respect to their respective hyperparameters and show that the OLT representation requires much less trial-and-error effort than the traditional parametrizations studied in DPS.
- In Section 5.6, we study the robustness of LT, standard DPS, and OLT when the initial states of the testing set differ from the initial states of the training set and show that OLT policies are quite robust with respect to both small and large perturbations of the initial state.
- Finally, in Section 5.7, we examine the robustness of OLT when the budget of allowed node expansions used during evaluation of a policy differs from the budget that was used to learn it.

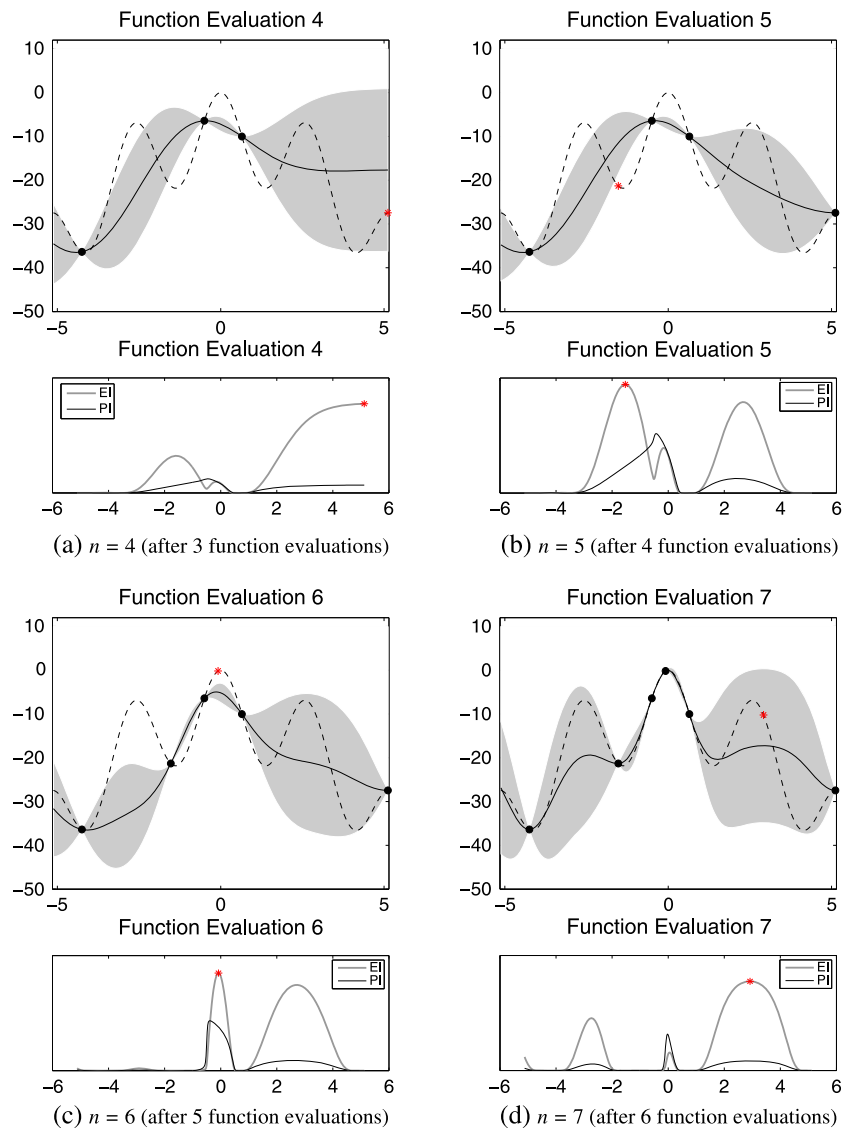


Figure 4. A concrete example illustrating Gaussian process optimization over four stages $n = 4, 5, 6, 7$. Each panel consists of two subplots. Each upper subplot shows the true (unknown) function we want to optimize as a dashed curve and the locations at which it was previously evaluated as black filled dots. Using these as training data \mathcal{D}_n in a GP, the black curve denotes the expected value (mean) of the resulting predictive distribution from Equation (22) and the shaded area denotes its variance from Equation (23), both evaluated at all locations $x \in [-5, 5]$. As one would expect, the predictive variance (and thus the uncertainty about the estimates) is close to zero at 'known' locations (it would be equal to zero if $\sigma_0^2 = 0$) and grows the farther one gets away from them. Each lower subplot shows the result of evaluating the acquisition function for one of the two possible choices EI or PI. Note that each of these would lead to a (slightly) different sampling behavior; here we have chosen the next sample location (marked in the figure by the star-shaped symbol) as the one that maximizes EI.

5.1. The benchmark domains

We consider four well-known challenging benchmark domains that exhibit some common characteristics: a continuous vector-valued state space, finite (discretized) actions, deterministic transitions, and rewards that to some extent are informative about the goal. Each domain corresponds to a real physical process (either mechanical or bio-chemical) internally described by a system of nonlinear differential equations, none of which can be solved by traditional control methods such as LQR.

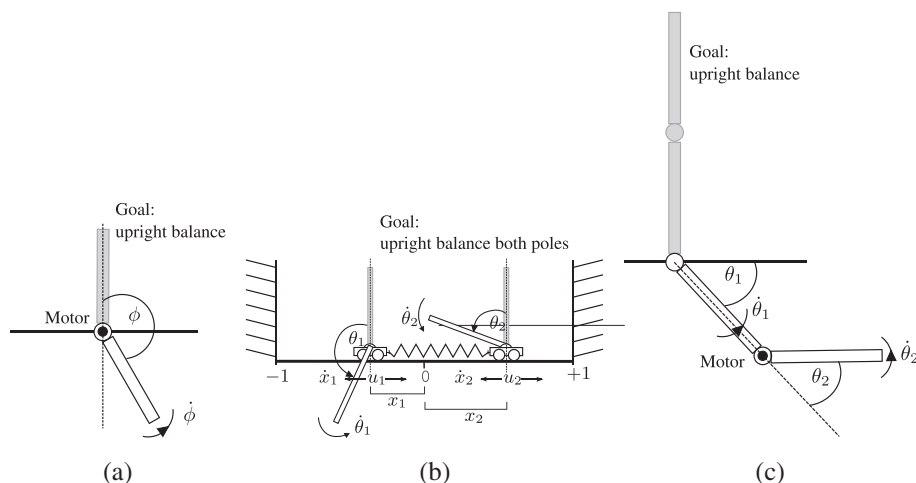


Figure 5. From left to right: the inverted pendulum task, the double inverted pendulum linked with a spring task, and the acrobot handstand task.

The transition function is obtained by discretizing in time and keeping the controls constant; the actions are obtained by discretizing the bounded control space.

Inverted pendulum. Our first domain is the inverted pendulum, a simple enough toy problem that is widely used in benchmarking different algorithms. The goal is to swing up and stabilize a single-link inverted pendulum as is shown in Figure 5(a). As the motor does not provide enough torque to push the pendulum up in one single rotation, the pendulum needs first to be swung back and forth to gather energy before then being pushed up and balanced. This creates a nonlinear control problem. The state space is two dimensional, $x = (x_1, x_2) \equiv (\phi, \dot{\phi})$, with $\phi \in [-\pi, \pi]$ being the angle and $\dot{\phi} \in [-10, 10]$ being the angular velocity. The control force is discretized to $a \in \{-5, -2.5, 0, +2.5, +5\}$ and held constant for $\Delta t = 0.2$ s. The dynamics of the system and physical parameters we used to instantiate the problem are specified in Appendix A.

To formulate this task as an optimal control problem of the form given in Equation (5), we used the following settings: the set of initial states is the singleton $X_0 = \{(\pi, 0)\}$, the reward is defined as $q(x(t), u(t)) := 1 - 0.1(\phi(t)^2 - 0.1\dot{\phi}(t)^2 - 0.1u(t)^2)$, the discount factor is set to $\gamma = 0.99$, and each policy is evaluated for $H = 500$ steps.

Double inverted pendulum. The next domain is a more complex variant of the inverted pendulum given above. This time we have two poles mounted each to a separate cart as depicted in Figure 5(b). The two carts are linked by a spring and allowed to move some distance horizontally on the x -axis (until they collide with a wall, which leads to a failure). Each cart is controlled separately; however, because of the spring, their dynamics is coupled. As in the inverted pendulum, the goal is to swing up and stabilize the poles as quickly as possible, but now by moving the carts back and forth. This creates a rather challenging nonlinear control problem. The state space is eight dimensional, $x \equiv (x_1, x_2, \dot{x}_1, \dot{x}_2, \theta_1, \theta_2, \dot{\theta}_1, \dot{\theta}_2)$, with $x_i \in [-1, 1]$ being the position and $\dot{x}_i \in [-10, 10]$ the velocity of the i th cart, and with $\theta_i \in [0, 2\pi]$ being the angle and $\dot{\theta}_i \in [-5, 5]$ being the angular velocity of the i th pole. The two-dimensional control vector is discretized to the four actions $a \in \{(-2, -2), (-2, +2), (+2, -2), (+2, +2)\}$ and held constant for $\Delta t = 0.1$ s. The dynamics of the system and physical parameters we used to instantiate the problem are specified in Appendix B. To formulate this task as an optimal control problem of the form given in Equation (5), we used the following settings: the set of initial states is the singleton $X_0 = \{(0, 0.5, \pi, \pi, 0, 0, 0, 0)\}$, the reward is defined as $q(x(t), u(t)) := [(1 + \cos \theta_1(t)) + (1 + \cos \theta_2(t))]/4$, the discount factor is set to $\gamma = 0.999$, and each policy is evaluated for $H = 250$ steps. In order to achieve a high reward in this

domain, the policy has to both balance the poles and not collide with one of the walls. Low rewards usually occur because of such collisions.

Acrobot. Our third domain is the acrobot from [14]: a two-link robot that resembles a gymnast swinging up above a high bar (Figure 5(c)). The acrobot freely swings around the first joint (the hands grasping the bar) and can exert force only at the second joint (bending the hips). The acrobot is an underactuated system; the task we consider here is the inverted ‘handstand’ position, which is hard to solve using standard RL methods. The state space is four dimensional, $x \equiv (\theta_1, \theta_2, \dot{\theta}_1, \dot{\theta}_2)$, with θ_1 and θ_2 being the angle of the upper and lower links, and $\dot{\theta}_1$ and $\dot{\theta}_2$ their angular velocity, respectively. The continuous control is discretized to $\{-1, +1\}$ and held constant for $\Delta t = 0.2$ s. To facilitate staying stable in the inverted handstand position (a highly unstable equilibrium), we also include a third non-primitive ‘balance’ action, which chooses control values derived from an LQR controller obtained from linearizing the system dynamics about the handstand position. Note that this balance action produces meaningful outputs $\in [-1, +1]$ only very close to the unstable equilibrium and thus cannot be used to bring the acrobot from the initial state to the goal region. The dynamics of the system and physical parameters we used to instantiate the problem are specified in Appendix C.

