# Improving the performance of Learned Controllers in Behavior Trees using Value Function Estimates at Switching Boundaries

Mart Kartašev and Petter Ögren

*Abstract*—Behavior trees represent a modular way to create an overall controller from a set of sub-controllers solving different sub-problems. These sub-controllers can be created using various methods, such as classical model based control or reinforcement learning (RL). If each sub-controller satisfies the preconditions of the next sub-controller, the overall controller will achieve the overall goal. However, even if all sub-controllers are locally optimal in achieving the preconditions of the next, with respect to some performance metric such as completion time, the overall controller might still be far from optimal with respect to the same performance metric. In this paper we show how the performance of the overall controller can be improved if we use approximations of value functions to inform the design of a sub-controller of the needs of the next one. We also show how, under certain assumptions, this leads to a globally optimal controller when the process is executed on all sub-controllers. Finally, this result also holds when some of the sub-controllers are already given, i.e., if we are constrained to use some existing sub-controllers the overall controller will be globally optimal given this constraint.

*Index Terms*—Behavior trees, Reinforcement learning, Autonomous systems, Artificial Intelligence

## I. Introduction

**B**EHAVIOR TREES (BTs) are receiving increasing attention in robotics [1], [2] where they are used to create modular reactive controllers from a set of sub-controllers solving different sub-problems. In this paper, we show how to improve the performance of such modular designs when incorporating RL [3] by iteratively computing value functions of sub-controllers and using those value functions in the design of the sub-controllers executing before them.

BTs were originally conceived in the game AI domain [4], in an effort to make the controllers of in-game characters more modular, and were later shown to be optimally modular, in the sense of having so-called essential complexity equal to one [5].

Modularity is an important property in many engineering disciplines that enables designers to solve problems by dividing them into sub-problems, that are combined into a solution for the overall problem. In the context of robot control, this might correspond to creating sub-controllers such as *Move to*, *Grasp*, *Push* etc. These controllers will then be executed in a sequence until a primary goal is reached.

However, just reaching the goal is sometimes not enough. Instead we might want to reach it in a way that is near optimal

The authors are with the Robotics, Perception and Learning Lab., School of Electrical Engineering and Computer Science, Royal Institute of Technology (KTH), SE-100 44 Stockholm, Sweden, kartasev@kth.se

Fig. 1: The agent first goes from room 1 to room 2, and then goes to either object A or object B, see the map in (a). The value functions for going from room 1 to room 2 can be seen in (b), the value function for going to A can be seen in (c) and the value function for going to B can be seen in (d).

with respect to some metric such as completion time, energy or safety. If the overall task is associated with such a performance metric it might be that all sub-controllers are locally optimal with respect to this metric, but the overall controller is still far from globally optimal.

An example of this is shown in Figure 1. Here the overall goal is to reach either position A or position B in room 2, starting from room 1. This problem is divided into three sub-controllers, *Go to room 2* (1b), *Go to position A* (1c) and *Go to position B* (1d). The numbers in the figures show the value function in terms of the remaining time to completion (assuming one step takes one time unit). Let's assume that the controllers are to take the shortest path to their goal. If we connect *Go to room 2* with *Go to position A*, the combined controller will first leave room 1 and then reach object A. However if we start in the lower half of room 1 we will exit through the lower door, from which there are 11 steps left to reach A, if we instead would have exited through the upper door, it would have taken longer time to exit room 1, but the path to A in room 2 would only be 2 steps long, leading to an overall shorter path. In this case the overall controller would not be globally optimal, even though the sub-controllers are locally optimal.

As seen in the example above, to achieve global optimality

Fig. 2: Knowing which object to go to, we can use the value function of that action, or an approximation of it, as a boundary value for the first action. Using the boundary values from object A in Figure 1(c) we get the value function of *Goto room 2* shown in (a). Similarly, using the boundary values from object B in Figure 1(d) we get the value function of *Goto room 2* shown in (b).

each controller needs to be aware of how good different states are from the perspective of the next controller. This information can be found using the value function (estimating the expected accumulated future reward) of the next controller at the switching boundary. The results of doing this for *Go to room 2* is shown in Figure 2.

As an example where such issues arise in practice, consider a hypothetical agent navigating indoors with a "Move To" action. If completion time is penalized, a locally optimal solution will induce high speeds and accelerations. However, this might cause problems after the switching boundary, where the robot might collide with an object, or frighten a human collaborator. With our method, the training of Move to can also take the considerations of the subsequent actions into account.

The main contributions of this paper are as follows:

1) For a set of two local RL-controllers (controllers created using RL) with given execution order, we show that if we use the value function of the second controller as a reward during the final step of the first controller, the resulting value function will satisfy the Bellman equation across the switching boundary, making the combined controller globally optimal.
2) We extend this result to a larger set of local RL-controllers, under certain assumptions, making the combined controller globally optimal.
3) Given a mixed set of controllers designed using RL or with manual design principles, we extend the result above so that the global controller is optimal, under the constraint that the non-RL controllers are not changed.

Before moving on we note two things: First, the proposed approach is applicable in any context where the switching boundaries between a set of policies are fixed, and some or all of those policies are implemented using RL. This includes settings such as using RL with the sequential behavior compositions in [6] and the consecutive policies of [7]. However, BTs are one of the most common setups that specifically lead to such static switching boundaries, and the computation needed to find these boundaries for an arbitrary BT can be found in [8]. Therefore, even though the results might be applicable outside a BT setting, we choose to present and motivate our work in the context of BTs.

