

SYMPOL: SYMBOLIC TREE-BASED ON-POLICY REINFORCEMENT LEARNING

Anonymous authors
 Paper under double-blind review

ABSTRACT

Reinforcement learning (RL) has seen significant success across various domains, but its adoption is often limited by the black-box nature of neural network policies, making them difficult to interpret. In contrast, symbolic policies allow representing decision-making strategies in a compact and interpretable way. However, learning symbolic policies directly within on-policy methods remains challenging. In this paper, we introduce SYMPOL, a novel method for SYMBOLic tree-based on-POLicy RL. SYMPOL employs a tree-based model integrated with a policy gradient method, enabling the agent to learn and adapt its actions while maintaining a high level of interpretability. We evaluate SYMPOL on a set of benchmark RL tasks, demonstrating its superiority over alternative tree-based RL approaches in terms of performance and interpretability. In contrast to existing methods, SYMPOL allows a gradient-based end-to-end learning of interpretable, axis-aligned decision trees within existing on-policy RL algorithms. Therefore, SYMPOL can become the foundation for a new class of interpretable RL based on decision trees.

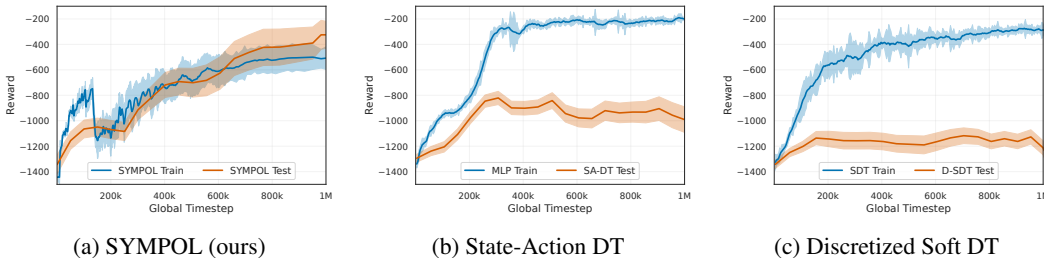


Figure 1: **Information Loss in Tree-Based Reinforcement Learning on PD-C.** Existing methods for symbolic, tree-based RL (see Figure 1b and 1c) suffer from severe information loss when converting the differentiable policy used for training (e.g., the MLP for SA-DT) into the symbolic policy used for interpretation (i.e., the DT). Using SYMPOL (Figure 1a), we can directly optimize the symbolic policy with PPO and therefore have no information loss during the application.

1 INTRODUCTION

Reinforcement learning lacks transparency. Reinforcement learning (RL) has achieved remarkable success in solving complex sequential decision-making problems, ranging from robotics and autonomous systems to game playing and recommendation systems. However, the policies learned by traditional RL algorithms, represented by Neural Networks (NNs), often lack interpretability and transparency, making them difficult to understand, trust, and deploy in safety-critical or high-stakes scenarios (Landajuela et al., 2021).

Symbolic policies increase trust. Symbolic policies, on the other hand, offer a promising alternative by representing decision-making strategies in terms of RL policies as compact and interpretable structures (Guo et al., 2024). These symbolic representations do not only facilitate human understanding and analysis but also ensure predictable and explainable behavior, which is crucial for building trust and enabling effective human-AI collaboration. Moreover, the deployment of sym-

054 bolic policies in safety-critical systems, such as autonomous vehicles or industrial robots, could
055 significantly improve their reliability and trustworthiness. By providing human operators with a
056 clear understanding of the decision-making process, symbolic policies can facilitate effective moni-
057 toring, intervention, and debugging, ultimately enhancing the safety and robustness of these systems.
058 In this context, decision trees (DTs) are particularly effective as symbolic policies for RL, as their
059 hierarchical structure provides natural interpretability.

060 **Existing challenges.** Despite these promising prospects, the field of symbolic RL faces several
061 challenges. One main reason is given by the fact that many symbolic models, like DTs, are non-
062 differentiable and cannot be integrated in existing RL frameworks. Therefore, traditional methods
063 for learning symbolic policies often rely on custom and complex training procedures (Costa et al.,
064 2024; Vos & Verwer, 2023; Kanamori et al., 2022), limiting their applicability and scalability. Al-
065 ternative methods involve pre-trained NN policies combined with some post-processing to obtain
066 an interpretable model (Silva et al., 2020; Liu et al., 2019; 2023; Bastani et al., 2018). However,
067 post-processing introduces a mismatch between the optimized policy and the model obtained for
068 interpretation, which can lead to loss of crucial information, as we show in Figure 1.