To formulate this task as an optimal control problem of the form given in Equation (5), we used the following settings: the set of initial states is the singleton $X_0 = \{(-\pi/2, 0, 0, 0)\}$, the reward is defined by the height of the end of the second link (the feet) as $\varrho(x(t), u(t)) = 2 + \cos(\theta_1(t) - \pi/2) + \cos(\theta_1(t) + \theta_2(t) - \pi/2)$ (+100 if handstand), the discount factor is set to $\gamma = 1$, and each policy is evaluated for $H = 500$ steps.^{‡‡}

HIV drug treatment. Our last problem domain is taken from a real-world application in medical control [15]. The aim is to optimize the treatment of a patient infected by HIV over a period of a few years using what is known as structured treatment interruption (STI). The treatment of the patient consists of choosing a combination of two drugs, yielding four possible choices (including the case where no drug at all is taken), to administer every 5 days. Administering the correct cocktail with the correct timing can hinder the spread of HIV-infected cells and will eventually bring the patient into a healthy state (a locally stable equilibrium); however, the drugs also have side effects on the patient’s health, and thus their use should be kept to a minimum. Finding an optimal treatment strategy is considered a challenging optimal control problem with highly nonlinear transition dynamics [16]. The system is represented by a six-dimensional state vector $x \equiv (T_1, T_2, T_1^*, T_2^*, V, E)$, where $T_1 \geq 0$ and $T_2 \geq 0$ are the count of healthy type 1 and type 2 cells, $T_1^* \geq 0$ and $T_2^* \geq 0$ are the count of infected type 1 and type 2 cells, $V \geq 0$ is the number of free virus copies, and $E \geq 0$ is the number of immune response cells. The two-dimensional control vector $u \equiv (\varepsilon_1, \varepsilon_2)$ consists of the dosage of two drugs, which is discretized to the four actions $\{(0.3, 0.7), (0.7, 0), (0, 0.3), (0, 0)\}$ and held constant for $\Delta t = 5$ days. The dynamics of the system and physical parameters we used to instantiate the problem are specified in Appendix D.

To formulate this task as an optimal control problem of the form given in Equation (5), we used the following settings: the set of initial states is the singleton $X_0 = \{(163573, 5, 11, 945, 46, 63919, 24)\}$, which corresponds to the unhealthy locally stable equilibrium (i.e., a high number of HIV-infected cells). The reward is composed of the number of infected and uninfected cells plus an additional term reflecting the cost for using a drug: $\varrho(x(t), u(t)) = -0.1V(t) + 10,000E(t) - 20,000\varepsilon_1(t) - 20,000\varepsilon_2(t)$. The discount factor is set to $\gamma = 0.98$, and each policy is evaluated for $H = 300$ steps. Note that in this domain, unlike in all the previous ones, the reward is not upper bounded (and thus the optimistic tree-development strategy from Equation (12) cannot be applied). Moreover, the values of the state variables can vary over a large range, from 0 up to the order of 10^6 ; to counter any unwanted scaling effects in our learning methods, we transformed the state variables by taking their logarithm.

^{‡‡}Note that as $\gamma = 1$, the criterion that we optimize for this problem is the (time-independent) finite sum of rewards. Changing the criterion in this way does not make any technical difference for any of the DPS techniques.

5.2. Contestant methods

We compare OLT policies with both pure DPS approaches and pure LT approaches. From the point of view of DPS, we consider three alternative policy representations (neural networks, adaptive radial basis functions, and OLTs) and two alternative optimizers (CE and GPO). Our pure LT approaches are composed of the uniform tree-development strategy, the optimistic development strategy proposed in [2], and two greedy tree-development strategies.

DPS representations. We consider the following two common policy representations:

- *Neural network.* This representation is probably the most widely used in the DPS literature [1]. The policy is represented by a fully connected feed-forward neural network, which has one hidden layer, one input layer with one neuron per state variable, and one output layer with one neuron per possible action (plus bias neurons). Given the current state, the neural network computes one activation score per action, and the policy returns an action with maximal score. Hidden nodes have *tanh* activations, and their number is a hyperparameter that enables to control the complexity of the policy; output nodes have a linear activation function. The entries of the state vector fed to the input layer are scaled to lie in $[-1, 1]$.
- *Adaptive radial basis functions.* In this representation proposed in [17], a policy is encoded through a set of radial basis functions defined over the state space, where each basis function is attached to a particular action. Given the current state, the policy works by searching for the nearest basis function (“nearest” being measured under the Mahalanobis distance metric) and by returning the action attached to it. Each single basis function is parameterized by the location of its center, the length scales along each dimension, and the recommended action. The entries of the state vector fed to the policy are raw values because the scale is incorporated into the adaptive length scales. The number of basis functions used to encode the policy is a complexity hyperparameter that has to be tuned externally.

Table I summarizes the characteristics of the three policy representations that we consider in this paper. Each representation has a hyperparameter that controls the complexity of the policy. Note that OLT is unique in the sense that the number of parameters in the representation—the dimensionality of the search space over which global optimization is performed—does not depend on the value of this hyperparameter.

DPS optimizers. CE [7] is a versatile global optimization technique that is widely used in DPS [1]. In particular, this method was used in previous work [3], under the alternative name of *estimation of distribution algorithm*. The chief advantage of CE (and the main reason for it being highly popular) is that it is very easy to implement and, because it is population based, can deal with a potentially very large number of sample locations (i.e., CE operates only on small batches of sample locations the size of which is constant throughout, whereas GPO has to consider all current sample locations to determine the next one). The algorithm works by fitting a distribution to the best currently found

Table I. Summary of the policy representations compared in this paper.

Name	Policy representation	Hyperparameter	Number of parameters
Neural network (NN)	Feed-forward neural network with tanh activation and one hidden layer	Number of hidden neurons (<i>nHidden</i>)	$(n_x + 1) \times nHidden + (nHidden + 1) \times K$
Radial basis functions (RBF)	Adaptive radial basis functions with Mahalanobis distance as in [17]	Number of basis functions (<i>nBF</i>)	$2 \times n_x \times nBF$
Optimized look-ahead trees (OLT)	Look-ahead tree (algo 2) with optimized node scoring heuristic $e(\cdot; \theta)$ as in Equation (13)	Budget of node expansions (h_{max})	$3 \times n_x$

n_x denotes the dimensionality of the state space, and K is the number of actions.

solutions and by using this distribution to sample new candidate solutions. In our case, we use a simple variant of CE that relies on a multi-variate Gaussian distribution with *diagonal* covariance matrix. This distribution is first initialized to cover the whole search space: we start with a Gaussian with zero mean and a diagonal covariance matrix, the entries of which are equal to the square of the half-length of an interval centered at zero and containing the corresponding coordinate of the search space. The algorithm then draws a number N_{CE} of observations from this distribution, where each observation is a parameter vector and represents a possible policy, evaluates the N_{CE} resulting policies, and sorts them according to their performance. It then picks a number M_{CE} of the top best-performing policies and uses them to update the generating distribution: the old mean is replaced by the sample mean, and each diagonal entry of the covariance matrix is replaced by the per-coordinate variance of the M_{CE} best parameter vectors. These two steps are iterated until a stopping condition, in our case a predefined maximum number of iterations, is reached.

LT strategies. Classical LT policies rely on fixed generic node scoring functions. In order to demonstrate the benefits of learning the node scoring function in a problem-driven way, we compare our approach OLT to LT policies by using e^{uniform} (uniform tree development, see Equation (11)), $e^{\text{optimistic}}$ (the method proposed in [2], see Equation (12)) and two forms of greedy tree development:

$$e^{\text{greedy-1}}(n_{i,h}) := Q_{t+h} \qquad e^{\text{greedy-2}}(n_{i,h}) := \gamma^h Q_{t+h}. \quad (26)$$

5.3. Performance comparison

We start by comparing the performance we obtain with the different approaches discussed previously. Note that, contrarily to OLT, NN and RBF typically involve solving challenging global optimization problems with tens or hundreds of parameters. Whereas GPO is quite efficient for problems that have a reasonable number of parameters, this approach requires sophisticated approximations to scale to higher-dimensional problems (e.g., see [18]). In this paper, we use a naive textbook implementation of GPO that is able to solve OLT optimization problems, but that suffers from scaling problems when the number of samples increases, which happens to be problematic in high-dimensional problems. Therefore, we use CE to learn NN-based and RBF-based policies. To make the comparison fair, we also use CE to optimize the parameters of OLT policies in this first part of our empirical study. A comparison between CE and GPO for learning OLT policies is provided in Section 5.4.

Experimental protocol. We use the same test procedure for all policies, and the performance is measured as the discounted sum of rewards obtained when executing the policy for H steps (Section 5.1), starting from the initial state^{§§} x_0 .

The parameters of CE were tuned by hand so as to give in each case a good result with a reasonable amount of computation. The result of this tuning is given in Table II. Note that, as OLT involves lower-dimensional optimization problems than the two other representations, the required optimization budget ($N_{\text{CE}} \times n_{\text{Iter}}$) is lower on all problems for this representation. For the RBF representation, the CE has to optimize both continuous parameters (RBF centers and length scales) and discrete parameters (action assignments). To optimize these discrete parameters, we chose a multinomial distribution with Dirichlet prior where the initial count of each action was set to 10.

The hyperparameters of the three policy representations were tuned by grid search. For both the number of hidden nodes in the NN representation and the number of basis functions in the RBF representation, we tested the following values: $\{1, 2, 3, 4, 5, 7, 10, 12, 15, 20, 30, 40, 50\}$. As in [2, 3], we tested budget values for LT policies that correspond to the number of nodes of fully developed trees of varying depth $d = 1, \dots, 8$. These budgets values are as follows: $h_{\text{max}} \in \{1, 1 + K, 1 + K + K^2, 1 + K + K^2 + K^3, \dots\}$.

^{§§}Recall that for each of our four benchmark domains, the set X_0 only contains one single initial state x_0 . Yet as we will see below, robustness with respect to perturbations of the initial state is one of the strengths of LT policies, so that optimizing only over one initial state is justified.

Table II. Cross-entropy parameters for the four problems and the three policy representations.

	Inverted pendulum			Double inverted pendulum			Acrobot handstand			HIV drug treatment		
	N_{CE}	M_{CE}	$nIter$	N_{CE}	M_{CE}	$nIter$	N_{CE}	M_{CE}	$nIter$	N_{CE}	M_{CE}	$nIter$
NN	400	20	100	400	20	200	1000	20	50	400	20	200
RBF	200	20	100	500	10	200	500	10	200	500	20	200
OLT	100	10	25	100	10	50	100	10	50	100	10	50

Table III. Quality of learned policies with neural networks, radial basis functions, and optimized look-ahead trees.

		Inverted pendulum	Double inverted pendulum	Acrobot handstand	HIV drug treatment
		NN	$nHidden$	50	50
	Performance	89.9	134.9	3.89e4	0.88e9
RBF	nBF	30	6	30	15
	Performance	73.4	122.5	4.09e4	1.39e9
OLT	Budget	31	1365	40	85
	Performance	93.2	145.2	4.07e4	4.22e9

For each problem, the best performance is shown in bold. For each method and each problem, we also display the value of the tuned hyperparameter.

Performance comparison. Table III reports the best scores obtained by the tree kinds of policies as well as the values of the tuned hyperparameters. We observe that OLT works significantly better than the other DPS methods on three out of the four benchmark domains, which demonstrates the usefulness of the LT-based representation over the flat representations used in traditional DPS. On the HIV domain, OLT enables us to obtain even better results than the current state-of-the-art [16] (4.22e9 against 4.16e9), whereas both of the two other representations only manage to reach policies with very moderate scores ($\approx 1e9$). RBF slightly outperforms OLT on the acrobot domain. However, as we will see later, this result holds thanks to a very careful tuning of the number of radial basis functions, whereas OLT works well for a wide range of budget values.

Learning the node scoring function and online efficiency. We have seen that the OLT policy representation enables to reach policies outperforming those obtained with neural networks and radial basis functions. One could wonder whether this result could be obtained by using LTs without learning. We therefore performed a series of experiments that, for various budget values, compare the performance of LT policies, with and without learning. Note that in these experiments, the budget is set before optimization; hence, we optimize one OLT policy per tested budget value.

