# Adding Neural Network Controllers to Behavior Trees without Destroying Performance Guarantees

Christopher Iliffe Sprague, Petter Ögren

*Abstract*— In this paper, we show how Behavior Trees that have performance guarantees, in terms of safety and goal convergence, can be extended with components that were designed using machine learning, without destroying those performance guarantees.

Machine learning approaches such as reinforcement learning or learning from demonstration can be very appealing to AI designers that want efficient and realistic behaviors in their agents. However, those algorithms seldom provide guarantees for solving the given task in all different situations while keeping the agent safe. Instead, such guarantees are often easier to find for manually designed model based approaches. In this paper we exploit the modularity of Behavior trees to extend a given design with an efficient, but possibly unreliable, machine learning component in a way that preserves the guarantees. The approach is illustrated with an inverted pendulum example.

*Index Terms*— Neural Networks, Machine Learning, Performance guarantees, Safety, Behavior Trees

## I. INTRODUCTION

Over the last decade, Behavior Trees (BTs) have become a very important tool in the design of game AI, and are widely appreciated for their modularity and reactivity, [1]–[6]. At the same time machine learning approaches have continued to show remarkable success in game AI applications, [7]–[9]. However, one problem with learning approaches is that they seldom provide guarantees for ending up at the desired state, or for avoiding some unsafe states that might harm the agent. The absence of guarantees is often not a problem, but sometimes the narrative of a game requires a given result for moving on to the next scene, or the experienced realism of a game might be spoiled by a game character failing to accomplish a seemingly straightforward task. In this paper, we will show how, and when, the BT design of Figure 1c can combine the safety guarantees of an emergency subtree with the performance guarantees of a model based subtree and the efficiency of a machine learning subtree.

The basic idea is very straightforward and relies on the modularity provided by the BT structure. Note that the rest of the BT can be arbitrarily complex, but we focus on what is going on when the given subtree of interest is executed. Looking at Figure 1c, the first priority is safety, and whenever the safety constraint is in risk of being violated we invoke the Emergency subtree, this might e.g., correspond to moving away from the edge of a cliff. If safety is ok, the BT checks if the current execution time is ok, that is, if there is reason to believe that the learning subtree is not going to complete
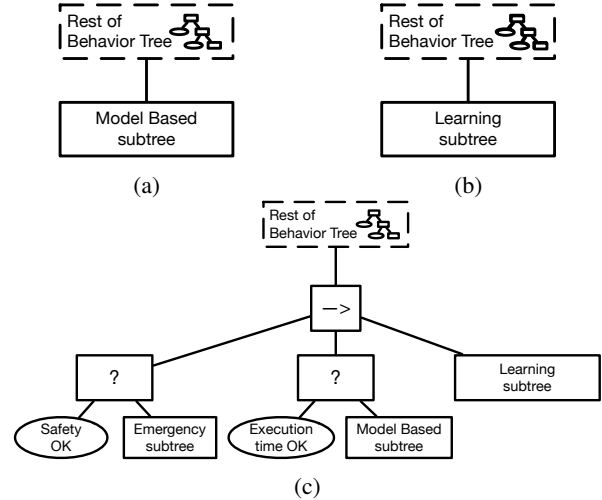
Fig. 1: Given a BT in (a) with a model based subtree, you would get the BT in (b) by replacing the model based subtree with the learning subtree. In order to guarantee safety and goal convergence, one can instead use the combined subtree in (c). If Safety is not ok, the emergency subtree will be executed. Else, if the execution time is not ok, the model-based subtree will be executed. If neither of those problems occur, the learning subtree will be executed.

the task in time, or at all. If the execution time is not ok, the previously designed model based subtree is invoked. Finally, if both of the conditions above are satisfied, we allow the learning subtree to be executed.

Note that switching between subtrees like this might induce undesired behaviors where one subtree counteracts another one, therefore, the rest of this paper is devoted to finding explicit formal conditions for when the approach outlined above will indeed provide the desired guarantees, building upon the theoretical tools proposed in [10], and illustrating the approach with a commonly known example in terms of an inverted pendulum swingup problem.

It is well known that learning algorithms might cause unsafe behavior. Both during training, and possibly even after training, as it can be hard to guarantee performance in all cases. Therefore, safety of learning approaches is a very active research area, [11]–[15].

In [13] Constrained Markov Decision Problems (CMDPs) were used, and the cumulative cost was replaced by a stepwise one, which was then transferred into the admissible control set leading back to a standard MDP formulation.

There is also a set of approaches using Lyapunov ideas,

originating in control theory. In [14] a Lyapunov approach was used to guarantee stability of an RL agent. The agent was allowed to switch between a given set of controllers that were designed to be safe no matter what switching strategy was used. Then, in [15], Lyapunov concepts were used to iteratively estimate the region of attraction of the policy, i.e., the region that the state is guaranteed not to leave, when applying the controller at hand. At the same time, while being in this safe region, the estimate of the region, as well as performance, was improved. Finally, Chow et al. used the CMDPs to construct Lyapunov functions using linear programming, [12]. The approach is guaranteed to provide feasible, and under certain assumptions, optimal policies.

Our approach differs from [11]–[15] by not trying to build the performance guarantees into the learning controller, but leveraging the reactivity of the surrounding BT structure and the existing model based controller to create a combination with the required guarantees.

BTs were invented in the gaming industry [1] and are currently spreading throughout the fields of robotics and AI [6]. Significant effort to combine BTs with learning from experience as well as demonstrations can be found in the literature [3]–[5], [16]–[21].

In [3], a complete sub-tree is learned using a genetic algorithm applied to the Mario AI environment. Similar ideas were explored in [16]. Furthermore, grammatical evolution was used in [4], to optimize the structure of a BT playing a platform game, while constraining the design to an and-or tree structure deemed efficient for the problem at hand.

Classical reinforcement learning was applied to BTs in [17], where the idea of replacing a given action (leaf) node with an RL policy was proposed. Replacing non-leaf nodes with an RL policy deciding which child to execute was explored in [18], and [19].

In [20] the BT for performing pick and place operations were learned from human demonstration, using logic and decision trees. A related idea was used in [5], where a game designer first controlled game characters to create a database of trajectories that are then used to learn controllers. Finally, in [21], a framework for end user instruction of a robot assistant was proposed.

Our approach differs from [3]–[5], [16]–[21] by not focusing on how to integrate learning into a BT, but instead providing safety and performance guarantees when such learning has been integrated. Thus, the proposed approach can be combined with any of the methods described in [3]–[5], [16]–[21].