069 **Contribution.** To mitigate the impact of information loss, a direct optimization of the policy is
070 crucial. Existing methods that employ a direct optimization of a tree-based policy, such as those
071 described by Silva et al. (2020) learn differentiable, soft decision trees which do not provide a high
072 level of interpretability. To obtain interpretable, axis-aligned DTs, these methods require post-hoc
073 distillation or discretization and therefore suffer from information loss (see Figure 1). In this paper,
074 we introduce SYMPOL, SYMBolic tree-based on-POLicy RL, a novel method employing a direct
075 optimization axis-aligned DT policies end-to-end. Our contributions are as follows:

- 076 • We integrate GradTree (Marton et al., 2024a) into existing RL frameworks via a separate
077 actor-critic architecture to directly optimize DT policies and extend it to continuous action
078 spaces (Section 4.1).
- 079 • We propose a dynamic rollout buffer to enhance exploration stability and a dynamic batch
080 size through gradient accumulation to improve gradient stability (Section 4.2) to mitigate
081 the instability of DT training in dynamic environments.
- 082 • We propose using weight decay on a subset of parameters to support a dynamic adjustment
083 of the model parameters when optimizing DTs with gradient descent (Section 4.1).

084 As a result, SYMPOL does not depend on pre-trained NN policies, complex search procedures, or
085 post-processing steps, but can be seamlessly integrated into existing RL algorithms (Section 3).

086 **Results.** Through extensive experiments on benchmark RL environments, we demonstrate that
087 SYMPOL does not suffer from information loss and outperforms existing tree-based RL approaches
088 in terms of interpretability and performance (Section 5.2), providing human-understandable expla-
089 nations. In most environments, SYMPOL’s performance is comparable to full-complexity models,
090 while in categorical environments, it even surpasses them. Furthermore, we provide a case study
091 (Section 6) to show how interpretable policies help in detecting misbehavior and misgeneralization
092 which might remain unnoticed with commonly used black-box policies.

093 2 RELATED WORK

094 Recently, the integration of symbolic methods into RL has gained significant attention. Symbolic
095 RL does cover different approaches including program synthesis (Trivedi et al., 2021; Penkov &
096 Ramamoorthy, 2019; Verma et al., 2018), concept bottleneck models (Ye et al., 2024), piecewise
097 linear networks (Wabartha & Pineau, 2024), logic (Delfosse et al., 2024b), mathematical expres-
098 sions (Landajuela et al., 2021; Guo et al., 2024; Luo et al., 2024; Xu et al., 2022). Another line
099 of work aims to synthesize symbolic policies using logical rules, leveraging differentiable inductive
100 logic programming for gradient-based optimization (Jiang & Luo, 2019; Cao et al., 2022). In con-
101 trast to first-order rules, DTs offer greater flexibility by not only combining atomic conditions but
102 also comparing features against thresholds — a critical capability for handling continuous observa-
103 tion spaces. In this paper, we focus exclusively on tree-based methods for symbolic RL. Several
104 approaches have been proposed to leverage the strengths of interpretable, tree-based representations

108 within RL frameworks. However, each approach comes with its own critical limitations. We sum-
 109 marize existing methods into three streams of work:

110
 111 **(1) Post-processing.** One line learns full-complexity policies first and then performs some kind of
 112 post-processing for interpretability. One prominent example is the VIPER algorithm (Bastani et al.,
 113 2018). In this case, a policy is learned using NNs before DTs are distilled from the policy. However,
 114 distillation methods often suffer from significant performance mismatches between the training and
 115 evaluation policies (see Figure 1b). To mitigate this mismatch, existing methods often learn large
 116 DTs (VIPER learns DTs with 1,000 nodes) and therefore aim for systematic verification rather than
 117 interpretability. In contrast, SYMPOL is able to learn small, interpretable DTs (average of only 50
 118 nodes) without information loss. Following VIPER, various authors proposed similar distillation
 119 methods (Li et al., 2021; Liu et al., 2019; 2023; Jhunjunwala et al., 2020). Furthermore, Kohler
 120 et al. (2024) propose a novel distillation method that distills interpretable and editable programmatic
 121 tree policies. In contrast to SYMPOL, the extracted trees are not considered axis-aligned as they
 allow for linear combinations and multiple features within the internal nodes.