The results of our comparison between OLT policies and traditional LT policies are given in Figure 6. OLT policies achieve a better score on three domains out of the four (inverted pendulum, acrobot, and HIV). Furthermore, thanks to learning, OLT requires a significantly smaller budget and hence online computational resources to reach a given level of performance, again on three domains out of the four (double inverted pendulum, acrobot, and HIV). The most impressive results are obtained on the HIV domain, for which an OLT policy with a budget of two node expansions already performs better than all other LT policies with a budget up to $h_{max} = 87,381$ node expansions (which corresponds to fully developed trees of depth 8). An OLT with such a small budget consists in first expanding the node corresponding to the current state and then expanding one of its successor nodes. Our results show that carefully selecting this successor state is much more efficient to solve the HIV problem than developing large trees in a generic way. In the same spirit, on the acrobot problem, an OLT policy with only eight node expansions performs slightly better than a uniformly developed tree with a budget of 9841 nodes (corresponding to a tree of depth of 7).

On the double inverted pendulum problem, OLTs are ultimately outperformed by LT policies using u -score. However, we observe that our approach is much better able to deal with a constrained

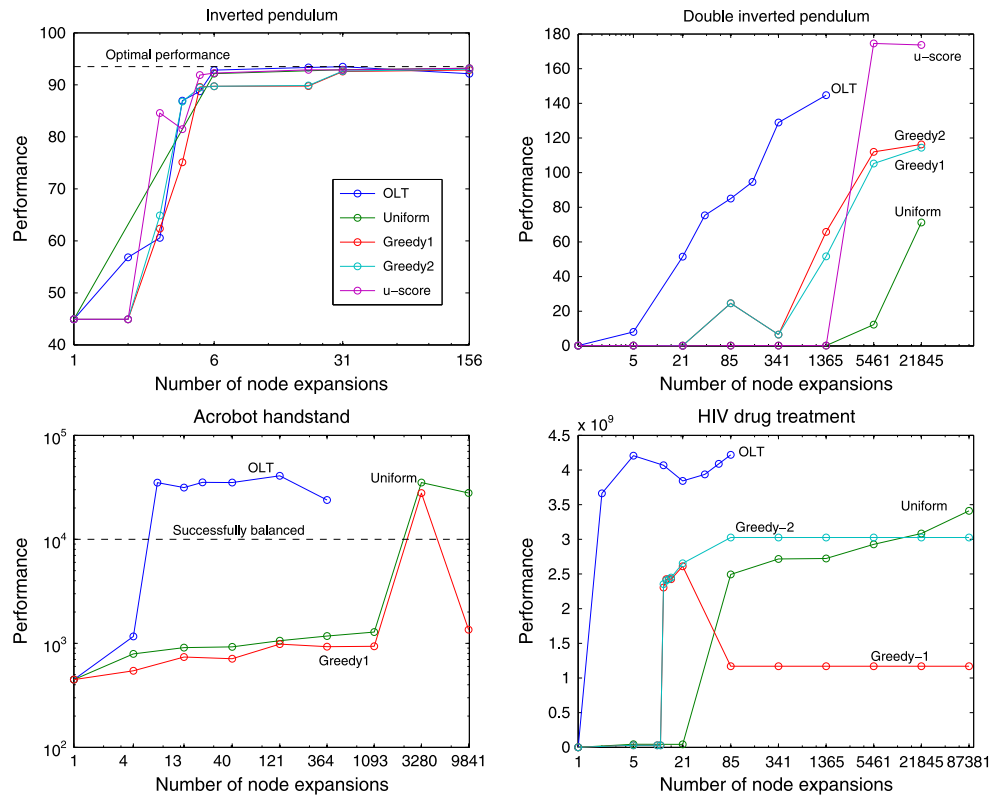


Figure 6. Performance of optimized look-ahead tree policies (OLT) versus baseline look-ahead tree policies for various budget values. As $\gamma = 1$ in the acrobot domain, the greedy-2 policy degenerates to greedy-1 and the u -score-based policy degenerates to the uniform policy. As the reward is not upper bounded in the HIV domain, the u -score is not defined for this problem. As explained in the text, because the time required to make a single decision increases linearly in the number of node expansions and because we have to evaluate many more policies to produce a single OLT curve than we have to produce an LT curve, we did not evaluate OLT on the same large budget values as LT.

computational budget. For example, using only 341 node expansions, we obtain a reasonably well-performing policy, that outperforms all other generic tree-development strategies even when they have very large budgets.

On the inverted pendulum problem, we see that all LT policies achieve a near-optimal performance (which we can compute for this domain, e.g., see [19]) already with a budget of 5. It thus seems that there is little interest in learning a specific node scoring function for this problem.

Illustration of the HIV policy. In order to allow a direct comparison between the performance of our method on the HIV domain with the performance given in the earlier related work [16, 17, 20], we plot in Figure 7 the trajectory that we obtain on the HIV system. These results show that our policy applies reverse transcriptase inhibitor (RTI) and protease inhibitor (PI) drugs in a way which is very similar to what other state-of-the-art policies do (which, however, are obtained in a fundamentally different manner, e.g., by fitted Q -value iteration using millions of sample transitions).

5.4. Offline complexity and sample efficiency

We have seen that the OLT representation enables to reach high-performance policies, which often outperform alternative DPS representations. Besides performance, another aspect of major

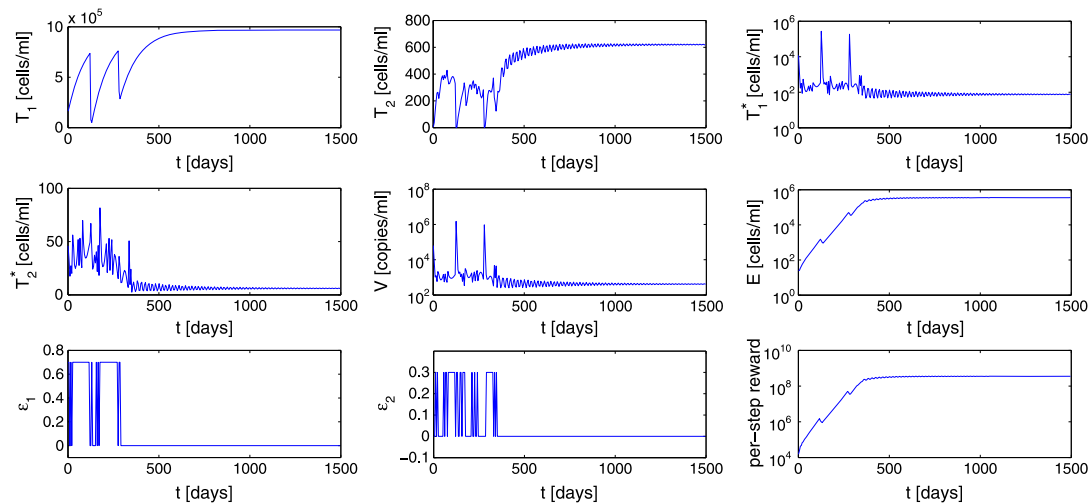


Figure 7. Trajectory of the HIV system when controlled by an OLT policy with a budget of 85 node expansions. The first six panels show the development of the six system states, the following two panels show the dosage of the RTI and PI drugs applied, and the last panel shows the reward obtained in each step.

importance in DPS is the computational complexity of the learning process, that is, how fast good policies can be obtained. We study the offline computational complexity by looking at the performance of our different methods as a function of two metrics: the number of policy evaluations performed and the number of transitions simulated.

Experimental protocol. The most common solution to measure sample efficiency in a DPS scheme is to look at the number of policy evaluations required to reach a certain level of performance. This measure corresponds to the number of different parameter values that have been tried by the optimization algorithm. Here, we consider two different optimization algorithms: CE and GPO. CE is tuned as previously and GPO is instantiated as follows: as kernel, we chose, as it is standard for GP regression, the squared exponential given in Equation (18), where the hyperparameters are found for each batch of data via marginal likelihood optimization [11]. In this optimization, the best setting of the hyperparameters found in the previous iteration is used as the mean of the hyperprior in the next iteration. To generate an initial batch of training data, we generated 10 samples (100 samples in the double inverted pendulum domain) via Latin hypercube sampling; these initial samples are taken into account when we compare sample efficiency. To find at each iteration of GPO the best next sample location, we optimize the EI acquisition function from Equation (25) with $\zeta = 0.01$ using DIRECT.

We are interested in two questions: ‘How does GPO compare against CE when using the OLT representation?’ and ‘How does the choice of representation impact the sample efficiency?’ To answer these questions, we selected the three setups (OLT with GPO, OLT with CE, and RBF with CE) and plotted the corresponding learning curves for different values of the h_{\max} and nBF hyperparameters. These plots are given in Figure 8.

Impact of the optimizer. We observe that GPO requires about one order of magnitude fewer policy evaluations than CE on all four problems (note that the GPO curves stop at 500 policy evaluations, whereas we display CE results until 5000 evaluations). We observe no notable difference between the performance of the final policies obtained with GPO and CE, and the overall shapes of the learning curves are similar. Note that CE is much easier to implement than GPO. Therefore, in the light of these results, we suggest using CE when sample complexity is not a big problem and selecting GPO in the case where the number of policy evaluations is a strong applicative constraint.

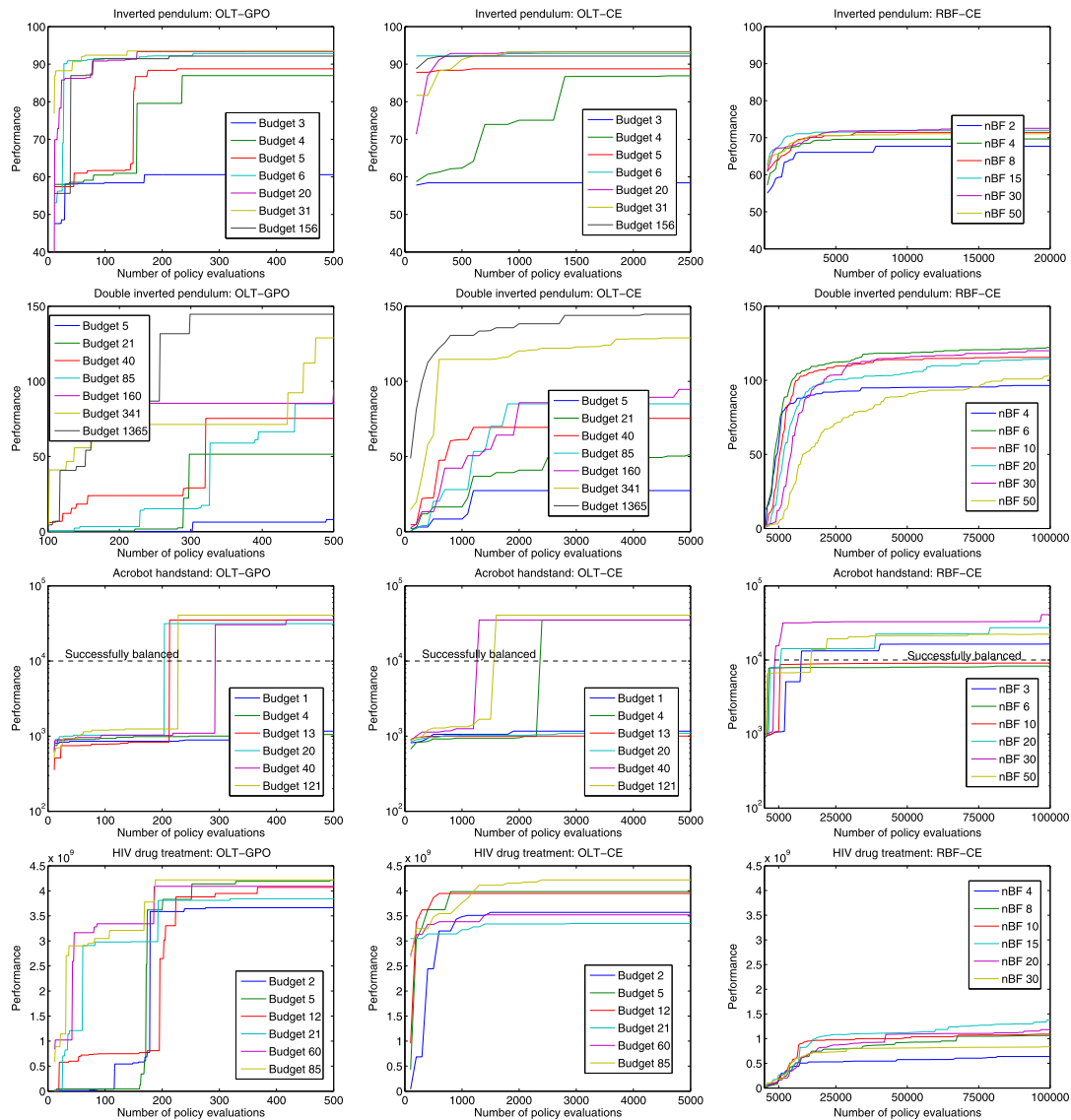


Figure 8. Performance versus number of policy evaluations. From left to right: OLT with GPO, OLT with CE, and RBF with CE. From top to bottom: inverted pendulum, double inverted pendulum, acrobot handstand, and HIV drug treatment.