The main contribution of this paper is thus that we show how the surrounding BT structure can be used to combine the efficiency of a possibly unreliable learning controller with the convergence of a model-based controller, and the safety of a hand crafted controller. Thereby any learning controller can be used, not restricting the choices to the elaborate approaches described above.

The outline of the paper is as follows. In Section II we present background material. Then, in Section III we formulate the main result of the paper, showing when performance guarantees can be made. A detailed inverted pendulum example is presented in Section IV and the performance of the solution for that example is explored in Section V. Finally, conclusions are drawn in Section VI.

## II. BACKGROUND

In this section we will first provide a background on classical BTs, followed by a state space formulation and a description of the corresponding theoretical tools used in this paper.

### A. Classical Formulation of BTs

A BT [6], [22]–[27] is a directed rooted tree where the internal nodes are called *control flow nodes* and leaf nodes are called *execution nodes*. For each connected node we use the common terminology of *parent* and *child*. The root is the node without parents; all other nodes have one parent. The control flow nodes have at least one child. Graphically, the children of a node are placed below it, see Figure 1.

A BT starts its execution from the root node, that generates signals called *Ticks* with a given frequency. These signals allow the execution of a node and are propagated to one or several of the children of the ticked node. A node is executed if and only if it receives Ticks. The child immediately returns *Running* to the parent, if its execution is under way, *Success* if it has achieved its goal, or *Failure* otherwise.

In the classical formulation, there exist two main categories of control flow nodes (Sequence and Fallback) and two categories of execution nodes (Action and Condition).

The Sequence node routes the Ticks to its children from the left until it finds a child that returns either *Failure* or *Running*, then it returns *Failure* or *Running* accordingly to its own parent. It returns *Success* if and only if all its children return *Success*. Note that when a child returns *Running* or *Failure*, the Sequence node does not route the Ticks to the next child (if any). The symbol of the Sequence node is a box containing the label "→".

The Fallback node[1] routes the Ticks to its children from the left until it finds a child that returns either *Success* or *Running*, then it returns *Success* or *Running* accordingly to its own parent. It returns *Failure* if and only if all its children return *Failure*. Note that when a child returns *Running* or *Success*, the Fallback node does not route the Ticks to the next child (if any). The symbol of the Fallback node is a box containing the label "?".

When an Action node receives Ticks, it executes a command such as moving the agent. It returns *Success* if the action is successfully completed or *Failure* if the action has failed. While the action is ongoing it returns *Running*.

When a Condition node receives Ticks, it checks a proposition. It returns *Success* or *Failure* depending on if the proposition holds or not. Note that a Condition node never returns a status of *Running*.

---

[1]Fallback nodes are sometimes also called *selector* or *priority selector* nodes.

## B. Theoretical Results and Statespace Formulation of BTs

In this section we will briefly review some results from [10] that will be needed to prove the overall properties of the behavior tree. To follow the notation of [10], in this and the following section, $x$ and $x_k$ denote the complete state of the system, whereas in the rest of the paper, $s$ denotes the complete state and $x$ denotes the position of the cart.

*Definition 1 (Behavior Tree):* A BT is a three-tuple

$$\mathcal{T}_i = \{f_i, r_i, \Delta t\}, \tag{1}$$

where $i \in \mathbb{N}$ is the index of the tree, $f_i : \mathbb{R}^n \to \mathbb{R}^n$ is the right hand side of an ordinary difference equation, $\Delta t$ is a time step and $r_i : \mathbb{R}^n \to \{\mathcal{R}, \mathcal{S}, \mathcal{F}\}$ is the return status that can be equal to either *Running* ($\mathcal{R}$), *Success* ($\mathcal{S}$), or *Failure* ($\mathcal{F}$). Let the Running region ($R_i$), Success region ($S_i$) and Failure region ($F_i$) correspond to a partitioning of the state space, defined as follows:

$$R_i = \{x : r_i(x) = \mathcal{R}\} \tag{2}$$
$$S_i = \{x : r_i(x) = \mathcal{S}\} \tag{3}$$
$$F_i = \{x : r_i(x) = \mathcal{F}\}. \tag{4}$$

Finally, let $x_k = x(t_k)$ be the system state at time $t_k$, then the execution of a BT $\mathcal{T}_i$ is a standard ordinary difference equation

$$x_{k+1} = f_i(x_k), \tag{5}$$
$$t_{k+1} = t_k + \Delta t. \tag{6}$$

The return status $r_i$ will be used when combining BTs recursively, as explained below.

*Definition 2 (Sequence compositions of BTs):* Two or more BTs can be composed into a more complex BT using a Sequence operator,

$$\mathcal{T}_0 = \text{Sequence}(\mathcal{T}_1, \mathcal{T}_2).$$

Then $r_0, f_0$ are defined as follows

$$\text{If } x_k \in S_1 \tag{7}$$
$$r_0(x_k) = r_2(x_k) \tag{8}$$
$$f_0(x_k) = f_2(x_k) \tag{9}$$
$$\text{else}$$
$$r_0(x_k) = r_1(x_k) \tag{10}$$
$$f_0(x_k) = f_1(x_k). \tag{11}$$

$\mathcal{T}_1$ and $\mathcal{T}_2$ are called children of $\mathcal{T}_0$. Note that when executing the new BT, $\mathcal{T}_0$ first keeps executing its first child $\mathcal{T}_1$ as long as it returns Running or Failure. The second child is executed only when the first returns Success, and $\mathcal{T}_0$ returns Success only when all children have succeeded, hence the name Sequence.

For notational convenience, we write

$$\text{Sequence}(\mathcal{T}_1, \text{Sequence}(\mathcal{T}_2, \mathcal{T}_3)) = \text{Sequence}(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3), \tag{12}$$

and similarly for arbitrarily long compositions.

*Definition 3 (Fallback compositions of BTs):* Two or more BTs can be composed into a more complex BT using a Fallback operator,

$$\mathcal{T}_0 = \text{Fallback}(\mathcal{T}_1, \mathcal{T}_2).$$

Then $r_0, f_0$ are defined as follows

$$\text{If } x_k \in F_1 \tag{13}$$
$$r_0(x_k) = r_2(x_k) \tag{14}$$
$$f_0(x_k) = f_2(x_k) \tag{15}$$
$$\text{else}$$
$$r_0(x_k) = r_1(x_k) \tag{16}$$
$$f_0(x_k) = f_1(x_k). \tag{17}$$

Note that when executing the new BT, $\mathcal{T}_0$ first keeps executing its first child $\mathcal{T}_1$ as long as it returns Running or Success. The second child is executed only when the first returns Failure, and $\mathcal{T}_0$ returns Failure only when all children have tried, but failed, hence the name Fallback.