122 **(2) Custom optimization.** The third line involves custom, tree-specific optimization techniques
 123 and/or objectives (Ernst et al., 2005; Roth et al., 2019; Gupta et al., 2015; Kanamori et al., 2022)
 124 and, hence, is more time-consuming and less flexible. As a result, their policy models cannot be
 125 easily integrated into existing learning RL frameworks. Examples are evolutionary methods (Costa
 126 et al., 2024; Custode & Iacca, 2023) and linear integer programming (Vos & Verwer, 2023). Topin
 127 et al. (2021) propose Iterative Bounding Markov Decision Process (IBMDP) that allow learning
 128 DT policies through a masking procedure and modified value updates by using arbitrary function
 129 approximators. However, using IBMDP, the learning problem becomes more complex compared
 130 to the base MDP, which can result in poor scalability and limits the applicability to very simple
 131 tasks (Milani et al., 2022; Kohler et al., 2024). In contrast, SYMPOL optimizes a DT policy directly
 132 on the base MDP, avoiding these limitations.

133 **(3) Soft Decision Trees (SDTs).** Methods optimizing SDTs (Silva et al., 2020; Silva & Gom-
 134 bolay, 2021; Coppens et al., 2019; Tambwekar et al., 2023; Liu et al., 2022; Farquhar et al., 2017)
 135 are difficult to interpret since they usually involve multiple features simultaneously at each decision
 136 node, creating complex, multidimensional splits rather than straightforward, single-feature thresh-
 137 olds. Nevertheless, the trees are usually not easily interpretable and techniques such as discretizing
 138 the learned trees into more interpretable representations are applied (Silva et al., 2020), occasionally
 139 resulting in high performance mismatches (Figure 1c). In contrast, SYMPOL directly optimizes
 140 hard, axis-aligned DTs and therefore does not exhibit a performance loss (Figure 1a).

141 **Distinction of SYMPOL from Differentiable and Soft Decision Trees.** In existing work like
 142 Silva et al. (2020), differentiable decision trees typically correspond to SDTs, achieving differenti-
 143 ability by relaxing discrete decisions in terms of feature selection at each internal node and path
 144 selection. This approach is fundamentally different from SYMPOL, which does *not* use differenti-
 145 able decision trees. Instead, SYMPOL leverages GradTree to optimize standard, non-differentiable
 146 decision trees through gradient descent, as we will show in Section 4.

147 Furthermore, trees have also been in used in other agentic components than the policy, such as
 148 reward functions (Milani et al., 2022; Kalra & Brown, 2023; 2022). Similarly, ensemble methods
 149 (Fuhrer et al., 2024; Min & Elliott, 2022) have been proposed. However, policies consisting of
 150 hundreds of trees and nodes lack interpretability and therefore are out of scope for this paper.

151 3 PRELIMINARIES

152
 153
 154 **Markov Decision Process.** We study a deterministic Markov decision process $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$
 155 where \mathcal{S} is a finite state space, \mathcal{A} is the finite action space, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ defines the tran-
 156 sition dynamics of the environment, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function and γ the discount factor.
 157 At each timestep t , an agent samples an action from policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ based on the current observa-
 158 tion $s_t \in \mathcal{S}$ and executes it in the environment. The environment transitions and the agent receives
 159 a reward r_t . In this context, the value function $\mathcal{V}^\pi(s) = \mathbb{E}_{a_t \sim \pi, s_{t+1} \sim \mathcal{P}} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_t = s]$
 160 approximates the expected return when starting in state s and then acting according to policy π .
 161 Similarly, the action-value function $\mathcal{Q}^\pi(s, a) = \mathbb{E}_{a_t \sim \pi, s_{t+1} \sim \mathcal{P}} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_t = s, a_t = a]$
 estimates the expected return when selecting action a in state s and then following policy π . Finally,

the advantage function $\mathcal{A}^\pi(s, a) = \mathcal{Q}^\pi(s, a) - \mathcal{V}^\pi(s)$ defines the difference between the expected return when choosing action a in state s and the expected return when following the policy π from state s . Overall, we aim for finding an optimal policy π^* that maximizes the expected discounted return $J(\pi) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$.