Using GPO requires a bit more work, but is profitable because it enables to roughly perform learning one order of magnitude faster.

Impact of the representation. Our results show that optimizing LTs is substantially more sample efficient than optimizing the parameters in a basis function representation, this difference being 5K versus 50K–100K policy evaluations when using the CE method. The same kinds of differences are observed when comparing with the NN representation. One obvious explanation for this is that optimizing the node scoring function involves fewer parameters, which results in the search space having lower dimensionality (cf. column 4 of Table I). As mentioned previously, a key advantage of OLT is that the complexity of the policy (i.e., the computational budget h_{\max}) can be arbitrarily scaled without increasing the number of parameters to be optimized. In contrast, for both the neural network and adaptive basis function representation, higher complexity can only be achieved at

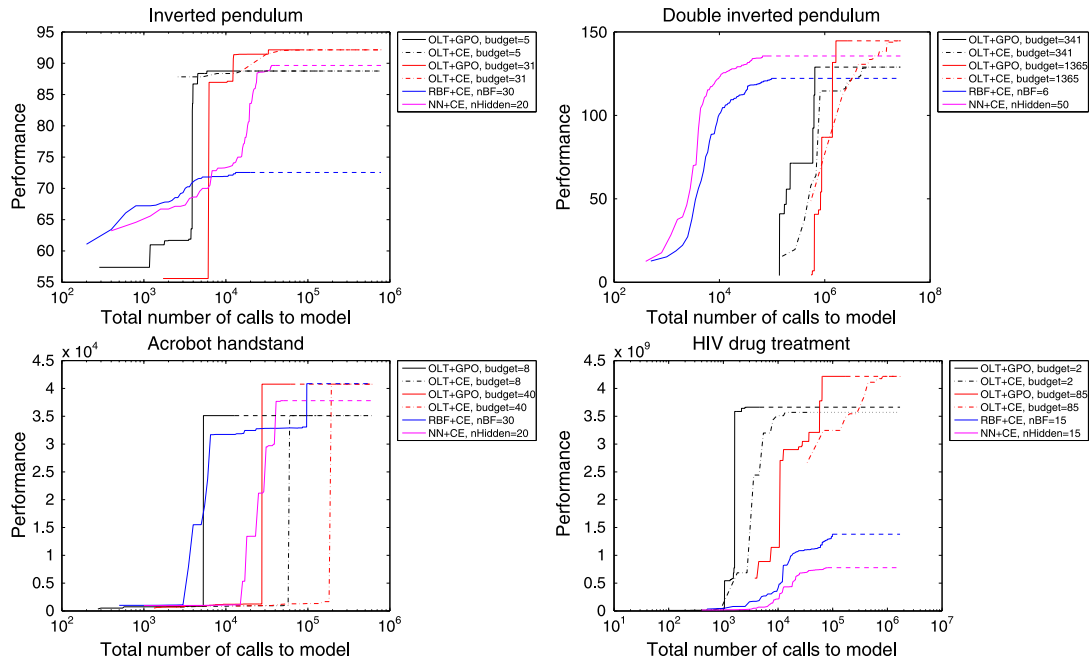


Figure 9. Performance versus number of simulated transitions.

the expense of a larger number of policy parameters, which generally involves harder optimization problems.

Performance versus number of simulated transitions. An important characteristic of OLT policies is that they use more online computational resources than NN and RBF policies to take their decisions. Specifically, in order to take a single decision, an LT policy requires simulating Kh_{\max} transitions using a generative model of the problem, because computing the policy involves expanding h_{\max} nodes and expanding a node invokes the model once per possible action. Hence, it is also important to study the learning performance with respect to the *number of simulated transitions*. This new sample complexity measure is computed in the following way. Note that one policy evaluation requires making a trajectory of H steps (assuming we have a unique training initial state) and that the simulator is called once per step of this trajectory. The number of simulated transitions per policy evaluation is thus H for NN and RBF policies. If we add the additional simulation cost of OLT policies, the number of simulated transitions per policy evaluation becomes $H \times (1 + Kh_{\max})$. We compare our three policy representations with the tuned hyperparameters given in Table III. As OLT policies with high budget values are strongly disadvantaged w.r.t. the number of simulated transitions, we consider an additional *low-budget* OLT setting, where the budget values were chosen by hand to be as small as possible while still producing reasonably good policies. Figure 9 reports the learning curves obtained by CE for the NN and RBF settings and by both GPO and CE for the two OLT settings.

We observe slightly different behaviors depending on the problem. In the first two problems, the best policies are obtained with NN and RBF for small number of simulated transitions, and after a given threshold, OLT policies become better. This was expected as evaluating a single OLT policy can already require simulating thousands of transitions. In the acrobot and HIV problems, OLT policies outperform the other ones on almost the whole x -axis range, which was more unexpected. This proves that the larger number of transitions required to take decisions is partly compensated by the smaller number of parameters to optimize.

On all four problems, we observe that (i) low-budget OLT policies outperform tuned OLT policies when the number of simulated transitions is small, (ii) policies optimized with GPO always

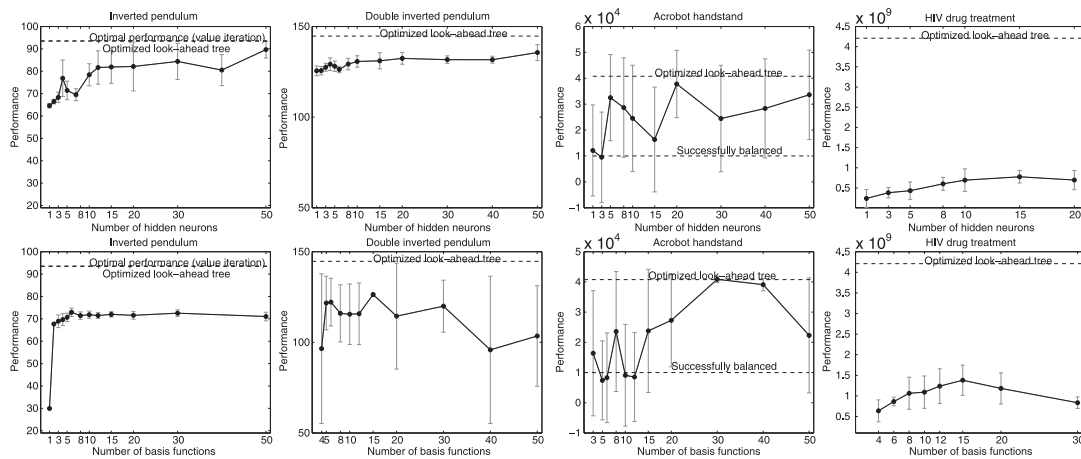


Figure 10. Performance versus hyperparameter value. Top: Performance of NN as a function of the number of hidden neurons. Bottom: Performance of RBF as a function of the number of basis functions.

outperform their counterparts optimized with CE, and (iii) the asymptotically best (or almost best in the case of acrobot) policies are based on OLTs.

5.5. Robustness with respect to hyperparameters

In addition to performance and learning efficiency, another property of DPS approaches particularly relevant for practical applications is that their hyperparameters should be easy to tune. Indeed, tuning these parameters is part of the whole policy learning process and contributes to the offline complexity; hence, the more trial and error is required, the longer it takes to obtain high-performance policies. The behavior of OLT policies as a function of the h_{\max} hyperparameter was shown in Figure 6. We report the performance of the other two kinds of policies as a function of their respective hyperparameters in Figure 10. As CE is a stochastic optimization algorithm, we performed these experiments ten times for each setting and report the empirical means and standard deviations.

Depending on the problem domain, NN and RBF can be quite sensitive to the setting of their complexity parameter. For both methods, there are some problems for which tuning the complexity hyperparameter is rather easy and others for which it is much harder. We observe that on the acrobot benchmark, which is the only one where NN and RBF are able to produce policies competitive with OLT policies, both the performance of NN and RBF policies have large variances: performance varies strongly both for increasing the complexity (number of nodes/basis functions) or across different runs in the same setting. To understand this large variance, it should be noted that acrobot is different from the previous domains in that there is no partial solution: a policy is either able to swing up and successfully balance in the handstand position indefinitely (the reward is $> 10^4$) or not at all (the reward is of order 10^3 at most). Successful balance is only possible if the system approaches the tiny region in the state space from where LQR can take over.

Although we observe in some cases that good policies can be found using a fairly compact representation, our results indicate that, seen across all domains, finding the best possible complexity parameter for NN and RBF requires a large amount of trial-and-error experimentation. On the other hand, as is shown in Figure 6, OLTs are rather easy to tune, because, in general, increasing the budget increases the quality of the policy or only slightly degrades it. In all our experiments, a default value corresponding to a fully developed tree of depth 3 or 4 yielded good policies.

5.6. Robustness with respect to initial states

In DPS, policies are optimized w.r.t. a particular distribution over initial states. An important issue is that this distribution may not perfectly match future unexpected usage of the learned policy. In

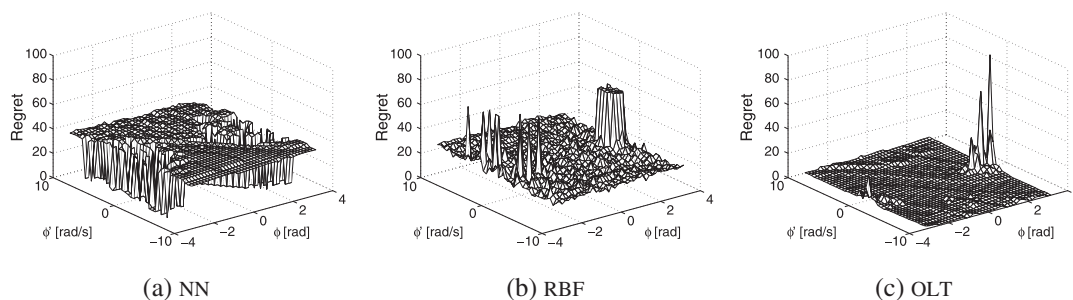


Figure 11. Regret for NN, RBF, and OLT policies on the inverted pendulum domain when the testing initial state differs from the training initial state.

general, preference should be given to policies that are robust w.r.t. such mismatches. To compare the robustness of the NN, RBF, and OLT policies and to study how well they are able to generalize when starting from initial states, which differ from those used during learning, we performed experiments by evaluating them on perturbed or completely different initial states on the inverted pendulum problem and the HIV drug treatment problem.