For notational convenience, we write

$$\text{Fallback}(\mathcal{T}_1, \text{Fallback}(\mathcal{T}_2, \mathcal{T}_3)) = \text{Fallback}(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3), \tag{18}$$

and similarly for arbitrarily long compositions.

*Definition 4 (Finite Time Successful):* A BT is Finite Time Successful (FTS) with region of attraction $R'$, if for all starting points $x(0) \in R' \subset R$, there is a time $\tau$, and a time $\tau'(x(0))$ such that $\tau'(x) \leq \tau$ for all starting points, and $x(t) \in R'$ for all $t \in [0, \tau')$ and $x(t) \in S$ for $t = \tau'$.

Note that exponential stability implies FTS, given the right choices of the sets $S, F, R$, since a BT for which $x_s$ is a globally exponentially stable equilibrium of the execution (5), and $S \supset \{x : ||x - x_s|| \leq \varepsilon\}$, $\varepsilon > 0$, $F = \emptyset$, $R = \mathbb{R}^n \setminus S$, is FTS.

*Lemma 1: (Robustness and Efficiency of Sequence Compositions)* If $\mathcal{T}_1, \mathcal{T}_2$ are FTS, with $S_1 = R'_2 \cup S_2$, then $\mathcal{T}_0 = \text{Sequence}(\mathcal{T}_1, \mathcal{T}_2)$ is FTS with $\tau_0 = \tau_1 + \tau_2$, $R'_0 = R'_1 \cup R'_2$ and $S_0 = S_1 \cap S_2 = S_2$.

*Proof:* See [10]. ∎

By *safety* we denote the ability to avoid a particular part of the statespace, which we denote the *Obstacle Region* for simplicity.

*Definition 5 (Safe):* A BT is safe, with respect to the obstacle region $O \subset \mathbb{R}^n$, and the initialization region $I \subset R$, if for all starting points $x(0) \in I$, we have that $x(t) \notin O$, for all $t \geq 0$.

In order to make statements about the safety of composite BTs, we also need the following definition.

*Definition 6 (Safeguarding):* A BT is safeguarding, with respect to the step length $d$, the obstacle region $O \subset \mathbb{R}^n$, and the initialization region $I \subset R$, if it is safe, and FTS with region of attraction $R' \supset I$ and a success region $S$, such that $I$ surrounds $S$ in the following sense:

$$\{x \in X \subset \mathbb{R}^n : \inf_{s \in S} ||x - s|| \leq d\} \subset I, \tag{19}$$

where $X$ is the reachable part of the state space $\mathbb{R}^n$.

This implies that the system, under the control of another BT with maximal state space step-length $d$, cannot leave $S$ without entering $I$, and thus avoiding $O$; see Lemma 2 below.

*Lemma 2 (Safety of Sequence Compositions):* If $\mathcal{T}_1$ is safeguarding, with respect to the obstacle $O_1$ initial region $I_1$, and margin $d$, and $\mathcal{T}_2$ is an arbitrary BT with $\max_x ||x - f_2(x)|| < d$, then the composition $\mathcal{T}_0 = \text{Sequence}(\mathcal{T}_1, \mathcal{T}_2)$ is safe with respect to $O_1$ and $I_1$.

*Proof:* See [10]. ∎

In the next section, the results above will be used and extended.

## III. MAIN RESULT

In this paper we will address the following problem:

*Problem 1:* Given a BT policy where the current design includes a reliable model based subtree $\mathcal{T}_{mb}$, as shown in Figure 1a, that one would like to replace with a learning based subtree $\mathcal{T}_{learning}$, as in Figure 1b, how can this be done while preserving performance guarantees in terms of avoidance of an unsafe region and guaranteed convergence to some desired goal region of the given subtree?

The solution we propose to Problem 1 above is illustrated in Figure 1c. Informally, the main result is that such a combination does provide the required guarantees when

- $\mathcal{T}_{safe}$ is indeed safe, that is it guarantees safety in terms of avoiding a given part $O_1$ of the state space by taking over control of the agent whenever the state comes close enough to $O_1$ and then bringing it sufficiently far away from $O_1$ before handing back control to the other controllers.
- $\mathcal{T}_{mb}$ is indeed guaranteed to bring the system to the desired part of the statespace $S_3$ in finite time. Furthermore, it does so without coming close enough to $O_1$ and thereby invoking $\mathcal{T}_{safe}$ which will stop the progression towards $S_3$.
- Finally, $\mathcal{T}_{learning}$ is allowed to control the system as long as it does not violate the safety requirement and does not take too long. Unless $\mathcal{T}_{learning}$ reaches $S_3$ within some finite time, it is no longer allowed to control the system, instead, only $\mathcal{T}_{mb}$ and $\mathcal{T}_{safe}$ are active.

Now we will state this result in two theorems that can be proved formally. The second one is more easily applied as it does not require explicit knowledge of the region of attraction of the subtrees involved. But first we make an assumption that formalise Figure 1c and some corresponding details, and will be used in both versions of the theorem.

*Assumption 1:* $\mathcal{T}_{Tot}$ is a BT on the form

$$\mathcal{T}_{Tot} = \text{Sequence}(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3),$$
$$= \text{Sequence}(\text{Fallback}(\mathcal{C}_{safe}, \mathcal{T}_{safe}),$$
$$\text{Fallback}(\mathcal{C}_{time}, \mathcal{T}_{mb}),$$
$$\mathcal{T}_{learning}),$$

Furthermore, $O_1 \subset \mathbb{R}^n$ is an obstacle region, and $I_1 \subset \mathbb{R}^n$ is an initialization region such that $\mathcal{T}_1$ is safe with respect to $I_1$ and $O_1$, and safeguarding with respect to the step length $d_1 \in \mathbb{R}$. Finally, $\mathcal{T}_2, \mathcal{T}_3$ respect the step length $d_1$, $\mathcal{T}_1, \mathcal{T}_2$ are

FTS, the success regions of $\mathcal{T}_{mb}$ and $\mathcal{T}_3$ are equal, $S_{mb} = S_3$, and $\mathcal{C}_{time}$ returns False when $t > \tau_2$ for some $\tau_2 > 0$.

*Theorem 1:* If Assumption 1 holds and $S_1 = R'_{mb} \cup S_{mb}$, then $\mathcal{T}_{Tot}$ is safe with respect to $O_1$ and FTS.