Proximal Policy Optimization (PPO). PPO (Schulman et al., 2017) is an on-policy, actor-critic RL method designed to enhance the training stability. The algorithm introduces a clipped surrogate objective to restrict the policy update step size. The main idea is to constrain policy changes to a small trust region, preventing large updates that could destabilize training. Formally, PPO optimizes:

$$\mathcal{L}^{\text{CLIP}}(\theta) = \mathbb{E}_{a_t \sim \pi_{\theta_{\text{old}}}, s_{t+1} \sim \mathcal{P}} \left[\min \left(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \hat{A}_t, \text{clip} \left(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right] \quad (1)$$

where $\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ is the probability ratio between the new policy π_θ and old policy $\pi_{\theta_{\text{old}}}$. \hat{A}_t is an estimate of the advantage function at time step t and ϵ is a hyperparameter for the clipping range.

4 SYMPOL: SYMBOLIC ON-POLICY RL

In the following, we formalize the online training of hard, axis-aligned DTs with the PPO objective. In contrast to existing work on RL with DTs, this allows an optimization of the DT on-policy without information loss. The main conceptual difference to existing work that learn symbolic policies end-to-end (Delfosse et al., 2024b; Fuhrer et al., 2024; Delfosse et al., 2024c; Topin et al., 2021; Luo et al.) is that SYMPOL does *not* require any modification of the RL framework itself, making the proposed method framework-agnostic. As a result, interpretable policies with SYMPOL are learned in the same way as NN policies are commonly learned. In the main paper, we focus on PPO as the, we believe, most prominent on-policy RL method. To support our claim of seamless integration, we provide additional results using Advantage Actor-Critic (A2C) (Mnih et al., 2016) in Appendix A.1 To efficiently learn DT policies with SYMPOL, we employed several crucial (see ablation study in Table 5) modifications, which we will elaborate below.

4.1 LEARNING DTs WITH POLICY GRADIENTS

SYMPOL utilizes GradTree (Marton et al., 2024a) as a core component to learn a DT policy directly from policy gradients as we will show in the following.

Arithmetic DT policy formulation. Traditionally, DTs involve nested concatenations of rules. In GradTree, DTs are formulated as arithmetic functions based on addition and multiplication to facilitate gradient-based learning. Therefore, our resulting DT policy is fully-grown (i.e., complete, full) and can be pruned post-hoc. Our basic pruning involves removing redundant paths, which significantly reduces the complexity. We define a path as redundant if the decision is already determined either by previous splits or based on the range of the selected feature. More details are given in Appendix A.4. Overall, we formulate a DT policy π of depth d with respect to its parameters as:

$$\pi(s|\mathbf{a}, \boldsymbol{\tau}, \boldsymbol{\iota}) = \sum_{l=0}^{2^d-1} a_l \mathbb{L}(s|l, \boldsymbol{\tau}, \boldsymbol{\iota}) \quad (2)$$

where \mathbb{L} is a function that indicates whether a state $s \in \mathbb{R}^{|\mathcal{S}|}$ belongs to a leaf l , $\mathbf{a} \in \mathcal{A}^{2^d}$ denotes the selected action for each leaf node, $\boldsymbol{\tau} \in \mathbb{R}^{2^d-1}$ represents split thresholds and $\boldsymbol{\iota} \in \mathbb{N}^{2^d-1}$ the feature index for each internal node.

Dense architecture. To support a gradient-based optimization and ensure an efficient computation via matrix operations, we make use of a dense DT representation. Traditionally, the feature index vector $\boldsymbol{\iota}$ is one-dimensional. However, as in GradTree, we expand it into a matrix form. Specifically, this representation one-hot encodes the feature index, converting $\boldsymbol{\iota} \in \mathbb{R}^{2^d-1}$ into a matrix $\mathbf{I} \in \mathbb{R}^{(2^d-1) \times |\mathcal{S}|}$. Similarly, for split thresholds, instead of a single value for all features, individual values for each feature are stored, leading to $\mathbf{T} \in \mathbb{R}^{(2^d-1) \times |\mathcal{S}|}$. The dense representation is visualized in Figure 2. Please note, that in contrast to SDTs, the dense representation of SYMPOL corresponds to an equivalent standard DT representation at each point in time, ensuring that the underlying model is a hard, axis-aligned DT. By enumerating the internal nodes in breadth-first order,