Second, note that in the example above, tailoring *Goto room 2* to *Goto A* makes the overall performance worse when combined with *Goto B*, and vice versa. But what if you wanted to create a policy that performs well for both A and B? This can be done in two ways. If the knowledge of what will come next can be included in the state, *Goto room 2* knows whether A or B is next, and the suggested optimization will produce a policy that handles both cases optimally. If on the other hand the outcome of what will come next is unknown or random, we can optimize over an (estimated) distribution of subsequent policies. In the example above, if there is a 30% probability of going to A and a 70% probability of going to B, a weighted average of the value functions in Figures 1c and 1d, $(0.3v_A + 0.7v_B)$ can be used to optimize *Goto room 2*, yielding a slight preference for the lower door, but not as strong as if it were known that B was always the target. Similarly, in the previous example of indoor navigation, we would learn that a high velocity on the boundary leads to a low reward with high probability as the value prediction decreases due to experienced collisions. Then, over time we would learn to slow down.

Thus, our approach can create a reusable policy, that is designed to optimize the expected reward over a given distribution of possible subsequent policies. This allows for modular use of actions within specific contexts, while improving expected performance, albeit with a possible trade-off in performance with respect to other contexts, that are not experienced in training or represented in the state. As seen in the examples above, our method yields the greatest performance improvement in scenarios where the locally optimal solution causes global performance degradation.

The remainder of this paper is organized as follows. The related work is described in Section II. A brief background is provided in Section III. Then the proposed approach including the theoretical results is given in Section IV, followed by a numerical example in Section V. Finally, conclusions are drawn in Section VI.

## II. RELATED WORK

RL has previously been used in BTs for training both the control switching mechanisms [3], [9]–[12], and the individual actions [3] that constitute a BT. Many of these works handle this in the framework of Hierarchical Reinforcement Learning (HRL). However HRL is not applicable to our problem as the hierarchical controller we consider, with it's switching boundaries, is fixed, as given by the BT.

In [9] the authors build a so-called Q-condition for the lowest level sequences of a BT, that can estimate the Q-value for each action in a sequence. Every tick of the BT evaluates the Q-conditions and reorders the actions in the sequences. The highest utility action of each sequence is used to determine the utility of the sequence, which can then be recursively used to reorder the entire tree, ensuring that the right action is executed at the right time. The concepts of *learning action node* and *learning composite node* are introduced in [3], using the Options framework [13] as the theoretical foundation. The learning action node encapsulates an independent RL problem with a complete definition of states, actions and rewards. In the

case of the composite node, the authors use branches of the BT as actions in an RL problem to create control flow nodes. In a similar fashion, [10]–[12] all use RL in slightly different ways to train fallbacks that optimize the choice between alternative options that achieve the same sub-goal. Our work is different from [10]–[12] in that we use RL to update the policies, whereas they use RL to decide when to use different policies.

As mentioned above, [3] included an approach for using RL to create individual policies in a BT using a *learning action node*, where a user defined MDP was created for each policy. However, these policies were solved separately, leading to the problem of locally optimal policies possibly being far from globally optimal, as discussed above. Our work can be seen as an extension of [3], where, under certain assumptions, we can produce globally optimal policies by letting the reward on the switching boundary be given by the value function of the subsequent policy.

The most similar related work can be found in [7], where the authors address the problem of local/global optimality in RL applied in different regions with fixed switching boundaries. However, their work is different from ours in the following aspects: (a) they provide no theoretical analysis, while we show convergence to the globally optimal policy. (b) they combine critics linearly across all states, while we use the next task's value function as reward only at the boundary. (c) they rely on a user-defined parameter ($\eta$), while our approach does not. (d) Their method cannot handle cases where preference information needs to flow across multiple switching boundaries, whereas our approach can.

Another paper that addresses a related problem is [14]. There, it is noted that switching between different controllers is sometimes problematic, as one controller might finish in a state that makes the subsequent controller fail. The proposed solution records success and failure data from experienced trajectories, and uses this to train transition policies to reach successful starting states before invoking the original controller. The approach in [14] does remove failures, but cannot optimize further than that. Our method can achieve near optimal performance, as measured by the MDP reward, whereas they focus on feasibility. We adopt the final part of a controller execution to the subsequent one, where they let the first controller finish and then execute a transition to the next skill.

## III. Background

In this section we will give a very brief description of how BTs induce a partitioning of the statespace into different operation regions, as well as listing some key concepts and results from reinforcement learning.

### A. Behaviour trees

A BT represents a way of combining a set of sub-controllers into one overall controller in a way that is hierarchical and have been shown to be optimally modular [5]. A recent survey can be found in [1], and technical overview in [8].

A given BT induces a partition of the state space into a hierarchical set of operating regions $\Omega_i$, where controller $i$

is executing when the state $s \in \Omega_i$ as seen in the following Lemma from [8].

**Lemma 1.** *For a given node $j$, with operating region $\Omega_j$, controller $j$ is executed when $s \in \Omega_j$, and the operating regions of the children of $j$ is a partitioning of $\Omega_j$.*

*Proof.* This is a concatenation of Lemmas 7 and 8 of [8]. $\square$

The computation of the operating regions is a straightforward, but somewhat complex recursive operation on the tree structure, given by Definitions 10 and 11 in [8].

An illustration of this partitioning can be found in Figure 4 below. In this paper we will investigate the interactions between controllers across the operating regions.