Stability on inverted pendulum. To analyze the stability of the various policies on this domain, we started by learning NN, RBF, and OLT policies with the single initial state $x_0 = (-\pi, 0)$ and then evaluated these policies across the whole domain by systematically varying the initial state. Remember that in the pendulum domain, other than in the higher-dimensional domains, we are actually able to compute the optimal policy (using a high-resolution grid). We use this latter policy to report the *regret* for each initial state, that is, the difference between the optimal performance and performance obtained from the learned policy. The results of these experiments are given in Figure 11. It can be seen that an OLT policy (budget 20) performs well and incurs close to zero regret across large parts of the state space, with an exception being the region around $(\pi, 0)$ (which can be explained by the choice of the parametrization in Equation (13) in which π and $-\pi$ are maximally separated). When comparing with NN and RBF policies, we see that both representations are less able to generalize and can incur comparatively large regrets even for small perturbations.

Stability on HIV drug treatment. Figure 12 illustrates the robustness of an OLT policy (budget 85) in the HIV domain. This time we are no longer able to exhaustively sample the state space; instead, we perturb each state variable in turn by multiplying each variable of x_0 independently by a factor ranging from 0.1 to 10. The figure then shows that the performance stays largely stable even for large perturbations.^{¶¶} In particular, even for large perturbations, the performance of the OLT policy remains significantly higher than the best results obtained by NN and RBF ($>3e9$ vs. $\approx 1e9$).

5.7. Robustness with respect to the budget

With DPS, OLT policies are optimized w.r.t. a fixed budget of allowed node expansions. We now examine how robust the solution is when the budget used during evaluation of a policy differs from the budget used to learn it. For each of the four domains, we took two or three representative OLT policies and evaluated them with various evaluation budgets. Figure 13 shows the results of this experiment.

^{¶¶}Note that in doing this we did not pay attention to the physical meaning of the variables; thus by increasing the value of the sixth state variable, it appears as if the policy begins to perform even better. However, as this state variable denotes the count of healthy cells in the blood of a patient, increasing its value by an order of magnitude alters the whole problem.

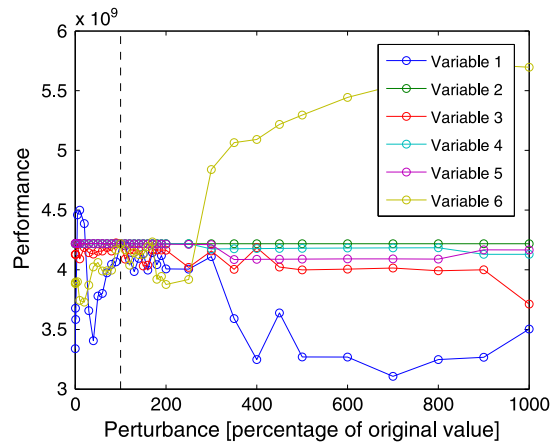


Figure 12. Robustness of an OLT policy for the HIV drug treatment domain. The plot examines how the performance of the policy (with budget 85) decreases when the initial state for which the policy was optimized is perturbed: each curve corresponds to the perturbation of one state variable independently. The vertical line at $x = 100$ denotes the unperturbed value.

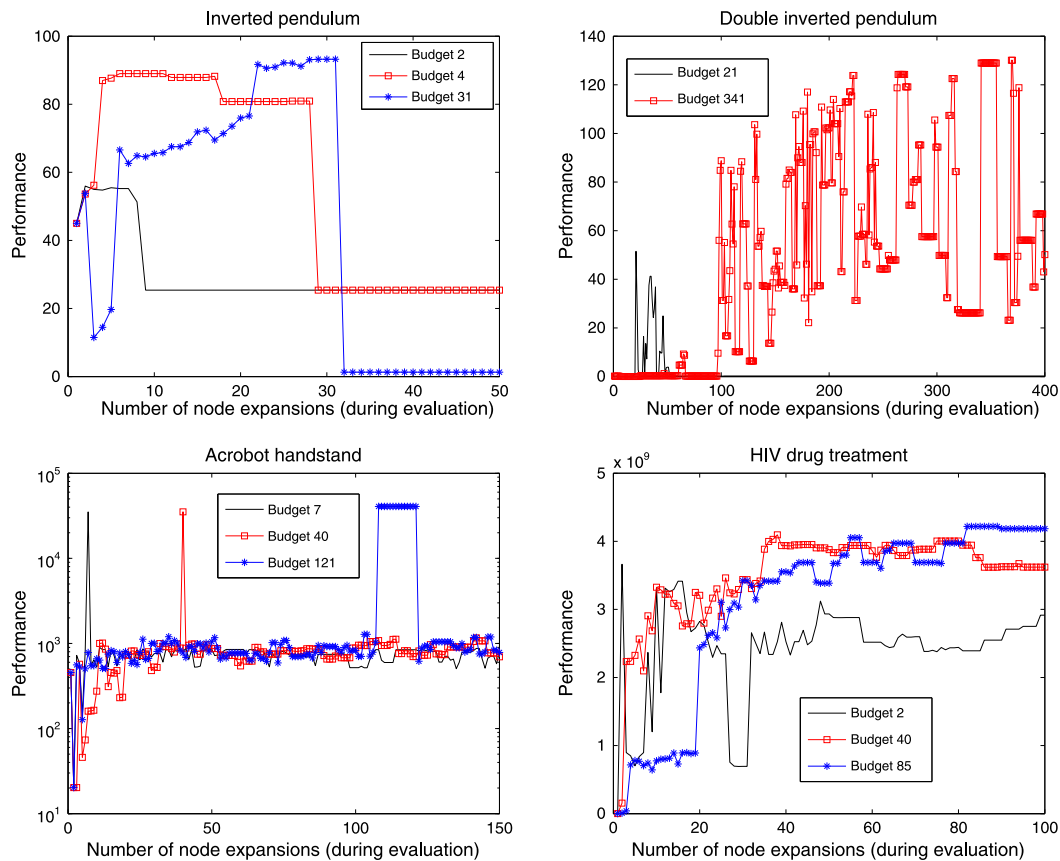


Figure 13. Robustness of OLT when the budget of allowed node expansions used during evaluation of a policy differs from the budget that was used to learn it. Each curve corresponds to a policy that was first optimized for a fixed budget of node expansions (as indicated by the legend) and then later evaluated with the number of node expansions varying as indicated by the x -axis.

We observe different kinds of behaviors. In some cases, OLT policies are rather robust w.r.t. changes in the evaluation budget. As an example, the HIV policy learned with a budget of 40 works reasonably well for evaluation budgets ranging from 10 to 100. In some other cases, small changes in the budget lead to fast performance drops. The acrobot domain leads to an extreme case of this phenomenon, in which, as an example, the OLT policy trained for a budget of 40 only works well if the evaluation budget is also 40. These results thus stress the fact that OLT policies are optimized for a given online budget and that the same budget should be available during both learning and evaluation.

6. RELATED WORK

The approach proposed in this paper lies at the intersection of two families of solutions for sequential decision making. We first overview DPS in Section 6.1 and then discuss relevant work on LT search in Section 6.2. Section 6.3 positions our approach with respect to model predictive control (MPC) techniques. Finally, Section 6.4 suggests that OLT policies belong to a larger emergent class of techniques: *parameterized algorithms for decision making*.

6.1. Direct policy search

Direct policy search is a widely used class of solutions that comes from RL. For a general overview over the field of RL, refer to one of the books in [1, 6, 21, 22]. Over the years, many DPS techniques have been proposed and giving credit to every one of them would be hardly possible. Individual techniques differentiate themselves by their policy parametrization and the optimization method used for identifying parameters leading to a high-performing policy. The main distinction is whether the policy search is gradient free or gradient based.

Gradient-based DPS, which very often also goes by the name of policy-gradient method, follows an iterative optimization scheme and uses the gradient of the value function to adapt and change the policy parameters such that performance increases. This requires the value function to be a smooth function of the policy parameters, which is typically achieved by considering compatible stochastic policies. Often however, the gradient cannot be computed analytically because the value function itself is not available in closed form. Instead, the gradient has to be estimated using, for example, finite difference approximation, Monte Carlo rollouts, or value function approximation, which leads to the actor-critic methods. One of the best known examples is probably the natural policy-gradient method described in [23]; another more recent example is the work in [24]. Although impressive results have been obtained in particular for learning controllers in robotics, policy-gradient methods do have some weaknesses: the most notable ones being a high sensitivity to the initial value of the policy parameters (the starting point of the iteration), abundance of local minima, and difficulties when the return is noisy or stochastic. Note that in this paper we do not deal with gradient-based policy search.

Gradient-free DPS on the other hand finds the best policy parameters via derivative-free global optimization. Its main strength is simplicity and generality; the latter meaning that, because it is derivative free, the representation of a policy by any algorithm with adjustable parameters is admissible, and not only those for which the value function will be differentiable. As gradient-free DPS techniques perform a global search in the parameters space, they do not suffer from local minima or from the problem of having to guess a good initial solution; these techniques were also shown to cope well with stochastic returns and hidden states. For certain domains in RL such as Tetris, the best performing policies known today have been obtained by gradient-free DPS (see [25] and follow-up work). The weakness of this approach is that conceptually it is less sample efficient than policy-gradient methods and thus will require a substantially higher number of policy evaluations. This comes as a consequence of having to solve a global optimization problem over a potentially high-dimensional search space; however, by carefully designing the policy parametrization and choosing powerful optimizers, sample efficiency is often improved.

Early examples of DPS came under the guise of evolutionary approaches for RL, or neuroevolution, and consisted of neural networks as policy representation, with the weights making up the

policy parameters, and variants of genetic algorithms acting as global optimizer. Examples can be found in [26, 27], where later work also considered optimization of the network structure [28] or using recurrent neural networks to better cope with hidden states [29]. A recent comparison of these methods can also be found in [30, 31].

As an alternative to genetic algorithms, more recent work started to explore the use of the CE method [7], or variants such as the covariance matrix adaptation evolution strategy CMA-ES [8] as global optimizer for policy search. Examples include [32], where the policy was represented by a simple linearly parametrized function; [31], where the policy was represented by a neural network; or [17], where the policy was represented by adaptive radial basis functions. These latter two were also used in this paper as baseline methods during the experimental evaluation of our LTs. Another example for a policy parametrization is to use domain-specific building blocks, such as motor primitives, as was done in [5, 33] to optimize the gait of the AIBO quadrupedal robot. A different kind of policy representation is used in [34] to learn a policy for the game of Ms. Pac-Man; here the policy is represented by a list of domain-specific parameterized rules.

Section 4 described a third option for the global optimization part: GPO [10, 13]. GPO can achieve very good sample efficiency and was previously considered for policy search in [4, 5]. In the end, however, the choice of which global optimizer one uses will always be secondary; what is more important is how the policy is represented and to what extent this representation facilitates the optimization process by shaping the ‘fitness landscape’.

6.2. Look-ahead tree search

The algorithm studied in this paper can also be related to the larger field of tree-based planning and search. One of the most seminal works in this field is the A^* algorithm [35], which uses a best-first search to find the shortest path from a source state/configuration to a goal state/configuration. The conceptual difference between A^* and related methods on the one side and what we are presenting here on the other side is that our method implements an *online* mechanism and a *finite computational budget* mechanism: a search tree is grown at each decision-making step, and only a finite number of node expansions is allowed during tree development (the larger the computational budget, the better the decision). In A^* , the search tree needs to be grown until a goal state is reached. Our method is thus capable of online planning and producing closed-loop policies, whereas A^* is not. More specifically, the algorithm studied in this paper can also be interpreted as a method for learning an exploration strategy in a tree. In A^* , the function used for evaluating the nodes is the sum of two terms: the length of the so far shortest path from the source to the current node and an optimistic estimate of the shortest path from this node to the goal (a so-called admissible heuristic). Several authors have sought to learn good admissible heuristics for the A^* algorithm. For example, we can mention the LRTA* algorithm [36], which is a variant of A^* and which learns over multiple trials an optimal admissible heuristic. More recent work for learning strategies to efficiently explore graphs have focused on the use of supervised regression techniques by using various approximation structures to solve this problem (e.g., linear regression, neural networks, and k -nearest neighbors); for example, see [35, 37–39].