*Proof:* Safety of $\mathcal{T}_{Tot}$ follows from Lemma 2 and the observation $\mathcal{T}_{Tot} = \text{Sequence}(\mathcal{T}_1, \text{Sequence}(\mathcal{T}_2, \mathcal{T}_3))$.

Now, since $\mathcal{C}_{time}$ returns False when $t > \tau_2$, and $S_2 = S_3$, the transient behaviour of $\mathcal{T}_{Tot} = \text{Sequence}(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3)$ is the same as $\text{Sequence}(\mathcal{T}_1, \mathcal{T}_2)$. By Lemma 1 we have that $\text{Sequence}(\mathcal{T}_1, \mathcal{T}_2)$ is FTS and therefore so is $\mathcal{T}_{Tot}$ with an upper bound increased by $\tau_2$. ∎

One limitation with Theorem 1 is that the assumption $S_1 = R'_{mb} \cup S_{mb}$ is sometimes hard to both achieve and verify. This assumption is needed to make sure that the safety controller $\mathcal{T}_1$ moves the state into the region of attraction of the model based controller $\mathcal{T}_2$, but then refrains from interrupting it and thereby preventing it from reaching the goal. By requiring $S_1 = R'_{mb} \cup S_{mb}$ we make sure that the trajectory moves into $S_1$ and then stays there, where the safety controller is not active.

One way to loosen this constraint is to introduce a kind of hysteresis into $\mathcal{T}_1$, giving it two success regions $S'_1$ and $S''_1$ with $S'_1 \subset S''_1$, where $S'_1$ is used to move the state into $R'_{mb}$, after which the second, larger region $S''_1$ is used to make sure $\mathcal{T}_1$ only interferes with $\mathcal{T}_2$ when needed. Instead of $S_1 = R'_{mb} \cup S_{mb}$ we can then require $S'_1 \subset R'_{mb} \cup S_{mb} \subset S''_1$, which is easier to both achieve and verify. Below we will formalise this approach.

*Theorem 2:* Let Assumption 1 hold and $\mathcal{T}_1$ have a hysteresis with two different success regions $S'_1 \subset S''_1$ such that $S'_1 \subset R'_{mb} \cup S_{mb} \subset S''_1$. The hysteresis is such that $S_1 = (S'_1 \text{ or } S''_1)$, $F_1 = \emptyset$ and $R_1 = \{x \notin S_1\}$.

The switching is done such that

$$S_1 \leftarrow S'_1 \text{ if } x \notin S''_1$$
$$S_1 \leftarrow S''_1 \text{ if } x \in S'_1,$$

that is, when $x$ moves outside the large success region $S''_1$, it is forced into the smaller success region $S'_1$ after which the larger region is once more activated. Finally, let $\mathcal{T}_{mb}$ be such that for $x \in S''_1, x \notin R'_{mb} \cup S_{mb}$ the state will leave $S''_1$ in finite time.

Under the assumptions above, $\mathcal{T}_{Tot}$ is safe with respect to $O_1$ and FTS.

*Proof:* As in Theorem 1 above, safety of $\mathcal{T}_{Tot}$ follows from Lemma 2 and the transient behaviour of $\mathcal{T}_{Tot} = \text{Sequence}(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3)$ is the same as $\text{Sequence}(\mathcal{T}_1, \mathcal{T}_2)$.

First, assume we are in a situation where $S_1 = S''_1$, the large success region. We know that $x \in S''_1$ and $S''_1 \supset R'_{mb} \cup S_{mb}$, so if we also happen to be inside $R'_{mb} \cup S_{mb}$ we will reach $S_{mb}$ in finite time. If not, we will by assumption leave $S''_1$ in finite time, switching the safety hysteresis to $S_1 = S'_1$. This will in turn drive the state into $S'_1$ and switch back to $S_1 = S''_1$. Now $x \in S'_1 \subset R'_{mb} \cup S_{mb}$ which implies that the state will reach $S_{mb}$ in finite time returning success for all of $\mathcal{T}_{Tot}$.

The case where we start with $S_1 = S'_1$ is covered in the scenario above.

Thus we have that Sequence$(\mathscr{T}_1, \mathscr{T}_2)$ is FTS and therefore so is $\mathscr{T}_{Tot}$ with an upper bound increased by $\tau_2$. ∎

## IV. EXAMPLE: INVERTED PENDULUM

To illustrate the results above we use the well known inverted pendulum problem. The normalised state equations of motions, adapted from [28], can be written as follows

$$\dot{s} = \begin{bmatrix} \dot{x} \\ \dot{v} \\ \dot{\theta} \\ \dot{\omega} \end{bmatrix} = \begin{bmatrix} v \\ u \\ \omega \\ -u\cos(\theta) + \sin(\theta) \end{bmatrix} = f(s, u), \qquad (20)$$

where the variables are: $x$ cart position, $v$ cart velocity, $\theta$ pole angle (clockwise from upright), $\omega$ pole angular velocity, and $u \in [-\bar{u}, \bar{u}]$ control input with $\bar{u}$ its maximum magnitude.

The proposed approach is built upon the structure depicted in Figure 1 and described in Theorems 1 and 2. Before going through the details of the conditions $\mathscr{C}_{\text{safe}}, \mathscr{C}_{\text{time}}$, and the actions/subtrees $\mathscr{T}_{\text{safe}}, \mathscr{T}_{\text{mb}}, \mathscr{T}_{\text{learning}}$, we provide some background on their design.

### A. Model-based controller for inverted pendulum

For the model based controller $\mathscr{T}_{\text{mb}}$ we implement the globally asymptotically stable controller of [29]. Using the normalized pendulum dynamics of Equation (20), the control law reads:

$$u = \frac{\omega k_\omega + k_\theta (\theta - \theta_r) + k_\theta \tan^{-1}(k_v v + k_x x) + \sin(\theta)}{\cos(\theta)} \quad (21)$$

$$u = \min(\max(u, -\bar{u}), \bar{u}) \qquad (22)$$

where the reference angle is

$$\theta_{r,i+1} = \begin{cases} \theta_{r,i} - 2\pi & \text{if } -\pi \leq \theta_i - \theta_{r,i} \leq -\tilde{\theta} \\ \theta_{r,i} & \text{if } -\tilde{\theta} < \theta_i - \theta_{r,i} < \tilde{\theta} \\ \theta_{r,i} + 2\pi & \text{if } \tilde{\theta} \leq \theta_i - \theta_{r,i} \leq \pi \end{cases}, \qquad (23)$$

and $\tilde{\theta}$ is the smallest positive value that satisfies

$$k_\theta(-\theta_r + \tilde{\theta}) + k_\theta \tan^{-1}(\pi k_v u_m + \pi^2 k_x(u_m + \pi)) - \pi u_m \cos(\tilde{\theta}) + \sin(\tilde{\theta}) = 0. \quad (24)$$

With a proper choice of controller gains

$$k_\theta > (1 + k_\omega) \qquad (25)$$

$$k_\theta, k_\omega \gg k_y, k_v > 0 \qquad (26)$$

the state $\theta = \theta_r$, $\omega = x = v = 0$ can be show to be globally asymptotically stable [29]. Finally, upper bounds on the deviation from $x = 0$ in terms of $|x| \leq d_{safe} = \pi^2(\pi + \bar{u})$ are also provided. For the simulations below we used the following values of the parameters: $k_x = 0.45, k_v = 0.49, k_\theta = 10$ and $k_\omega = 8$.