### B. MDPs and Reinforcement learning

Let a Markov Decision Process (MDP) be defined as follows [15, p232]:

**Definition 1.** *(Markov Decision Process). An MDP is a 4-tuple*

$$(S, A, p, r), \tag{1}$$

*where $S$ is a set of states, $A$ is a set of actions, with $A(s) \subset A$ the set of actions available at state $s$, $p(s'|s,a) = p(s_{t+1} = s'|s_t = s, a_t = a)$ is the probability of state transitions, and $r(s, s', a)$ is the reward for transitioning to state $s'$ from state $s$ applying action $a$.*

A policy, $\pi : S \rightarrow A$, assigns an action to each state. The value function, $v^\pi : S \rightarrow \mathbb{R}$, of a policy $\pi$ is the expected cumulative reward gained by the policy,

$$v^\pi(s) = \mathbb{E}\{r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \ldots | s_t = s, \pi\}, \tag{2}$$

and satisfies the Bellman equation

$$v^\pi(s) = \sum_{s'} p(s'|s, \pi(s)) \left[ r(s, s', \pi(s)) + \gamma v^\pi(s') \right]. \tag{3}$$

A policy that maximizes $v^\pi$ is called an optimal policy $\pi^*$, and the corresponding value function, the optimal value function $v^*$, is the unique solution to the Bellman optimality equation [15, p233]

$$v^*(s) = \max_a \sum_{s'} p(s'|s, a)[r(s, s', a) + \gamma v^*(s')]. \tag{4}$$

An optimal policy can be found from the optimal value function as [15, p233]

$$\pi^*(s) = \text{argmax}_a \sum_{s'} p(s'|s, a)[r(s, s', a) + \gamma v^*(s')]. \tag{5}$$

## IV. Proposed approach

In this section we will provide the main result of the paper. The intended use of the results is to let the sets $\Omega_\alpha, \Omega_\beta$ below correspond to different operating regions $\Omega_i$ of a BT. Then, Lemma 2 shows how the policy in $\Omega_\beta$ can be designed independently of $\Omega_\alpha$, and exactly how the policy in $\Omega_\alpha$ must take $\Omega_\beta$ into account to provide overall optimality. Furthermore, Lemma 4 shows how to handle the case when the policy in $\Omega_\beta$ is already given. Finally, Lemma 5 shows how to

recursively extend these results to BTs with many operating regions.

**Definition 2** (MDP-neighbors). *Given an MDP, we say that two disjoint sets $\Omega_\alpha, \Omega_\beta \subset S$ are MDP-neighbors if there are two states $s_\alpha \in \Omega_\alpha, s_\beta \in \Omega_\beta$ and an action $a \in A$ such that the transition probability $p(s_\beta | s_\alpha, a) \neq 0$, i.e. the MDP can transition from $\Omega_\alpha$ to $\Omega_\beta$ in a single step.*

**Assumption 1.** *Assume that $\Omega_\alpha, \Omega_\beta$ are MDP-neighbors for some MDP and that every trajectory of an optimal policy*
  *1) ends in $\Omega_\beta$ with a finite accumulated reward,*
  *2) if it is in $\Omega_\beta$ it will not leave $\Omega_\beta$,*
  *3) if it starts in $\Omega_\alpha$ it will transition to $\Omega_\beta$ without entering some other part of $S$ first.*

Note that the above assumption is satisfied for many problems with a large positive reward for transition to some states inside $\Omega_\beta$ ending the episode, and a smaller negative reward for all other transitions inside $\Omega_\alpha \cup \Omega_\beta$. This might correspond to a problem where the policy should reach a goal region inside $\Omega_\beta$ using e.g., minimum time or minimum energy. The two regions might correspond to a *move to* action and a *pick/push object* action, where it is clear that you have to move to the object before picking/pushing it. However, there are also many combinations of MDPs and regions that do not satisfy Assumption 1, i.e., when the optimal policy passes the switching boundary multiple times, so it needs to be checked for each case.

We will use the following definition to first solve a smaller MDP in $\Omega_\beta$ and then solve another smaller MDP in $\Omega_\alpha$, using the value function of the first one as part of the reward. We call the two smaller MDPs *restrictions* of the original MDP. We want them to be very similar to the larger MDP, but since an MDP cannot have transitions out of the state set, we need to add a set of states $S_{add}$ along the boundary that are absorbing (no transitions out), and give no rewards, see below.

**Definition 3** (Restriction). *By the restriction of a MDP $P_0 = (S, A, p, r)$ to $\hat{S} \subset S$ we mean a new MDP $\bar{P} = (\bar{S}, \bar{A}, \bar{p}, \bar{r})$ with a smaller set of states $\bar{S} = \hat{S} \cup S_{add}$, where*

$$S_{add} = \{s' \in S \setminus \hat{S} : \exists s \in \hat{S}, a \in A(s) \wedge p(s'|s,a) \neq 0\}.$$

*The available actions are the same $\bar{A}(s) = A(s)$. The transition probability is given by*

$$\bar{p}(s'|s,a) = \begin{cases} p(s'|s,a), & \text{if } s \in \hat{S} \\ 1, & \text{if } s = s', s \in S_{add} \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

*Note that this makes the states in $S_{add}$ absorbing, i.e. they can never be exited. The reward is given by*

$$\bar{r}(s,s',a) = \begin{cases} r(s,s',a), & \text{if } s, s' \in \hat{S} \\ r(s,s',a) + v_+(s'), & \text{if } s \in \hat{S}, s' \in S_{add} \\ 0, & \text{if } s \in S_{add}, \end{cases} \quad (7)$$

*where $v_+(s')$ is a given function impacting the reward when transitioning from $\hat{S}$ to $S_{add}$.*

Below we will use $v_+$ to penalise undesired transitions, and reward desired transitions based on the value function in the destination state.