Another noteworthy example for combining learned heuristic functions and LT search (in this case IDA* and a bootstrapping procedure) is the recent work in [40]. However, it should be noted that this approach is not directly comparable as it is applied to search problems and puzzles, and not, as we do here, to continuous optimal control tasks with its more general concept of per-step costs and rewards.

6.3. Model predictive control

Model predictive control techniques have originally been introduced as ways to stabilize large-scale systems with constraints around equilibrium points (or around a reference trajectory) [41–43]. They exploit an explicitly formulated model of the problem and solve in a receding horizon manner a series of finite-time, open-loop deterministic optimal control problems. In such, they are very much related to LT techniques. Actually, an MPC technique that searches for the first action of the optimal sequence of actions over the finite optimization horizon through (clever) exploration of a tree is an

LT technique. However, most of the control techniques labeled as MPC techniques do not use tree exploration for identifying (near-) optimal sequences of actions. They rather assume strong regularity assumptions on the system dynamics and the reward function (e.g., linearity)—that we do not make here—which are exploited to reformulate the search for an optimal open-loop sequence of action as a standard mathematical programming problem (e.g., convex optimization problem and mixed integer programming problem). This reformulation of the problem as a mathematical programming problem often leads to techniques that are able to scale well to very large state–action spaces, provided of course that the regularity assumptions hold.

6.4. Parameterized algorithms for decision making

We see our work here as an application of a more generic meta-algorithm and research methodology: for a given targeted class of problems, (i) identify an algorithm skeleton that is believed to provide promising solutions to problem instances, (ii) parameterize this algorithm, and (iii) optimize the parameters in a problem-driven way, through DPS, that is, through direct optimization of the algorithm performance. This methodology can be applied to a wide range of problem kinds and has already inspired several authors. The system proposed in [44] places an optimization layer on top of an approximate value iteration algorithm to optimize the location of its basis functions. In [45] (Section 5.3), the authors consider multi-stage stochastic programming techniques and optimize the scenario trees by using Monte Carlo methods. Closer to our work, it is proposed in [46] to parameterize a tree-search technique for decision making: *upper confidence trees*. In this work, the parameters enable to control the simulation policy used to estimate long-term returns within the upper confidence tree algorithm.

Parameterized algorithms have also been shown to be relevant to solve various kinds of exploration/exploitation dilemma in a problem-driven way. References [47, 48] propose to learn exploration/exploitation strategies for multi-armed bandit problems either by using the same kind of parameterizations (a simple linear function) and the same kind of optimizers (derivative-free global optimizers) as ours or by searching in a space of formulas. Reference [49] extends this idea to the exploration/exploitation dilemma that occurs in single-trajectory RL. Here, the parameters are no longer real-valued vectors but rather small formulas that depend on the internal variables of the policy. Given a target class of Markov decision problems, this method is able to learn high-performance policies that outperform state-of-the-art generic RL algorithms.

We believe that with the recent progress in derivative-free global optimization (with algorithms such as CE, CMA-ES, and GPO) and the continuously growing available computing power, the class of such parameterized algorithms optimized in a problem-driven way is likely to play a more prominent role in the near future, as they offer a systematic way to improve upon generic solutions, by exploiting problem-dependent characteristics through learning.

7. DISCUSSION AND FUTURE WORK

This paper has focused on a particular kind of sequential decision-making problems: finite actions, deterministic transitions, known transition and reward function (i.e., a generative black-box model, which allows us to simulate arbitrary transitions), and an ‘informative’ reward. However, it relies on a generic idea—learning online strategies in an offline way—that we believe to be relevant to several other kinds of sequential decision-making problems.

Optimized LTs can be extended in many different ways by applying the following general methodology: (i) choose a tree-based decision algorithm adapted to the target class of problems, (ii) parameterize its exploration strategy, and (iii) learn the parameters through DPS. We now detail various possible extensions of OLTs.

- *Stochastic transitions.* Reference [20] extended the work on optimistic planning [2] to deal with sparsely stochastic systems (sparsely meaning that there is only a small number of possible successor states for each state/action pair). In this extension, expanding a node means simulating all possible one-step transitions for every action. Expanding OLTs in a similar way

is rather immediate and mainly impacts the computation of ℓ -score of nodes used for taking final decisions (Section 3.2).

- *Sparse rewards.* In OLTs, it is assumed that short-term sums of rewards are sufficiently informative to take high-quality decisions. This assumption can be problematic in many decision problems. For example, in games, a single reward, such as +1 in case of win or -1 in case of loss, is only observed at the end of the game. In the field of game playing, Monte Carlo tree-search (MCTS) techniques [50] have raised a huge amount of interest recently thanks to breakthrough results in computer Go [51]. The tree-development strategy of OLTs is the equivalent of the so-called selection phase of MCTS. The successful results presented in this paper thus suggest to parameterize the selection policy used in MCTS and to learn it through DPS.
- *Large number of actions.* Large number of actions means that we can no longer exhaustively generate all successor states when predicting one step ahead; instead, one would have to use sampling techniques or progressive widening techniques, which is also an active research topic in tree-based search [51–53]. Some of these techniques have also been proposed in the context of continuous actions for RL [54–57].
- *Model learning.* To what extent OLTs can be useful for solving sequential decision-making problems in which the model is unknown (the classic RL setting) is an open question. We believe that OLTs could in principle deal with imperfect models, by learning a node scoring function that compensates for the weaknesses of the model. We could thus imagine combining OLTs with an online model learner (e.g., see [19]).
- *Choice of parameterization.* Although the current work already achieves good results with the node scoring heuristic being just linearly parameterized as in Equation (13), one can easily imagine using more powerful nonlinear parameterizations, for example, neural networks or radial basis function networks. Similarly, one could also explore the benefit of including additional features when computing the node scores, for example, the u -score.
- *OLT and RL.* Another promising direction that remains to be explored is how OLTs can be combined with value function-based RL. As the value function estimated by RL could also act as a node scoring heuristic (in fact, the node scoring heuristic can be seen as a surrogate for the value function with the difference being that the former is ranking nodes for expanding an LT over multiple steps into the future, whereas the latter is directly ranking nodes in a one-step look-ahead), combining OLTs with RL could produce further improvements.
- *Varied applications.* The current paper explores OLT from the point of view of DPS and studies its properties in the context of deterministic optimal control tasks, such as pole balancing. Coming from the point of view of LT, another interesting application area would be planning and search, and in particular the domain of games (cf. item *sparse rewards* above).

8. CONCLUSION

This paper has introduced OLT policies, a novel model-based technique that bridges the gap between two major families of solutions for solving optimal control problems: DPS and LT policies. Our approach uses a new way of representing a policy by using parameterized LTs. In this representation, the parameters of the policy over which optimization takes place encode the node scoring function by which the tree is grown (until a predefined computational budget is exhausted) every time an action is required from the system. This approach manages to combine the best of both DPS and LT approaches. From the point of view of DPS, OLT policies provide a more generic way of representing policies that largely avoids the trial and work required to find a suitable functional form for a parameterized policy. From the point of view of LT, DPS is the key to significantly reduce the online computational requirement, through a principled offline procedure for learning the heuristic driving the construction of the LT.

We have shown through an experimental study that our approach has several desirable qualities: it (i) produces high-performance and robust policies that often outperform both pure LT and pure DPS policies, (ii) has lower offline complexity than pure DPS and lower online complexity than pure DP, and (iii) can, in particular, be optimized for a fixed online budget to meet the requirements of a specific application.

Table A.1. Physical parameters of the inverted pendulum domain.

Symbol	Value	Meaning
g	9.81 m/s ²	Gravitation
m	1 kg	Mass of link
l	1 m	Length of link
μ	0.05	Coefficient of friction

APPENDIX A: DYNAMIC MODEL OF THE INVERTED PENDULUM

Refer to the schematic representation of the inverted pendulum given in Figure 5(a). The state variables are the angle measured from the vertical axis, $\phi(t)$ (rad) $\in [0, \pi]$, and the angular velocity $\dot{\phi}(t)$ (rad/s). The control variable is the torque $u(t)$ (N m) applied, which is restricted to the interval $[-5, 5]$. The motion of the pendulum is described by the differential equation:

$$\ddot{\phi}(t) = \frac{1}{ml^2} (-\mu\dot{\phi}(t) + mgl \sin \phi(t) + u(t)). \quad (\text{A.1})$$

The angular velocity is restricted via saturation to the interval $\dot{\phi} \in [-10, 10]$. The values and meaning of the physical parameters are given in Table A.1.

The solution to the continuous-time dynamic equation is obtained by writing Equation (A.1) as a first-order system and using a Runge–Kutta solver with $h = 40$ intermediate steps. The time step of the simulation is $\Delta t = 0.2$ s, during which the applied control is kept constant. The two-dimensional state vector is $x(t) \equiv (\phi(t), \dot{\phi}(t))$, the scalar control variable is $u(t)$. As our algorithm requires a finite set of possible actions, we discretized the continuous control space into five discrete action choices $a \in \{-5, -2.5, 0, 2.5, 5\}$.

APPENDIX B: DYNAMIC MODEL OF THE DOUBLE INVERTED PENDULUM

Refer to the schematic representation of the double inverted pendulum on spring-linked carts given in Figure 5(b). The state variables are, for each cart $i = 1, 2$, the angle of displacement measured from the vertical axis, $\theta_i(t)$ (rad), the angular velocity $\dot{\theta}_i(t)$ (rad/s), the position of the cart $x_i(t)$ (m) measured from the origin (note $x_1(t) < x_2(t)$), and its velocity $\dot{x}_i(t)$ (m/s). The vector-valued control is the force $u_i(t)$ (N m) applied to each cart, which is restricted to the interval $[-2, 2]$. The system as a whole is described by the system of differential equations (see online appendix of [2]):

$$\ddot{\theta}_i(t) = [b_i^2(t)a_i^{12}(t) - a^{22}b_i^1(t)] / [a_i^{12}(t)a_i^{21}(t) - a^{11}a^{22}], \quad i = 1, 2 \quad (\text{B.1})$$

$$\ddot{x}_i(t) = [b_i^1(t) - a^{11}\ddot{\theta}_i(t)] / a_i^{12}(t), \quad i = 1, 2, \quad (\text{B.2})$$

where

$$\begin{aligned} a_i^{12}(t) &:= -\cos \theta_i(t) \\ a_i^{21}(t) &:= lm_p \cos \theta_i(t) \\ b_i^1(t) &:= g \sin \theta_i(t) - \mu_p \dot{\theta}_i(t) / (lm_p) \\ b_i^2(t) &:= lm_p \dot{\theta}_i^2(t) \sin \theta_i(t) - f_i(t) + \mu_c \text{sign}(\dot{x}_i(t)) \\ f_i(t) &:= u_i(t) + K(l_s - |x_2(t) - x_1(t)|) \\ a^{11} &:= 4l/3 \\ a^{22} &:= -(m_c + m_p). \end{aligned}$$

The angular velocity is restricted via saturation to the interval $\dot{\phi}_i \in [-10, 10]$ and the velocity of the cart to $\dot{x}_i \in [-5, 5]$. The system knows two terminal conditions, which lead to a stopping of the process: collision between one of the carts and a wall, and collision between the carts. More

Table B.1. Physical parameters of the double inverted pendulum domain.