### B. Creating a database of optimal trajectories

The details of how the database of optimal trajectories was created can be found in the Appendix, but examples of the results are shown in Figure 2 and Figure 3.
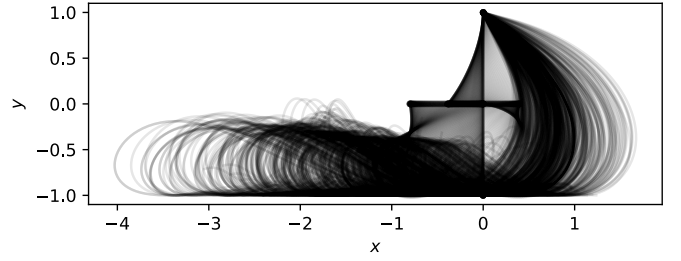


Fig. 2: A plot of $354{,}918$ state trajectories, illustrating the position of both the cart and the pole. Note that the cart moves horizontally at $y = 0$, and the tip of the pole ends at $x = 0, y = 1$ in an upright position for all trajectories.
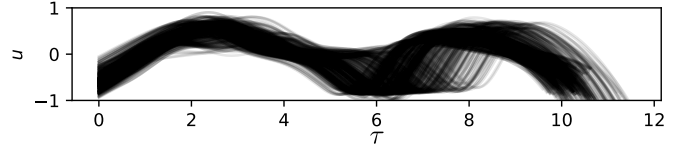


Fig. 3: The control trajectories over time for the $354{,}918$ state trajectories of Figure 2.

### C. Learning a controller from the database

For the learning controller $\mathscr{T}_{\text{learning}}$, we use a feedforward multilayer perceptron with 4 input features, corresponding to the system's state dimensionality $s \in \mathbb{R}^4$, and 1 output feature, corresponding to the control dimensionality $u \in \mathbb{R}^1$. For the hidden architecture, we use 3 hidden linear layers, each having 50 nodes. Before each hidden layer, we apply leaky rectified linear unit (LReLU), as opposed to regular ReLUs in order to avoid the "dying ReLU" problem, in which the slope of the activation function's output becomes zero, thus improving training performance. Through trial and error, we find that this architecture is sufficiently complex to model the optimal state-feedback controller.

With a learning rate of $1 \times 10^{-3}$, we train the neural network on $354{,}918$ state-control pairs for $10{,}000$ iterations, reserving 10% of the data for testing. Using the Adam optimisation algorithm [30] and the mean-squared error (*MSE*) loss function, our network achieves final training and testing losses of $2.1445 \times 10^{-3}$ and $2.2745 \times 10^{-3}$, respectively.

### D. The detailed implementations

The pseudocode for the the conditions $\mathscr{C}_{\text{safe}}, \mathscr{C}_{\text{time}}$, and the actions/subtrees $\mathscr{T}_{\text{safe}}, \mathscr{T}_{\text{mb}}, \mathscr{T}_{\text{learning}}$, can be found in Algorithms 1-5. The positive constants $\varepsilon_i > 0$ are used to define when the overall task is finished $||x, v, \theta, \omega|| \leq \varepsilon_2$ how close to the origin the safety controller should move the cart before handing over to the other controllers $||x, v|| \leq \varepsilon_1$, and how close to the origin the model based controller should always be used, due to poor stationary performance of the learning controller $||x, v, \theta, \omega|| \leq \varepsilon_3$. Note that $\varepsilon_1, \varepsilon_2$ are only there to account for the fact that a real system cannot be expected to hit exactly $x = 0$, while $\varepsilon_3$ accounts for the
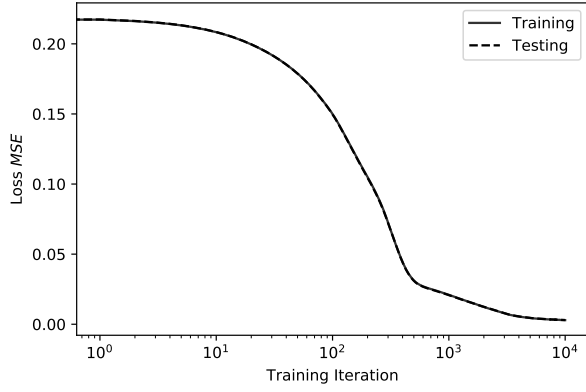
Fig. 4: The mean-squared error training and testing loss of the implemented neural network controller.

fact that the learned controller has poor performance for very small deviations from the origin, which is in line with the assumption that learning controllers are not completely reliable.

---

**Algorithm 1:** Safety OK condition: $\mathscr{C}_{\text{safe}}$

---
1 **Function** *Tick()*
2    **if** $||x|| \leq d_{safe}$ *and ResettingPosition == False* **then**
3      **return** *Success*
4    **else**
5      **return** *Failure*

---

**Algorithm 2:** Emergency Controller: $\mathscr{T}_{\text{safe}}$

---
1 **Function** *Tick()*
2    **if** $||x|| > d_{safe}$ **then**
3      $u = -\bar{u}\,\text{sgn}(x)$
4      ResettingPosition == True
5    **else**
6      $u = -k_x x - k_v v$
7      $k_x, k_v > 0$
8    **if** $||x, v|| \leq \varepsilon_1$ **then**
9      ResettingPosition == False
10      **return** *Success*
11    **return** *Running*

---

**Algorithm 3:** Execution time OK condition: $\mathscr{C}_{\text{time}}$

---
1 **Function** *Tick()*
2    **if** $t < \tau_2$ **then**
3      **return** *Success*
4    **else**
5      **return** *Failure*

---

**Algorithm 4:** Model-based controller: $\mathscr{T}_{\text{mb}}$

---
1 **Function** *Tick()*
2    Let $u$ be given by Equation (22)
3    **if** $||x, v, \theta, \omega|| \leq \varepsilon_2$ **then**
4      **return** *Success*
5    **else**
6      **return** *Running*

---

**Algorithm 5:** Learning controller: $\mathscr{T}_{\text{learning}}$

---
1 **Function** *Tick()*
2    Let $u$ be given by an ANN trained on the data of Section IV-B.
3    **if** $||x, v, \theta, \omega|| \leq \varepsilon_2$ **then**
4      **return** *Success*
5    **else if** $||x, v, \theta, \omega|| \leq \varepsilon_3$ **then**
6      Let $u$ be given by the model based controller, Equation (22)
7      **return** *Running*
8    **else**
9      **return** *Running*

---

### E. Theoretical Analysis

*Lemma 3:* The example specified in Algorithms 1-5 produces an overall BT that is FTS and Safe, with obstacle region $O_1$

*Proof:* We now need to verify the conditions stated in Theorem 2. Looking at Assumption 1 we have the following.