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

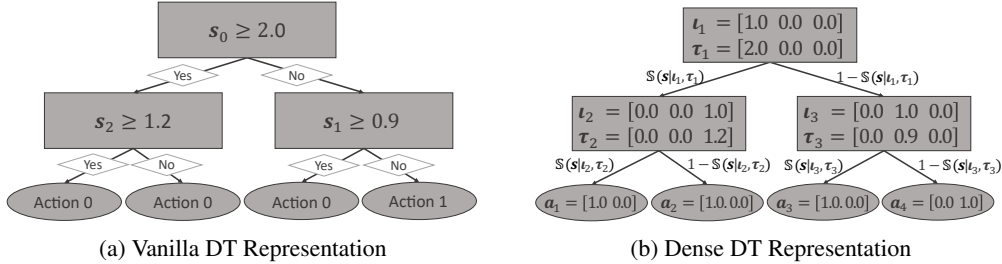


Figure 2: **Standard vs. Dense DT Representation.** A comparison between the standard decision tree representation and its dense equivalent, illustrated using an example decision tree of depth 2, with a state space of dimensionality 3 and two possible actions.

we can redefine the indicator function \mathbb{L} for a leaf l , resulting in

$$\pi(\mathbf{s}|\mathbf{a}, T, \mathbf{I}) = \sum_{l=0}^{2^d-1} a_l \mathbb{L}(\mathbf{s}|l, T, \mathbf{I}) \tag{3}$$

$$\text{where } \mathbb{L}(\mathbf{s}|l, T, \mathbf{I}) = \prod_{j=1}^d (1 - p(l, j)) \mathbb{S}(\mathbf{s}|\mathbf{I}_{i(l,j)}, \mathbf{T}_{i(l,j)}) + p(l, j) (1 - \mathbb{S}(\mathbf{s}|\mathbf{I}_{i(l,j)}, \mathbf{T}_{i(l,j)})) \tag{4}$$

Here, i is the index of the internal node preceding a leaf node l at a certain depth j and p indicates whether the left ($p = 0$) or the right branch ($p = 1$) was taken.

Axis-aligned splitting. Typically, DTs use the Heaviside function for splitting, which is non-differentiable. We use the split function introduced in GradTree to account for reasonable gradients:

$$\mathbb{S}(\mathbf{s}|\boldsymbol{\iota}, \boldsymbol{\tau}) = \lfloor S(\boldsymbol{\iota} \cdot \mathbf{s} - \boldsymbol{\iota} \cdot \boldsymbol{\tau}) \rfloor \tag{5}$$

where $S(z) = \frac{1}{1+e^{-z}}$ represents the logistic function, $\lfloor \cdot \rfloor$ stands for rounding and $\mathbf{a} \cdot \mathbf{b}$ denotes the dot product. We further need to ensure that $\boldsymbol{\iota}$ is a one-hot encoded vector to account for axis-aligned splits. This is achieved by applying a hardmax transformation before calculating \mathbb{S} . Both rounding and hardmax operations are non-differentiable and therefore, SYMPOL is *not* considered as a soft or differentiable DT method. Instead, to overcome non-differentiability, SYMPOL employs a straight-through operator (Bengio et al., 2013) during backpropagation. This allows the model to use non-differentiable operations in the forward pass while ensuring gradient propagation in the backward pass. As a result, we can directly learn an interpretable DT from policy gradient. This makes SYMPOL framework-agnostic and facilitates a seamless integration into existing RL frameworks.

Weight decay. In contrast to GradTree, which employs an Adam (Kingma & Ba, 2014) optimizer with stochastic weight averaging (Izmailov et al., 2018), we opted for an Adam optimizer with weight decay (Loshchilov & Hutter, 2017). In the context of SYMPOL, weight decay does not serve as a regularizer for model complexity, as the interpretation of model parameters differs. We distinguish between three types of parameters: the distributions in the leaves (\mathbf{a}), the split index encoding (\mathbf{I}), and the split values (\mathbf{T}). We do not apply weight decay to the split values because they are independent of magnitude. However, for the split indices and leaves, weight decay enhances exploration during training by penalizing large parameter values. As a result, the distribution of the split index selection and class prediction are narrow and have lower variance. This aids in dynamically adjusting which feature is considered at a split and in altering the predicted leaf distribution.

Actor-critic network architecture. Commonly, the actor and critic use a similar network architecture or even share the same weights (Schulman et al., 2017). While SYMPOL aims for a simple and interpretable policy, we do not have the same requirements for the critic. Therefore, we decided to only employ a tree-based actor and use a full-complexity NN as a value function. As a result, we can still capture complexity through the value function, without losing interpretability, as we maintain a simple and interpretable policy.