**Lemma 2** (Decoupled solutions). *Given an MDP $P_0 = (S, A, p_0, r_0)$, with optimal value function $v_0^*(s)$, let Assumption 1 hold for two sets $\Omega_\alpha, \Omega_\beta$.*

*Let the MDP $P_\beta = (S_\beta, A_\beta, p_\beta, r_\beta)$ be the restriction of $P_0$ to $\Omega_\beta$ with*

$$v_+(s') = -\infty, \quad (8)$$

*having the corresponding optimal value function $v_\beta^*$.*

*Let the MDP $P_\alpha = (S_\alpha, A_\alpha, p_\alpha, r_\alpha)$ be the restriction of $P_0$ to $\Omega_\alpha$ with*

$$v_+(s') = \begin{cases} \gamma v_\beta^*(s'), & \text{if } s' \in \Omega_\beta \\ -\infty, & \text{otherwise} \end{cases} \quad (9)$$

*having the corresponding optimal value function $v_\alpha^*$.*

*Then the optimal value functions of $P_\alpha$ and $P_\beta$ are identical to the optimal value function of $P_0$ in $\Omega_\alpha$ and $\Omega_\beta$ respectively. That is*

$$v_0^*(s) = \begin{cases} v_\alpha^*(s), & \text{if } s \in \Omega_\alpha \\ v_\beta^*(s), & \text{if } s \in \Omega_\beta \end{cases} \quad (10)$$

Note that the lemma above states that we can solve $P_0$ by first solving the smaller MDP $P_\beta$ over $\Omega_\beta$ and then solving $P_\alpha$ over $\Omega_\alpha$ using optimal value function $v_\beta^*(s')$ on the boundary as $v_+$. Also note that once we have the optimal value function, we can easily find an optimal policy through Equation (5).

*Proof.* The optimal value function that solves Equation (4) for a given MDP is known to be unique [15]. Therefore we can assume that $v_0^*$ of $P_0$ is known, and if we can construct solutions $v_\alpha^*$ of $P_\alpha$ and $v_\beta^*$ of $P_\beta$ that satisfies Equations (4) and (10) we are done.

We start by looking at $P_\beta$. Let

$$v_\beta^*(s) = \begin{cases} v_0^*(s), & \text{if } s \in \Omega_\beta \\ 0, & \text{if } s \in S_{add} \end{cases} \quad (11)$$

Clearly this satisfies Equation (10), so it remains to show that it satisfies the Bellman optimality equation (4) for $P_\beta$.

We will show that the equations for both $P_0$ and $P_\beta$ do not depend on the values outside $\Omega_\beta$, and since the values inside $\Omega_\beta$ are identical, (4) must hold for $P_\beta$ if it holds for $P_0$.

For $P_0$ we know by Assumption 1 that any trajectory of an optimal policy $\pi_0^*$ of $P_0$ will remain in $\Omega_\beta$. Therefore $\pi_0^*$ must be such that $p(s'|s, \pi_0^*(s)) = 0$ for all $s \in \Omega_\beta, s' \notin \Omega_\beta$.

We also know that $\pi_0^*$ satisfies (5), therefore, for $s \in \Omega_\beta$, the action $a$ that maximizes the right hand side of (5) is such that $p(s'|s,a) = 0$ for all $s \in \Omega_\beta, s' \notin \Omega_\beta$. In (4) the same sum is maximized over $a$, and therefore $p(s'|s,a) = 0$ in (4) as well.

Thus, by (4), the values of $v_0^*(s)$ for $s \in \Omega_\beta$ do not depend on the values of $v_0^*(s)$ for $s \notin \Omega_\beta$.

For $P_\beta$ we have that states outside $\Omega_\beta$ are the absorbing states $S_{add}$. The reward of transferring to those is $-\infty$, since $\bar{r}(s,s',a) = r(s,s',a) + v_+(s') = r(s,s',a) - \infty = -\infty$ by (7). We know that there exists a policy that avoids leaving $\Omega_\beta$ and results in a finite accumulative reward, since the available actions are the same, $\bar{A}(s) = A(s)$. Since leaving $\Omega_\beta$ has a

Fig. 3: The basic case with $\Omega_\alpha = \Omega_1$ and $\Omega_\beta = \Omega_2$. The arrow indicates that trajectories of an optimal policy will always move from $\Omega_1$ to $\Omega_2$ and never in the opposite direction.

reward of $-\infty$ we conclude that the optimal policy for $P_\beta$ does not leave $\Omega_\beta$, and by the argument above, the values of $v_\beta^*(s)$ for $s \in \Omega_\beta$ do not depend on the values of $v_\beta^*(s)$ for $s \notin \Omega_\beta$. This concludes the proof regarding $P_\beta$.

For $P_\alpha$ we will explore how it depends on values outside $\Omega_\alpha$. First we note that $v_\alpha^*(s) = 0$ for $s \in S_{add}$ since $S_{add}$ are absorbing (no transitions out) and have reward 0 for staying by (7). If $s \in \Omega_\alpha, s' \in \Omega_\beta$ we note the following

$$p_\alpha(s'|s,a)[r_\alpha(s,s',a) + \gamma v_\alpha^*(s')], \tag{12}$$

$$p_\alpha(s'|s,a)[(r_0(s,s',a) + v_+(s')) + \gamma v_\alpha^*(s')] = \tag{13}$$

$$p_\alpha(s'|s,a)[(r_0(s,s',a) + \gamma v_\beta^*(s')) + 0] = \tag{14}$$

$$p_0(s'|s,a)[r_0(s,s',a) + \gamma v_0^*(s')] \tag{15}$$

where we have used (7) in the first equality, $v_\alpha^*(s') = 0$ and (9) in the second equality, and finally (6) and (11) to get the last equality.