Symbol	Value	Meaning
g	9.81 m/s ²	Gravitation
L	1.0 m	Half-length of the track
l	0.5 m	Half-length of a pole
m_c	1.0 kg	Mass of a cart
m_p	0.1 kg	Mass of a pole
μ_c	0.0005	Coefficient of friction of a cart
μ_p	0.000002	Coefficient of friction of a pole
K	2.0	Coefficient K of the spring
l_s	0.5 m	Relaxed length of the spring
l_{smin}	0.1 m	Minimum length of the spring before deformation
l_{smax}	1.5 m	Maximum length of the spring before deformation

specifically, these conditions are implemented as follows: the temporal evolution of the system is halted if at any time t at least one of the following is true:

- $|x_i(t)| > L$ (first cart collides with left wall or second cart collides with right wall);
- $x_2(t) \leq x_1(t)$ (first cart has passed the second cart);
- $|x_2(t) - x_1(t)| \notin [l_{smin}, l_{smax}]$ (outside minimal and maximal length of spring before deformation).

The values and meaning of the physical parameters are given in Table B.1.

The solution to the continuous-time dynamic equations is obtained by writing Equations (B.1) and (B.2) as a first-order system and using a Runge–Kutta solver with $h = 10$ intermediate steps. The time step of the simulation is $\Delta t = 0.1$ s, during which the applied control is kept constant. The eight-dimensional state vector is $x(t) \equiv (x_1(t), x_2(t), \dot{x}_1(t), \dot{x}_2(t), \theta_1(t), \theta_2(t), \dot{\theta}_1(t), \dot{\theta}_2(t))$, and the control vector is $u(t) \equiv (u_1(t), u_2(t))$. As our algorithm requires a finite set of possible actions, we discretized the control space into four discrete action choices $a \in \{(-2, -2), (-2, +2), (+2, -2), (+2, +2)\}$.

APPENDIX C: DYNAMIC MODEL OF THE ACROBOT

Refer to the schematic representation of the acrobot domain given in Figure 5(c). The state variables are the angle of the first link measured from the horizontal axis, $\theta_1(t)$ (rad), the angular velocity $\dot{\theta}_1(t)$ (rad/s), the angle between the second link and the first link $\theta_2(t)$ (rad), and its angular velocity $\dot{\theta}_2(t)$ (rad/s). The control variable is the torque $\tau(t)$ (N m) applied at the second joint. The dynamic model of the acrobot system is [14]

$$\ddot{\theta}_1(t) = -\frac{1}{d_1(t)} \left(d_2(t)\ddot{\theta}_2(t) + \phi_1(t) \right) \quad (C.1)$$

$$\ddot{\theta}_2(t) = \frac{1}{m_2 l_{c2}^2 + I_2 - \frac{d_2(t)^2}{d_1(t)}} \left(\tau(t) + \frac{d_2(t)}{d_1(t)} \phi_1(t) - m_2 l_1 l_{c2} \dot{\theta}_1(t)^2 \sin \theta_2(t) - \phi_2(t) \right), \quad (C.2)$$

Table C.1. Physical parameters of the acrobot domain.

Symbol	Value	Meaning
g	9.8 m/s ²	Gravitation
m_i	1 kg	Mass of link i
l_i	1 m	Length of link i
l_{ci}	0.5 m	Length to center of mass of link i
I_i	1 kg m ²	Moment of inertia of link i

where

$$\begin{aligned}
 d_1(t) &:= m_1 l_{c1}^2 + m_2 (l_1^2 + l_{c2}^2 + 2l_1 l_{c2} \cos \theta_2(t)) + I_1 + I_2 \\
 d_2(t) &:= m_2 (l_{c2}^2 + l_1 l_{c2} \cos \theta_2(t)) + I_2 \\
 \phi_1(t) &:= -m_2 l_1 l_{c2} \dot{\theta}_2(t)^2 \sin \theta_2(t) - 2m_2 l_1 l_{c2} \dot{\theta}_2(t) \dot{\theta}_1(t) \sin \theta_2(t) \\
 &\quad + (m_1 l_{c1} + m_2 l_1) g \cos \theta_1(t) + \phi_2(t) \\
 \phi_2(t) &:= m_2 l_{c2} g \cos(\theta_1(t) + \theta_2(t)).
 \end{aligned}$$

The angular velocities are restricted via saturation to the interval $\theta_1 \in [-4\pi, 4\pi]$, and $\theta_2 \in [-9\pi, 9\pi]$. The values and meaning of the physical parameters are given in Table C.1; we used the same parameters as in [6].

The solution to the continuous-time dynamic equations in Equations (C.1) and (C.2) is obtained using a Runge–Kutta solver with $h = 20$ intermediate steps. The time step of the simulation is $\Delta t = 0.2$ s, during which the applied control is kept constant. The four-dimensional state vector is $x(t) \equiv (\theta_1(t), \theta_2(t), \dot{\theta}_1(t), \dot{\theta}_2(t))$, and the scalar control variable is $\tau(t)$.

The motor was allowed to produce torques τ in the range $[-1, 1]$. As our algorithm requires a finite set of possible actions, we discretized the continuous control space. Here we use three actions: the first two correspond to a bang-bang control and take on the extreme values -1 and $+1$. However, a bang-bang control alone does not allow us to keep the acrobot in the inverted handstand position, which is an unstable equilibrium. As a third action, we therefore introduce a more complex balance action, which is derived via LQR. First, we linearize the acrobot’s equation of motion about the unstable equilibrium $(-\pi/2, 0, 0, 0)$, yielding

$$\dot{x}(t) = Ax(t) + Bu(t),$$

where, after plugging in the physical parameters of Table C.1,

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 6.21 & -0.95 & 0 & 0 \\ -4.78 & 5.25 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 \\ 0 \\ -0.68 \\ 1.75 \end{bmatrix}, \quad x(t) = \begin{bmatrix} \theta_1(t) - \pi/2 \\ \theta_2(t) \\ \dot{\theta}_1(t) \\ \dot{\theta}_2(t) \end{bmatrix} \quad u(t) = \tau(t).$$

Using MATLAB, an LQR controller was then computed for the cost matrices $Q = I_{4 \times 4}$ and $R = 1$, yielding the state feedback law

$$u(t) = -Kx(t), \tag{C.3}$$

with constant gain matrix $K = [-189.28, -47.46, -89.38, -29.19]$. The values resulting from Equation (C.3) were truncated to stay inside the valid range $[-1, 1]$. Note that the LQR controller works as intended and produces meaningful results only when the state is already in a close neighborhood of the handstand state; in particular, it is incapable of swinging up and balancing the acrobot on its own from the initial state $(0, 0, 0, 0)$.

APPENDIX D: DYNAMIC MODEL OF THE HIV DRUG TREATMENT DOMAIN

The HIV infection dynamics are described by a six-dimensional nonlinear system with the state vector $x(t) \equiv (T_1(t), T_2(t), T_1^*(t), T_2^*(t), V(t), E(t))$, where

1. $T_1(t) \geq 0$ ($T_1^*(t) \geq 0$) is the number of non-infected (infected) CD4⁺ T-lymphocytes (in cells/ml);
2. $T_2(t) \geq 0$ ($T_2^*(t) \geq 0$) is the number of non-infected (infected) macrophages (in cells/ml);
3. $V(t) \geq 0$ is the number of free HI virus particles (in copies/ml); and
4. $E(t) \geq 0$ is the number of cytotoxic T-lymphocytes (in cells/ml).

The dynamics is described by the following system of first-order differential equations (see [15]):

$$\dot{T}_1(t) = \lambda_1 - d_1 T_1(t) - (1 - \varepsilon_1(t)) k_1 V(t) T_1(t) \tag{D.1}$$

$$\dot{T}_2(t) = \lambda_2 - d_2 T_2(t) - (1 - f \varepsilon_1(t)) k_2 V(t) T_2(t) \quad (\text{D.2})$$

$$\dot{T}_1^*(t) = (1 - \varepsilon_1(t)) k_1 V(t) T_1(t) - \delta T_1^*(t) - m_1 E(t) T_1^*(t) \quad (\text{D.3})$$

$$\dot{T}_2^*(t) = (1 - f \varepsilon_1(t)) k_2 V(t) T_2(t) - \delta T_2^*(t) - m_2 E(t) T_2^*(t) \quad (\text{D.4})$$

$$\begin{aligned} \dot{V}(t) = & (1 - \varepsilon_2(t)) N_T \delta (T_1^*(t) + T_2^*(t)) - c V(t) - \\ & [(1 - \varepsilon_1(t)) \rho_1 k_1 T_1(t) + (1 - f \varepsilon_1(t)) \rho_2 k_2 T_2(t)] V(t) \end{aligned} \quad (\text{D.5})$$

$$\dot{E}(t) = \lambda_E + \frac{b_E (T_1^*(t) + T_2^*(t))}{T_1^*(t) + T_2^*(t) + K_b} E(t) - \frac{d_E (T_1^*(t) + T_2^*(t))}{T_1^*(t) + T_2^*(t) + K_d} E(t) - \delta_E E(t). \quad (\text{D.6})$$

The vector-valued control variable is $u(t) = (\varepsilon_1(t), \varepsilon_2(t))$, where ε_1 and ε_2 corresponds to the dosage of the reverse transcriptase inhibitor (RTI) and the protease inhibitor (PI) drug, respectively. In STI, drugs are either fully administered (they are ‘on’) or not at all (they are ‘off’). A fully administered RTI drug corresponds to the value $\varepsilon_1 = 0.7$, whereas a fully administered PI drug corresponds to the value $\varepsilon_2 = 0.3$. This leads to a discrete actions space with four possible choices $a \in \{(0.7, 0.3), (0.7, 0), (0, 0.3), (0, 0)\}$. Because it is not clinically feasible to change the treatment daily, the state is measured and the drugs are switched on or off once every 5 days. Therefore, the system is controlled in discrete time with a sampling period of $\Delta t = 5$ days (during which the chosen controls are kept constant).

As shown in [15], in the absence of treatment (i.e., $\varepsilon_1 = \varepsilon_2 \equiv 0$), the system in Equations (D.1)–(D.5) exhibits three physical equilibrium points:

1. an unstable equilibrium point $(T_1, T_2, T_1^*, T_2^*, V, E) = (10^6, 3198, 0, 0, 0, 10)$, which represents an uninfected state;
2. a ‘healthy’ locally stable equilibrium point $(T_1, T_2, T_1^*, T_2^*, V, E) = (967839, 621, 76, 6, 415, 353108)$, which corresponds to a small viral load, a high CD4^+ T-lymphocytes count and a high HIV-specific cytotoxic T-cells count;
3. a ‘non-healthy’ locally stable equilibrium point $(T_1, T_2, T_1^*, T_2^*, V, E) = (163573, 5, 11, 945, 46, 63, 919, 24)$ for which T-cells are depleted and the viral load is very high.

Numerical simulations show that the basin of attraction of the healthy steady state is relatively small in comparison with the one of the non-healthy steady state. Furthermore, perturbation of the uninfected steady state by adding as little as one single particle of virus per milliliter of blood plasma leads to asymptotical convergence towards the non-healthy steady state.

The solution to the continuous-time dynamic equations in Equations (D.1)–(D.5) is obtained by using a Runge–Kutta solver with $h = 500$ intermediate steps. The values and meaning of the constants in the model are the same as in [15, 16]: $\lambda_1 = 10,000$, $d_1 = 0.01$, $k_1 = 8 \cdot 10^{-7}$, $\lambda_2 = 31.98$, $d_2 = 0.01$, $f = 0.34$, $k_2 = 1 \cdot 10^{-4}$, $\delta = 0.7$, $m_1 = 1 \cdot 10^{-5}$, $m_2 = 1 \cdot 10^{-5}$, $N_T = 100$, $c = 13$, $\varrho_1 = 1$, $\varrho_2 = 1$, $\lambda_E = 1$, $b_E = 0.3$, $K_b = 100$, $d_E = 0.25$, $K_d = 500$, $\delta_E = 0.1$.