$\mathscr{T}_1$ is FTS since the cart is just a double integrator and the control is a standard PD-controller making the origin globally asymptotically stable. $\mathscr{T}_2$ is FTS since it is globally asymptotically stable, with upper bounds on the deviation from $x = 0$ in terms of $|x| \leq d_{safe} = \pi^2(\pi + \bar{u})$ as described in [29].

Starting from zero velocity at one end of the $2d_{safe}$ long area the maximal velocity after applying the acceleration $\bar{u}$ is $2\sqrt{d_{safe}\bar{u}}$. Thus the initialization region must be at least of width $d = 2\sqrt{d_{safe}\bar{u}}\Delta t$.

Since the cart is a double integrator with bounded acceleration, the distance needed to stop the cart at a given speed is equal to the distance needed to accelerate it to that speed. In this case that distance is equal to $2d_{safe}$. Thus we let the initialization region $I_1$ be an area just outside $|x| \leq d_{safe}$ of width $d$, and the obstacle region $O_1$ needs to be at least $2d_{safe}$ away from $I_1$ to allow a reversal of the velocity before a collision occurs. $S_{\text{mb}} = S_3$ is clear from Algorithms 4 and 5.

Looking at the additional assumptions in Theorem 2 we have that the safety hysteresis is created in Algorithms 1 and 2 with $S_1' = \{s : |x| \leq d_{safe}\}$ and $S_1'' = \{s : ||x, v|| < \varepsilon_1\}$, clearly implying $S_1' \subset S_1''$, and $S_1' \subset R_{\text{mb}}' \cup S_{\text{mb}} \subset S_1''$.

To conclude, all prerequisites for Theorem 2 are satisfied and we conclude that $\mathscr{T}_{Tot}$ is indeed FTS and safe with respect to $O_1$. ∎
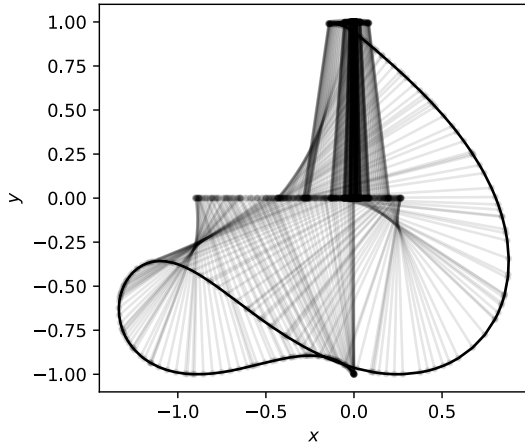
Fig. 5: Simulation of $\mathcal{T}_{\text{Learning}}$. The cart is moving horizontally at $y = 0$, starting at $x = 0$. The pole is initially hanging straight down $\theta = \pi$, with the tip at $(0, -1)$. The thin straight lines show the pole at different time steps, connecting the cart position with the tip of the pole. Initially the cart moves left, to $(-1, 0)$ and then right again after a small pause to get the pole rotating fast. Note how the ANN controller efficiently swings up the pole, and the model-based controller takes over near the goal state with small corrections, producing the thick dark region between $(0, 0)$ and $(0, 1)$.
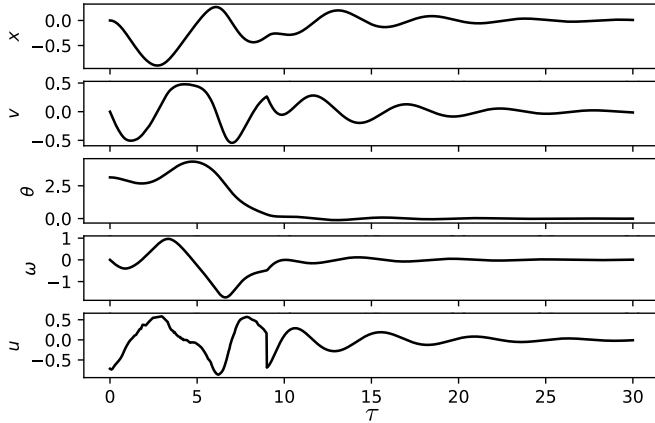


Fig. 6: The states over time for the scenario in Figure 5. Note how $\theta$ starts at $\pi$ and ends at 0.

## V. Simulations and Efficiency

First we run a few simulations to illustrate the approach and then derive a table to assess the efficiency of it.

In Figures 5 and 6, the overall behaviour of the proposed BT can be seen. In particular we see how the condition $||x, v, \theta, \omega|| \leq \varepsilon_3$ changes the behaviour close to the goal. In Figures 7 and 8, the effect of $\mathscr{C}_{\text{time}}$ and a switch to the model based subtree can be seen. The effects of $\mathscr{C}_{\text{safe}}$, and a switch to the safety subtree can be seen in Figures 9 and 10.