Thus we see that the right hand side of (4) is the same for $P_0$ and $P_\alpha$ for $s \in \Omega_\alpha, s' \in \Omega_\beta$.

For transitions inside $\Omega_\alpha$, i.e. $s, s' \in \Omega_\alpha$ the right hand sides are clearly the same. Transitions leaving $\Omega_\alpha \cup \Omega_\beta$ are handled as above as we know the optimal policy of $P_0$ never leaves $\Omega_\alpha \cup \Omega_\beta$, and the reward in $P_\alpha$ for leaving $\Omega_\alpha \cup \Omega_\beta$ is $-\infty$. Thus values outside $\Omega_\alpha \cup \Omega_\beta$ does not influence (4) for $s \in \Omega_\alpha$, and states $s' \in \Omega_\beta$ produce identical contributions by (15). Therefore, if $v_0^*(s)$ satisfies (4) of $P_0$ so does $v_\alpha^*(s)$ of $P_\alpha$. $\square$

**Example 1.** *Consider a ball-shaped agent that is to push a box to a given goal region, modelled as an MDP $P_0$. It is clear that each successful policy must first move to the box, and then push it into the goal region (as you cannot push without being close to something). Thus we can divide $S$ into $\Omega_1$ (not close to box) and $\Omega_2$ (close to box).*

*Lemma 2 now tells us that we can first solve the ball pushing MDP $P_2$, and then solve the move to MDP $P_1$ using the value function from $P_2$ (which basically describes what positions around the box are beneficial for pushing it to the goal), as input.*

*Having done this, Lemma 2 tells us that the optimal value functions of $P_0$ will be the same as $P_1, P_2$ on $\Omega_1, \Omega_2$ correspondingly. Since the optimal policy can be found from the optimal value function through (5), we have found the optimal policy to the original problem by solving the two smaller problems.*

*A detailed version of this example can be found in Section V below.*

The next result concerns the case where we have a MDP and want to use an existing policy for parts of the solution. Examples include a classical PID controller, inverse kinemat-

ics for bringing a robot arm to some configuration, or some other policy that we want to reuse.

**Lemma 3** (Constraining the optimal policy)**.** *Let the MDP $P_0 = (S, A(\cdot), p, r)$ and a policy $\pi_g : S \to A$ that we must execute in some domain $\Omega_g \subset S$ be given.*

*If we define a new function $\bar{A} : S \to A$ such that*

$$\bar{A}(s) = \begin{cases} \{\pi_g(s)\}, & \text{if } s \in \Omega_g \\ A(s), & \text{otherwise,} \end{cases} \tag{16}$$

*then the optimal solution $\pi_1^*(s)$ of the new MDP $P_1 = (S, \bar{A}(\cdot), p, r)$ is such that $\pi_1^*(s) = \pi_g(s)$ when $s \in \Omega_g$.*

*Proof.* Any policy in the new MDP must satisfy $\pi(s) = \pi_g(s)$ when $s \in \Omega_g$, as this is the only available action. Thus it also holds for the optimal policy.

$\square$

**Lemma 4** (Combining manual and learned sub-polices)**.** *Given a MDP $P_0 = (S, A, p, r)$ and a given policy $\pi_g : S \to A$ that we want to execute in some domain $\Omega_g \subset S$.*

*If we constrain the available actions $A(\cdot)$ as described in Lemma 3 we can still apply Lemma 2 if the new MDP also satisfies Assumption 1.*

*Proof.* This is clear since applying Lemma 3 just produces a new MDP, and Lemma 2 can be applied to any MDP satisfying Assumption 1. $\square$

**Remark 1.** *Note that this can be practical if there exists a manually designed controller, e.g., a PID- or LQR controller that we want to combine with a learned controller in an optimal way. If $\Omega_g = \Omega_\beta$, the optimal policy in all of $\Omega_\beta$ is already known, and we can compute the optimal value function $v_\beta^*(s)$ in $\Omega_\beta$ using e.g. policy evaluation, as suggested in [16].*

**Example 2.** *Looking back at the ball pushing problem in Example 1, we might have an existing pushing policy that we want to use. If that is the case we can constrain the available actions according to this policy, compute the value function for the given policy, and then find the optimal move to behavior using the value function of the manual push-behavior. A detailed version of this example can be found in Section V below.*

**Lemma 5** (Recursive application over many policies)**.** *Given a MDP $P_0 = (S, A, p, r)$ where $S$ is divided into a set of disjoint operating regions $\Omega_i$, such that $S = \cup_i \Omega_i$, as illustrated in Figure 4. Assume the optimal policy has a finite accumulated reward from all starting states.*

*Let $M \subset \mathbb{N}$ be the indices of existing policies $\pi_i : \Omega_i \to A$ for $i \in M$ we want to use.*

*First we constrain the MDP with respect to these controllers, according to Lemma 3.*

*If the $\Omega_i$ are numbered such that transitions of optimal policies will always happen from a lower index to a higher index in this constrained MDP, we improve the policy in any region $\Omega_i, i \notin M$ by letting $\Omega_\alpha = \Omega_i$ and $\Omega_\beta = \cup_{j>i}\Omega_j$, and applying Lemma 2.*

*If we recursively apply this strategy backwards from the highest index, we will recreate the globally optimal policy.*

Fig. 4: Illustration of the iterative application of the main result. If we know that trajectories of optimal policies will always move from $\Omega_i$ to $\Omega_j$, with $i < j$, and finally stay in $\Omega_7$, we can apply Lemma 5 with $\Omega_\alpha = \Omega_i$ and $\Omega_\beta = \cup_{j>i}\Omega_j$, starting from the back with $\Omega_\alpha = \Omega_6$ and $\Omega_\beta = \Omega_7$, then $\Omega_\alpha = \Omega_5$ and $\Omega_\beta = \Omega_6 \cup \Omega_7$, and so on.