ACKNOWLEDGEMENTS

The authors first wish to thank the anonymous reviewers for their hard but fair criticisms, which greatly helped to improve the positioning and tone of the paper. Tobias Jung acknowledges financial support from a research fellowship of ULg. This paper presents research results of the PASCAL2 European FP7 Network of Excellence, and of the BIOMAGNET Network funded by the Interuniversity Attraction Poles Programme initiated by the Belgian State, Science Policy Office.

REFERENCES

1. Busoniu L, Babuska R, De Schutter B, Ernst D. *Reinforcement Learning and Dynamic Programming Using Function Approximators*. Taylor & Francis CRC Press, 2010. ISBN: 978-1439821084.
2. Hren J-F, Munos R. Optimistic planning of deterministic systems. *Proceedings of European Workshop on Reinforcement Learning (EWRL 2008)*, Villeneuve d'Ascq, France, June 30–July 3, 2008; 151–164.
3. Maes F, Wehenkel L, Ernst D. Optimized look-ahead tree policies. *Proceedings of the 9th European Workshop on Reinforcement Learning (EWRL 2011)*, Athens, Greece, September 9–11, 2011; 5–17.
4. Frenn M, Boyle P. Using Gaussian processes to optimize expensive functions. In *AI 2008: Advances in Artificial Intelligence, vol 5360 LNCS*, Wobcke W, Zhang M (eds). Springer: Berlin, 2008; 258–267.
5. Lizotte D, Wang T, Bowling M, Schuurmans D. Automatic gait optimization with Gaussian process regression. *Proceedings of IJCAI 2007*, Hyderabad, India, January 6–12, 2007; 944–949.
6. Sutton R, Barto A. *Reinforcement Learning: An Introduction*. MIT Press: Cambridge, MA, 1998.
7. Rubinstein RY, Kroese DP. *The Cross Entropy Method: A Unified Approach to Combinatorial Optimization, Monte-Carlo Simulation, and Machine Learning*. Springer: Berlin, 2004.
8. Hansen N. The CMA evolution strategy: a comparing review. In *Towards a New Evolutionary Computation. Advances on Estimation of Distribution Algorithms*, Lozano JA, Larranaga P, Inza I, Bengoetxea E (eds). Springer: Berlin, 2006; 75–102.
9. Perttunen CD, Jones DR, Stuckman BE. Lipschitzian optimization without the Lipschitz constant. *Journal of Optimization Theory and Application* 1993; **79**(1):157–181.
10. Brochu E, Cora VM, de Freitas N. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *CoRR* 2010; **abs/1012.2599**.
11. Rasmussen CE, Williams CKI. *Gaussian Processes for Machine Learning*. MIT Press: Cambridge, MA, 2006.
12. Jones DR. A taxonomy of global optimization methods based on response surfaces. *Journal of Global Optimization* 2001; **21**:345–383.
13. Mockus J. *Bayesian Approach to Global Optimization*. Kluwer Academic Publishers: Dordrecht, Holland, 1989.
14. Spong M. The swing up control problem for the acrobat. *IEEE Control Systems Magazine* 1995; **15**:49–55.
15. Adams B, Banks H, Kwon H-D, Tran H. Dynamic multidrug therapies for HIV: optimal and STI control approaches. *Mathematical Biosciences and Engineering* 2004; **1**:223–241.
16. Ernst D, Stan G, Goncalves J, Wehenkel L. Clinical data based optimal STI strategies for HIV: a reinforcement learning approach. *Proceedings of the 45th IEEE Conference on Decision and Control CDC 2006*, San Diego, USA, December 13–15, 2006; 667–672.
17. Busoniu L, Ernst D, Babuska R, De Schutter B. Cross-entropy optimization of control policies with adaptive basis functions. *IEEE Transactions on Systems, Man, and Cybernetics-Part B: Cybernetics* 2011; **41**(1):196–209.
18. Osborne M, Garnett R, Roberts SJ. Gaussian processes for global optimization. *Proceedings of 3rd International Conference on Learning and Intelligent Optimization (LION 3)*, Trento, Italy, January 14–18, 2009; 1–15.
19. Jung T, Stone P. Gaussian processes for sample efficient reinforcement learning with RMAX-like exploration. *Proceedings of ECML 2010*, Barcelona, Spain, September 20–24, 2010; 601–616.
20. Busoniu L, Munos R, De Schutter B, Babuska R. Optimistic planning for sparsely stochastic systems. *Proceedings of IEEE International Symposium on Adaptive Dynamic Programming and Reinforcement Learning ADPRL 2011*, Paris, France, April 11–15, 2011; 48–55.
21. Bertsekas D. *Dynamic Programming and Optimal Control, Vol. II*. Athena Scientific, 2007.
22. Powell W. *Approximate Dynamic Programming*. Wiley: New York, NY, 2007.
23. Peters J, Schaal S. Natural actor-critic. *Neurocomputing* 2008; **71**(7–9):1180–1190.
24. Deisenroth MP, Rasmussen CE. Pilco: a model-based and data-efficient approach to policy search. *Proceedings of the International Conference on Machine Learning ICML 2011*, Montreal, Canada, 2011; 465–472.
25. Szita I, Lőrincz A. Learning Tetris using the noisy cross-entropy method. *Neural Computation* 2006; **18**(12):2936–2941.
26. Moriarty DE, Schultz AC, Grefenstette JJ. Evolutionary algorithms for reinforcement learning. *Journal of Artificial Intelligence Research* 1999; **11**:241–276.
27. Gomez F, Mikkulainen R. Active guidance for a finless rocket using neuroevolution. *Proceedings of the Genetic and Evolutionary Computation Conference GECCO 2003*, Chicago, July 9–11, 2003; 2084–2095.
28. Stanley KO, Mikkulainen R. Evolving neural networks through augmenting topologies. *Evolutionary Computation* 2002; **10**:99–127.
29. Gomez F, Schmidhuber J, Mikkulainen R. Accelerated neuroevolution through cooperatively coevolved synapses. *Journal of Machine Learning Research* 2008; **9**:937–965.
30. Whiteson S, Taylor ME, Stone P. Critical factors in the empirical performance of temporal difference and evolutionary methods for reinforcement learning. *Autonomous Agents and Multi-Agent Systems* 2010; **21**(1):1–35.
31. Kalyanakrishnan S, Stone P. Characterizing reinforcement learning methods through parameterized learning problems. *Machine Learning* 2011; **82**(3):205–247.
32. Heidrich-Meisner V, Igel C. Variable metric reinforcement learning methods applied to the noisy mountain car problem. *Proceedings of European Workshop on Reinforcement Learning EWRL 2008*, Villeneuve d'Ascq, France, June 30–July 3, 2008; 136–150.

33. Kohl N, Stone P. Machine learning for fast quadrupedal locomotion. *Proceedings of the 19th National Conference on Artificial Intelligence AAAI-04*, San Jose, California, July 25–29, 2004; 611–616.
34. Szita I, Lőrincz A. Learning to play using low-complexity rule-based policies: illustrations through Ms. Pac-Man. *Journal of Artificial Intelligence Research* 2007; **30**:659–684.
35. Hart P, Nilsson N, Raphael B. A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics* 1968; **4**(2):100–107.
36. Korf RE. Real-time heuristic search. *Artificial Intelligence* 1990; **42**(2–3):189–211.
37. Maes F. Learning in Markov decision processes for structured prediction. *Ph.D. Thesis*, Pierre and Marie Curie University, Computer Science Laboratory of Paris 6 (LIP6), October 2009.
38. Minton S. *Machine Learning Methods for Planning*. Morgan Kaufmann Publishers Inc.: San Francisco, CA, USA, 1994.
39. Yoon SW, Fern A, Givan R. Learning heuristic functions from relaxed plans. *International Conference on Automated Planning and Scheduling ICAPS 06*, The English Lake District, Cumbria, UK, June 6–10, 2006; 162–171.
40. Arfaee SJ, Zilles S, Holte RC. Learning heuristic functions for large state spaces. *Artificial Intelligence* 2011; **175**:2075–2098.
41. Morari M, Lee JH. Model predictive control: past, present and future. *Computers & Chemical Engineering* 1999; **23**(4):667–682.
42. Coen T, Anthonis J, De Baerdemaeker J. Cruise control using model predictive control with constraints. *Computers and Electronics in Agriculture* 2008; **63**(2):227–236.
43. Ernst D, Glavic M, Capitanescu F, Wehenkel L. Reinforcement learning versus model predictive control: a comparison on a power system problem. *IEEE Transactions on Systems, Man, and Cybernetics—Part B: Cybernetics* 2009; **39**(2):517–529.
44. Busoniu L, Ernst D, Babuska R, De Schutter B. Fuzzy partition optimization for approximate fuzzy q-iteration. *Proceedings of the 17th IFAC World Congress IFAC-08*, Seoul, Korea, 2008.
45. Defourny B, Ernst D, Wehenkel L. Multistage stochastic programming: a scenario tree based approach to planning under uncertainty. In *Decision Theory Models for Applications in Artificial Intelligence: Concepts and Solutions*, Morales EF, Sucar LE, Hoes J (eds), 2012. ISBN: 978-1609601652.
46. Couetoux A, Dohmen H, Teytaud O. Improving the exploration in upper confidence trees. *Proceedings of the 6th International Conference on Learning and Intelligent Optimization LION 6*, Paris, France, January 16; 366–371.
47. Maes F, Wehenkel L, Ernst D. Learning to play K-armed bandit problems. *Proceedings of International Conference on Agents and Artificial Intelligence*, Vilamoura, Algarve, Portugal, February 2012.
48. Maes F, Wehenkel L, Ernst D. Automatic discovery of ranking formulas for playing with multi-armed bandits. *Proceedings of the 9th European Workshop on Reinforcement Learning (EWRL 2011)*, 2011; 5–17.
49. Castronovo M, Maes F, Fonteneau R, Ernst D. Learning exploration/exploitation strategies for single trajectory reinforcement learning. *Journal of Machine Learning Research Workshop & Conference Proceedings 24(JMLR W&CP) (Proceedings of the EWRL 2012)*, Edinburgh, Scotland, June 30; 1–9.
50. Kocsis L, Szepesvári C. Bandit based Monte Carlo planning. *Proceedings of the 17th European Conference on Machine Learning ECML 2006*, Berlin, Germany, September 18–22, 2006; 282–293.
51. Coulom R. Efficient selectivity and backup operators in Monte-Carlo tree search. *Proceedings of the 5th International Conference on Computers and Games*, Turin, Italy, 2006; 72–83.
52. Sokolovska N, Teytaud O, Milone M. Q-learning with double progressive widening: application to robotics. *Proceedings of ICONIP*, Shanghai, China, 2011; 103–112.
53. Chaslot G, Winands M, Uiterwijk J, van den Herik H, Bouzy B. Progressive strategies for Monte-Carlo tree search. *New Mathematics and Natural Computation* 2008; **4**(3):343–357.
54. Rolet P, Sebag M, Teytaud O. Boosting active learning to optimality: a tractable Monte-Carlo, billiard-based algorithm. *Proceedings of ECML*, Bled, Slovenia, 2009; 302–317.
55. Xu X, Liu C, Hu D. Continuous-action reinforcement learning with fast policy search and adaptive basis function selection. *Soft Computing* 2011; **15**(6):1055–1070.
56. Lazaric A, Restelli M, Bonarini A. Reinforcement learning in continuous action spaces through sequential Monte Carlo methods. *Advances in Neural Information Processing Systems NIPS 2007*, Vancouver, Canada, 2007.
57. van Hasselt H, Wiering MA. Reinforcement learning in continuous action spaces. *Proceedings of IEEE Symposium on Approximate Dynamic Programming and Reinforcement Learning ADPRL 2007*, Honolulu, HI, USA, 2008; 272–279.