We asses the performance of our BT on the swing up task, and choose absolute impulse $\int_0^{\tau_f} |u| dt$ as our primary performance metric. The impulse, as well as convergence
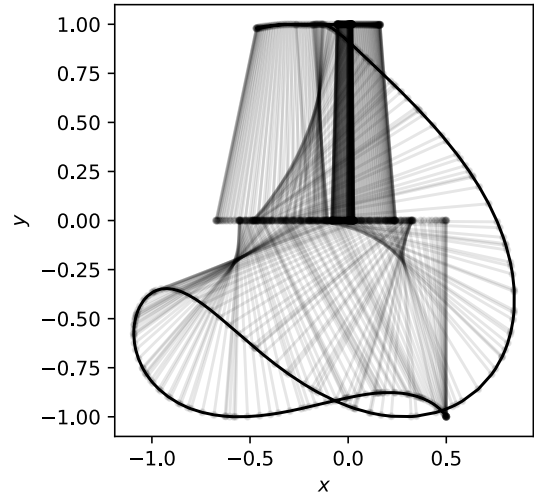


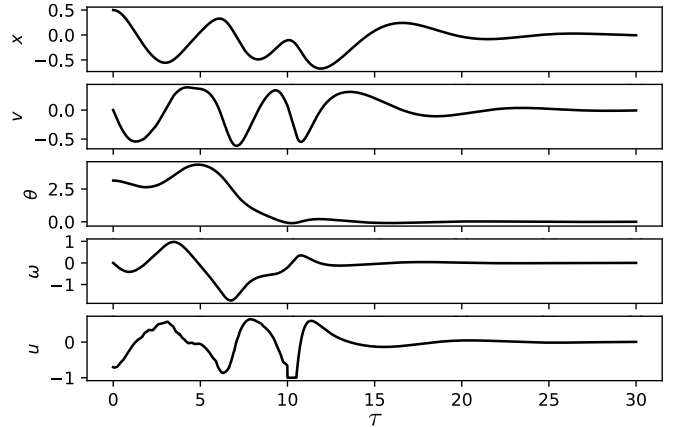Fig. 7: At $t > 10$ the model-based controller is invoked to guarantee convergence.



Fig. 8: The states over time for the scenario in Figure 7. Note the discontinuity in $u$ at $t = 10$.

time, are supplied for each controller of Figure 1 in Table I.

As expected, the learning subtree has the best performance, in fact it is nearly optimal; however it still lacks safety guarantees. The model-based subtree had the worst performance, although it is formally guaranteed to converge. The combined BT, including both the learning and model-based subtrees, not only has acceptable performance with respect to the nominal optimal trajectory, but is also guaranteed to converge and be safe.

| Controller | Convergence time $(s)$ | Impulse $(N \cdot s)$ |
|---|---|---|
| Optimal | 10.249 | 3.696 |
| Neural network | 10.253 | 3.734 |
| Model-based | 20.161 | 8.402 |
| Behaviour tree | 14.785 | 4.479 |

TABLE I: Swing up performance analysis of controllers with respect to nominal optimal control trajectory.
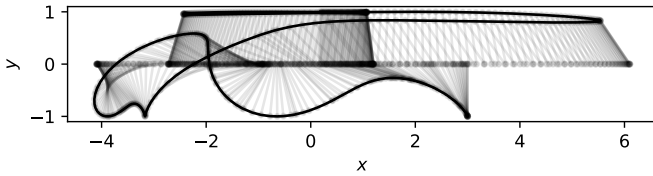
Fig. 9: At $x \leq -4$ the emergency controller is invoked to avoid getting too far from the origin.
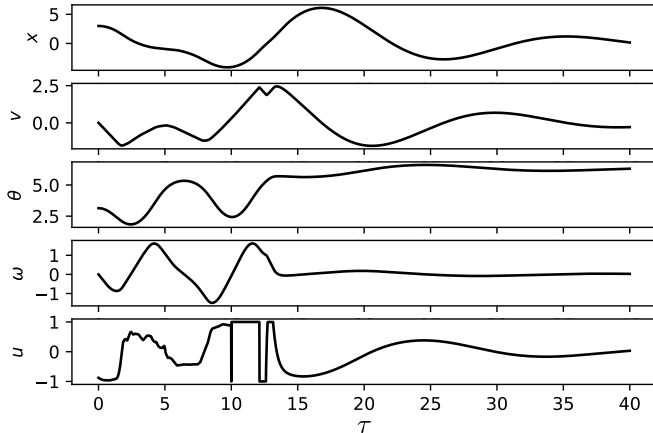


Fig. 10: The states over time for the scenario in Figure 9.

## VI. Conclusions

Previous works have separately explored different ways of adding learning components to a BT, and approaches for building safety guarantees into learning controllers.

In this paper we have shown how the reactivity and modularity of BTs enable a design where safety and convergence guarantees are provided on top of any learning controller. This was done using a natural combination of the learning controller with a safety controller and a convergent model based controller. Such a design might however introduce deadlocks where the different controllers work against each other. The paper presents a set of conditions that guarantee that the proposed design does not suffer from such problems, and instead achieves both safety and convergence to the desired goal states. An inverted pendulum example was used to illustrate the approach.

## Acknowledgement

## Appendix

The database of optimal control trajectories was created using the minimum principle of Pontryagin's [31]. We first define a trajectory cost functional

$$\mathscr{J} = \int_0^{\tau_f} u^2 d\tau \qquad (27)$$

and using the minimum principle [31] define the Hamiltonian

$$\mathscr{H} = -\lambda_\omega (u \cos(\theta) - \sin(\theta)) + \lambda_\theta \omega + \lambda_v u + \lambda_x v + u^2, \qquad (28)$$

where $\boldsymbol{\lambda}$ is the so-called costate. We minimise the Hamiltonian and retrieve the optimal control as a function of both states and costates

$$u^\star = \frac{\lambda_\omega \cos(\theta)}{2} - \frac{\lambda_v}{2}. \qquad (29)$$

Lastly, we compute the costate equations of motion

$$\dot{\lambda} = -\nabla_s \mathscr{H} = \begin{bmatrix} 0 \\ -\lambda_x \\ -\lambda_\omega (u \sin(\theta) + \cos(\theta)) \\ -\lambda_\theta \end{bmatrix}. \qquad (30)$$

Considering the task of swinging up the pendulum, we set initial and terminal state constraints

$$s(0) = \begin{bmatrix} 0 & 0 & \pi & 0 \end{bmatrix}^\mathsf{T} \qquad (31)$$

$$s(\tau_f) = \begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix}^\mathsf{T}, \qquad (32)$$

as well as the free-time condition

$$\mathscr{H}(s(\tau), \lambda(\tau), u(\tau)) = 0. \qquad (33)$$

The trajectory optimisation problem then becomes: choose the trajectory duration $\tau_f$ and initial costates $\lambda(0)$ such that the constraints are satisfied.

We first solve for a nominal trajectory, then execute *random-walks*, exploiting the homotopy of the state dynamics by resolving the problem from randomly perturbed states along the trajectory, see [32]. Through this process we are able to generate large databases of optimal control trajectories, whose state-control pairings can be used to train a universal function approximator, namely a neural network, to approximate the optimal state-feedback control policy. Examples of these databases are shown in Figure 2 and Figure 3.