*Proof.* Since we know that the optimal policy has a finite accumulated reward, and all transitions of the optimal policy will happen to a region with higher index $i$, items 1, 2 and 3 of Assumption 1 is satisfied for all $\Omega_\alpha, \Omega_\beta$ constructed as above. If the assumption is satisfied, we can apply the lemma to improve performance. Since the solution in $\Omega_\alpha$ depends on the solution in $\Omega_\beta$, we will get the optimal solution by starting with the highest index and going backwards. □

## V. NUMERICAL EXAMPLES

A simple numerical example was already described in Section I and Figures 1 and 2 above. To see how the theoretical results from Section IV apply to a more dynamic example, we implemented Examples 1 and 2 above using the Unity Engine and the RL framework called ML-agents [17]. This section will give an overview of the setting, configurations and results.

### A. Motivation of the setup

To illustrate the theoretical results above we picked a simple example with just two policies. This example shows how the approach works, but not why it is needed, as the example can easily be solved by using a standard single RL policy, as is done for performance comparison in the experiment. To motivate the approach, we note that this example is not a typical intended use case. As discussed above, BTs have been shown to be optimally modular [18], and modularity is important only when you have a complex system with many policies, such as the one in Figure 1 of [8]. This is a typical use case, with a BT having 13 conditions to elaborately switch between 8 different actions. However, the principles are the same, and to enable us to go into details, we pick the simple example above.

### B. Scenario setup

As in Examples 1 and 2 above, and illustrated in Figure 5, the environment consists of an agent in the shape of a ball, a target object in the shape of a box, and a rectangular goal area. The environment is initialized by setting a random position for the agent and the target object, as well as a randomly chosen edge for the goal area. The task of the agent is to push the



Fig. 5: A snapshot of the experimental environment including a blue spherical agent, a blue goal area rectangle and a red target box. The agent fails if the target or the agent leave the plane, and succeeds if the target reaches the blue goal area.



Fig. 6: The simple BT of the example on the left, with one condition and two policies, and the operating regions on the right. Move to box will execute in $\Omega_1$ and Push Box will execute in $\Omega_2$.

target to the goal area. It can fail by exiting the mission area, or moving the target outside of the mission area. The agent's controller is split into two sub-behaviors. A *Move To* behavior, which is active when the agent is far away from the box, and a *Push* behavior which is activated if the agent is close to the box. Thus $\Omega_\beta = \{s \in S : ||p_{agent} - p_{box}|| < d\}$ and $\Omega_\alpha = S \setminus \Omega_\beta$, with $d$ equal to 2.5 times the side of the box.

The target and the agent are subject to second-order dynamics, including friction, as implemented by the physics engine in Unity. The agent is controlled through applying a force in the horizontal plane. The agent observes it's own position and velocity, as well as the relative position of the goal and the target box, making the state space $S$ 8-dimensional.

### C. The different controllers

First we will create five different sub-behaviors, listed on the left of Table I. *Move To (Local)* will be created using RL aiming to reach the box as fast as possible. *Move To (VF)* will be created using RL aiming to quickly reach a position in $\Omega_\beta$ with a good value function of the following push behavior. *Push (RL)* will be created using RL, aiming to quickly get the box to the goal area when starting in $\Omega_\beta$. *Push (Manual)* will be created manually, aiming to do a decent job of pushing the box to the goal area when starting in $\Omega_\beta$. Finally, for comparison, we create a *Single Behavior* for completing the entire task in $\Omega_\alpha \cup \Omega_\beta$ as fast as possible, using RL. Note that the *Push (Manual)* controller will be intentionally suboptimal to highlight how this effects the overall system performance. Thus *Push (Manual)* first moves straight towards a position opposite the goal with respect to the target and then moves the ball straight towards the box.

Then we create five different combinations of the sub-behaviors, as listed on the left of Table II.

| Behavior | Reward every timestep | Completion | Fail |
|---|---|---|---|
| Move To (Local) | $-0.001$ | 1 | -1 |
| Move To (VF) | $-0.001$ | $\gamma v_\beta^*(s)$ | -1 |
| Push (RL) | $-0.001 + \Delta_{dist}$ | 1 | -1 |
| Push (Manual) | $-0.001 + \Delta_{dist}$ | 1 | -1 |
| Single Behavior | $-0.001 + \Delta_{dist}$ | 1 | -1 |

TABLE I: Rewards for the five different policies. In the case of the manually designed push behavior, the policy is not trained, but a critic network to estimate the value function (VF) $v_\beta^*$ is.

| Overall Policy | Accumulated Reward | Success Duration (steps) | Failure Duration (steps) | Success Rate (%) |
|---|---|---|---|---|
| $\pi_A$, Move To(Local) + Push (RL) | $1.39 \pm 0.48$ | $248.02 \pm 90.97$ | $171.12 \pm 58.85$ | 95.20 |
| $\pi_B$, Move To(Local) + Push (Manual) | $-0.07 \pm 1.20$ | $237.88 \pm 120.40$ | $148.31 \pm 57.84$ | 39.59 |
| $\pi_C$, Move To(VF) + Push (RL) | $\mathbf{1.57 \pm 0.18}$ | $204.01 \pm 56.21$ | $175.65 \pm 45.15$ | $\mathbf{99.86}$ |
| $\pi_D$, Move To(VF) + Push (Manual) | $1.40 \pm 0.51$ | $275.67 \pm 138.14$ | $280.56 \pm 153.34$ | 96.69 |
| $\pi_E$, Single Behavior | $1.52 \pm 0.17$ | $\mathbf{192.49 \pm 56.41}$ | $128.50 \pm 80.01$ | 98.76 |

TABLE II: Evaluated policies. Mean and standard deviation across 10 000 random episodes post training. Best results indicated by bold numbers.

Note that the two *Move To (VF)* are different, as they are trained with different value functions on the boundary, coming from *Push RL* and *Push Manual* respectively.

### D. RL formulation

The rewards of the RL are summarized in Table I. All sub-behaviors receive a small negative reward of $-0.001$ for each passing time step and a small reward $\Delta_{dist}$ for moving the box closer to the goal area, $\Delta_{dist} = (d_{goal}(s_{t-1}) - d_{goal}(s_t))/d_{start}$, with $d_{goal}$ being the distance from the box to the goal area, and $d_{start}$ this distance at the start, for normalization. At the end of the episode, all sub-behaviors receive a negative reward of $-1$ if the agent or the target leaves the mission area, and a positive reward of 1 or $\gamma v_\beta^*(s)$ for completion of the task. Above, $v_\beta^*(s)$ is the value function of the following sub-behavior, i.e., either the manual or RL-version of Push. All actions and observations are normalised to a range of $[-1, 1]$.

The training is done using the built-in implementation of PPO [19] from the Unity ML-Agents package [20]. In the case of the manual behavior *Push (Manual)*, we only train a critic network to estimate the value function, using the normal PPO loss function [19], with the policy loss removed.

The RL was executed by training each sub-behavior when active in a single environment, as opposed to using separate training environments for each behavior. The agent operates according to the policy of the currently active behavior, using the gathered experiences for training that specific behavior. Upon crossing a switching boundary, the ongoing episode ends, and a new one is started for the switched-to policy.

As suggested in Lemma 2, we first train *Move To (Local)* and *Push (RL)*, and estimate the value function of *Push (Manual)*. This creates the components of $\pi_A$ and $\pi_B$. Then we train the two versions of *Move To (VF)*, using the value functions from *Push (RL)* and *Push (Manual)* respectively, creating the components of $\pi_C$ and $\pi_D$. Finally, we train *Single Behavior*, for $\pi_E$.

### E. Results

We present the data in three parts. First, we discuss post-



Fig. 7: Training results for $\pi_A$, $\pi_C$ and $\pi_E$, using Push (RL).



Fig. 8: Training results for $\pi_B$, $\pi_D$ and $\pi_E$, using Push (Manual).

training evaluation results. Providing context to these results, we then analyze reward graphs and histograms of the training data, comparing outcomes with and without utilizing the value function as a reward signal. Finally, we measure the disparity between the value function of the single learned policy and those of the split models to illustrate the theory in Lemma 2.

Table II presents an overall evaluation of running five policies from 10 000 random starting states. As theory suggests, $\pi_C$ and $\pi_E$ perform best, with slight differences. $\pi_C$ yields marginally higher rewards and fewer failures, and $\pi_E$ completing tasks a bit faster when not failing. Both $\pi_C$ and $\pi_E$ outperform the locally optimal $\pi_A$. Comparing versions with manual push behavior, $\pi_B$ has a significant 60% failure rate, while despite the sub-optimal manual behavior, $\pi_D$ achieves success rates and accumulated rewards comparable to RL policies like $\pi_A$, albeit with slightly longer completion times.

The historical training reward, per training session, can be seen in Figures 7 and 8. Both graphs feature the single model behavior as a baseline for comparison. We can see that after an initial decrease in reward, all training configurations have a stable upwards trend that reaches an equilibrium, with the reward plateauing around a specific value.

In Figure 7 we see the training results of $\pi_A$, $\pi_C$ and $\pi_E$. $\pi_E$ is trained first to use as a reference. Then we train $\pi_A$ that does not use the value function to connect the two sub-behaviors. As can be seen, the mean reward converges quickly for *Move To (Local)*(green) and *Push (RL)*(red). Then we train $\pi_C$, by training *Move To (VF)*(blue) using $v_\beta^*(s)$ from *Push (RL)*. As can be seen, this achieves a higher mean reward than *Move To (Local)*. This is reasonable, as the final reward given by $v_\beta^*(s)$ is sometimes larger than one.

Fig. 9: Mean absolute difference between the value functions of the experimental configuration and single learned model value function with 95% confidence interval. Computed over 100 example states per model iteration, 625 samples per state.

In Figure 8 we see the training results of $\pi_B$, $\pi_D$ and $\pi_E$. When we train $\pi_B$ an interesting thing happens. As *Move To (Local)* gets better, the performance of the fixed manual controller *Push (Manual)*(red), decreases. This is due to the fact that the manual controller is not very robust, and as *Move To (Local)* learns to enter $\Omega_\beta$ in close to minimum time, using high velocities, the starting states handed to *Push (Manual)* become more difficult.