## References

[1] A. Champandard, M. Dawe, and D. Cerpa, "Behavior trees: Three ways of cultivating strong ai," in *Game Developers Conference, Audio Lecture*, 2010.

[2] G. Flórez-Puga, M. A. Gomez-Martin, P. P. Gomez-Martin, B. Díaz-Agudo, and P. A. González-Calero, "Query-enabled behavior trees," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 1, no. 4, pp. 298–308, 2009.

[3] M. Colledanchise, R. Parasuraman, and P. Ögren, "Learning of Behavior Trees for Autonomous Agents," *IEEE Transactions on Games, DOI 10.1109/TG.2018.2816806*, 2018.

[4] M. Nicolau, D. Perez-Liebana, M. O'Neill, and A. Brabazon, "Evolutionary behavior tree approaches for navigating platform games," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 9, no. 3, pp. 227–238, 2016.

[5] I. Sagredo-Olivenza, P. P. Gómez-Martín, M. A. Gómez-Martín, and P. A. González-Calero, "Trained behavior trees: Programming by demonstration to support ai game designers," *IEEE Transactions on Games*, 2017.

[6] M. Colledanchise and P. Ögren, *Behavior Trees in Robotics and AI, an Introduction.* Chapman and Hall/CRC, 2018.

[7] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015, Accessed on: Apr 20, 2017. [Online]. Available: http://dx.doi.org/10.1038/nature14236

[8] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, p. 484, 2016.

[9] N. Justesen, P. Bontrager, J. Togelius, and S. Risi, "Deep learning for video game playing," *IEEE Transactions on Games*, 2019.

[10] M. Colledanchise and P. Ögren, "How Behavior Trees Modularize Hybrid Control Systems and Generalize Sequential Behavior Compositions, the Subsumption Architecture, and Decision Trees," *IEEE Transactions on Robotics*, vol. 33, no. 2, pp. 372–389, 2017.

[11] W. Saunders, G. Sastry, A. Stuhlmueller, and O. Evans, "Trial without error: Towards safe reinforcement learning via human intervention," in *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2018, pp. 2067–2069.

[12] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh, "A lyapunov-based approach to safe reinforcement learning," *arXiv preprint arXiv:1805.07708*, 2018.

[13] M. El Chamie, Y. Yu, and B. Açıkmese, "Convex synthesis of randomized policies for controlled markov chains with density safety upper bound constraints," in *American Control Conference*, 2016, pp. 6290–6295.

[14] T. J. Perkins and A. G. Barto, "Lyapunov design for safe reinforcement learning," *Journal of Machine Learning Research*, vol. 3, no. Dec, pp. 803–832, 2002.

[15] F. Berkenkamp, M. Turchetta, A. Schoellig, and A. Krause, "Safe model-based reinforcement learning with stability guarantees," in *Advances in Neural Information Processing Systems*, 2017, pp. 908–918.

[16] C. Lim, R. Baumgarten, and S. Colton, "Evolving Behaviour Trees for the Commercial Game DEFCON," *Applications of Evolutionary Computation*, pp. 100–110, 2010.

[17] R. d. P. Pereira and P. M. Engel, "A Framework for Constrained and Adaptive Behavior-based Agents," *arXiv Preprint arXiv:1506.02312*, 2015.

[18] R. Dey and C. Child, "Ql-bt: Enhancing behaviour tree design and implementation with q-learning," in *2013 IEEE Conference on Computational Intelligence in Games (CIG)*. IEEE, 2013, pp. 1–8.

[19] Y. Fu, L. Qin, and Q. Yin, "A reinforcement learning behavior tree framework for game ai," in *2016 International Conference on Eco-nomics, Social Science, Arts, Education and Management Engineering*. Atlantis Press, 2016.

[20] French, Wu, Pan, Zhou, and Jenkins, "Learning Behavior Trees From Demonstration," in *Robotics and Automation (ICRA), 2019 IEEE International Conference on*. IEEE, 2019.

[21] C. Paxton, A. Hundt, F. Jonathan, K. Guerin, and G. D. Hager, "CoSTAR: Instructing Collaborative Robots with Behavior Trees and Vision," in *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2017, pp. 564–571.

[22] D. Perez, M. Nicolau, M. O'Neill, and A. Brabazon, "Evolving Behaviour Trees for the Mario AI Competition Using Grammatical Evolution," in *Proceedings of the 2011 International Conference on Applications of Evolutionary Computation - Volume Part I*, ser. EvoApplications'11. Berlin, Heidelberg: Springer, 2011.

[23] K. Y. W. Scheper, S. Tijmons, C. C. de Visser, and G. C. H. E. de Croon, "Behaviour Trees for Evolutionary Robotics," *CoRR*, vol. abs/1411.7267, 2014.

[24] B. Hannaford, D. Hu, D. Zhang, and Y. Li, "Simulation Results on Selector Adaptation in Behavior Trees," *arXiv Preprint arXiv:1606.09219*, 2016.

[25] B. Merrill, "Ch 10, Building Utility Decisions into Your Existing Behavior Tree," *Game AI Pro. A Collected Wisdom of Game AI Professionals*, 2014.

[26] S. Ocio, "Adapting AI Behaviors To Players in Driver San Francisco: Hinted-Execution Behavior Trees," in *Eighth Artificial Intelligence and Interactive Digital Entertainment Conference*, 2012.

[27] M. Nicolau, D. Perez-Liebana, M. O'Neill, and A. Brabazon, "Evolutionary Behavior Tree Approaches for Navigating Platform Games," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. PP, no. 99, pp. 1–1, 2016.

[28] F. Mazenc and L. Praly, "Adding integrations, saturated controls, and stabilization for feedforward systems," *IEEE Transactions on Automatic Control*, vol. 41, no. 11, pp. 1559–1578, Nov 1996.

[29] B. Srinivasan, P. Huguenin, K. Guemghar, and D. Bonvin, "A global stabilization strategy for an inverted pendulum," *IFAC Proceedings Volumes*, vol. 35, no. 1, pp. 133 – 138, 2002, 15th IFAC World Congress. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S1474667015386936

[30] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[31] L. S. Pontryagin, V. G. Boltyanshii, R. V. Gamkrelidze, and E. F. Mishenko, *The Mathematical Theory of Optimal Processes*. New York: John Wiley and Sons, 1962.

[32] D. Izzo, C. I. Sprague, and D. V. Tailor, *Machine Learning and Evolutionary Techniques in Interplanetary Trajectory Design*. Cham: Springer International Publishing, 2019, pp. 191–210. [Online]. Available: https://doi.org/10.1007/978-3-030-10501-3_8