When training $\pi_D$ we used the value function $v_\beta^*(s)$ from *Push (Manual)* as final reward. In this way, *Move To (VF)* is made aware of the capabilities of the following *Push (Manual)*. The average reward of *Move To (VF)* rises slower than *Move To (Local)*, but at the same time the average reward of the static *Push (Manual)* continues to increase. In the end, the combined result of the constrained $\pi_D$ is slower, but almost as reliable as the near optimal unconstrained $\pi_E$, see Table II.

### F. Comparing the Value Functions

Lemma 2 shows that the optimal value functions of the original MDP $P_0$ is identical to the optimal value functions of the two restricted MDPs $P_\alpha$ and $P_\beta$.

Running an RL algorithm we know that as time tends to infinity, the learned value function will converge towards the optimal one [16]. We measured the difference between the value function of $\pi_E$ and the two restricted ones for the three policies $\pi_A$, $\pi_C$ and $\pi_D$. The results, seen in Figure 9, show that the difference is smallest for $\pi_C$, as predicted by the Lemma. Note also that after the initial drop, the difference to $\pi_A$ and $\pi_D$ actually increases as training progresses, which is reasonable as they are not expected to converge to the same optimal value function.

The difference in value functions is more distinct in the Manual case, since the *Push* behavior is fixed causing the *MoveTo* to try to compensate. This means that the overall behavior will diverge more strongly from the optimal single behavior case, with both sub-behavior policies shaped differently. With $\pi_A$, the difference would primarily be in the *MoveTo* behavior, causing the error the be smaller.

## VI. CONCLUSIONS

We investigated solutions to the local optimality issue that occurs when trying to combine learned controllers into a handcrafted policy. Leveraging the BT structure's deterministic switching between policies, we used the value function estimate of the controller being switched to as a final reward of the previous controller. The proposed approach goes beyond the state of the art in that we provide theoretical guarantees of optimality, are not using an arbitrary design parameter, and can handle problems where information has to flow across more than one switching boundary.

## REFERENCES

[1] M. Iovino, E. Scukins, J. Styrud, P. Ögren, and C. Smith, "A survey of behavior trees in robotics and ai," *Robotics and Autonomous Systems*, vol. 154, p. 104096, 2022.

[2] M. Colledanchise and P. Ögren, *Behavior trees in robotics and AI: An introduction*. CRC Press, 2018.

[3] R. d. P. Pereira and P. M. Engel, "A framework for constrained and adaptive behavior-based agents," *arXiv preprint arXiv:1506.02312*, 2015.

[4] D. Isla, "Handling Complexity in the Halo 2 AI," in *Proceedings of the Game Developers Conference (GDC)*, 2005.

[5] O. Biggar, M. Zamani, and I. Shames, "On modularity in reactive control architectures, with an application to formal verification," *ACM Transactions on Cyber-Physical Systems (TCPS)*, May 2022.

[6] R. R. Burridge, A. A. Rizzi, and D. E. Koditschek, "Sequential Composition of Dynamically Dexterous Robot Behaviors," *The International Journal of Robotics Research*, vol. 18, no. 6, pp. 534–555, June 1999.

[7] J. Erskine and C. Lehnert, "Developing cooperative policies for multistage reinforcement learning tasks," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 6590–6597, 2022.

[8] P. Ögren and C. I. Sprague, "Behavior Trees in Robot Control Systems," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, no. 1, 2022.

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

[10] B. Hannaford, D. Hu, D. Zhang, and Y. Li, "Simulation results on selector adaptation in behavior trees," *arXiv preprint arXiv:1606.09219*, 2016.

[11] Y. Fu, L. Qin, and Q. Yin, "A reinforcement learning behavior tree framework for game ai," in *Proceedings of the 2016 International Conference on Economics, Social Science, Arts, Education and Management Engineering*, 2016.

[12] Q. Zhang, L. Sun, P. Jiao, and Q. Yin, "Combining behavior trees with maxq learning to facilitate cgfs behavior modeling," in *2017 4th International Conference on Systems and Informatics (ICSAI)*. IEEE, 2017, pp. 525–531.

[13] R. S. Sutton, D. Precup, and S. Singh, "Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning," *Artificial Intelligence*, vol. 112, no. 1-2, pp. 181–211, Aug. 1999.

[14] Y. Lee, S.-H. Sun, S. Somasundaram, E. S. Hu, and J. J. Lim, "Composing complex skills by learning transition policies," in *International Conference on Learning Representations*, 2018.

[15] T. G. Dietterich, "Hierarchical Reinforcement Learning with the MAXQ Value Function Decomposition," *Journal of Artificial Intelligence Research*, vol. 13, pp. 227–303, Nov. 2000.

[16] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, second edition ed., ser. Adaptive Computation and Machine Learning Series. Cambridge, Massachusetts: The MIT Press, 2018.

[17] A. Juliani, V.-P. Berges, E. Teng, A. Cohen, J. Harper, C. Elion, C. Goy, Y. Gao, H. Henry, M. Mattar, and D. Lange, "Unity: A general platform for intelligent agents," *arXiv preprint arXiv:1809.02627*, 2020.

[18] O. Biggar, M. Zamani, and I. Shames, "On modularity in reactive control architectures, with an application to formal verification," *arXiv preprint arXiv:2008.12515*, 2020.

[19] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," *arXiv*, Aug. 2017.

[20] A. Juliani, V.-P. Berges, E. Teng, A. Cohen, J. Harper, C. Elion, C. Goy, Y. Gao, H. Henry, M. Mattar, and D. Lange, "Unity: A general platform for intelligent agents," *arXiv preprint arXiv:1809.02627*, 2020.