Continuous action spaces. Furthermore, we extend the DT policy of SYMPOL to environments with continuous action spaces. Therefore, instead of predicting a categorical distribution over the

270 classes, we predict the mean of a normal distribution at each leaf and utilize an additional variable
 271 $\sigma_{\log} \in \mathbb{R}^{|\mathcal{A}|}$ to learn the log of the standard deviation.
 272

273 4.2 ADDRESSING TRAINING STABILITY 274

275 One main challenge when using DTs as a policy is the stability. While a stable training is also
 276 desired and often hard to achieve for a NN policy, this is even more pronounced for SYMPOL.
 277 This is mainly caused by the inherent tree-based architecture. Changing a split at the top of the
 278 tree can have a severe impact on the whole model, as it can completely change the paths taken for
 279 certain observations. This is especially relevant in the context of RL, where the data distribution can
 280 vary highly between iterations. To mitigate the impact of highly non-stationary training samples,
 281 especially at early stages of training, we made two crucial modifications for improved stability.

282 **Exploration stability.** Motivated by the idea that rollouts of more accurate policies contain in-
 283 creasingly diverse, higher quality samples, we implemented a dynamic number of environment steps
 284 between training iterations. Let us consider a pendulum as an example. While at early stages of
 285 training a relatively small sample size facilitates faster learning as the pendulum constantly flips,
 286 more optimal policies lead to longer rollouts and therefore more expressive and diverse experiences
 287 in the rollout buffer. Similarly, the increasing step counts stabilize the optimization of policy and
 288 critic, as the number of experiences for gradient computation grow with agent expertise and capture
 289 the diversity within trajectories better. Therefore, our novel collection approach starts with n_{init}
 290 environment steps and expands until n_{final} actions are taken before each training iteration. For
 291 computational efficiency reasons, instead of increasing the size of the rollout buffer at every time
 292 step, we introduce a step-wise exponential function. The exponential increase supports exploration
 293 in the initial iterations, while maintaining stability at later iterations. Hence, we define the number
 294 of steps in the environment n_t at time step t as

$$295 n_t = n_{init} \times 2^{\lfloor \frac{(t+1) \times i}{1+t_{total}} \rfloor - 1} \text{ with } i = 1 + \log_2 \left(\frac{n_{init}}{n_{final}} \right) \quad (6)$$

297 For our experiments, we define n_{init} as a hyperparameter (similar to the static step size for other
 298 methods) and set $n_{final} = 128 \times n_{init}$ and therefore $i = 8$ which we observed is a good default value.
 299

300 **Gradient stability.** We also utilize large batch sizes for SYMPOL resulting in less noisy gradients,
 301 leading to a smoother convergence and better stability. In this context, we implement gradient
 302 accumulation to virtually increase the batch size further while maintaining memory-efficiency. As
 303 reduced noise in the gradients also leads to less exploration in the parameter space, we implement a
 304 dynamic batch size, increasing in the same rate as the environment steps between training iterations
 305 (Equation 6). Therefore, we can benefit from exploration and fast convergence early on and increase
 306 gradient stability during the training.
 307

308 5 EVALUATION 309

310 We designed our experiments to evaluate whether SYMPOL can learn accurate DT policies without
 311 information loss and observe whether the trees learned by SYMPOL are small and interpretable.
 312 As mentioned above, we focus on PPO as the most prominent actor-critic, on-policy RL algorithm
 313 in our evaluation. To support our claim that SYMPOL can be seamlessly integrated into existing
 314 on-policy RL frameworks, we additionally provide results using A2C in Appendix A.1.

315 5.1 EXPERIMENTAL SETTINGS 316

317 **Setup.** We implemented SYMPOL in a highly efficient single-file JAX implementation that al-
 318 lows a flawless integration with highly optimized training frameworks (Lange, 2022; Weng et al.,
 319 2022; Bonnet et al., 2024). Our implementation is available in the supplementary material, and
 320 will be made publically available upon acceptance. We evaluated our method on several environ-
 321 ments commonly used for benchmarking RL methods. Specifically, we used control environments
 322 including CartPole (CP), Acrobot (AB), LunarLander (LL), MountainCarContinuous (MC-C) and
 323 Pendulum (PD-C), as well as the MiniGrid (Chevalier-Boisvert et al., 2023) environments Empty-
 Random (E-R), DoorKey (DK), LavaGap (LG) and DistShift (DS).

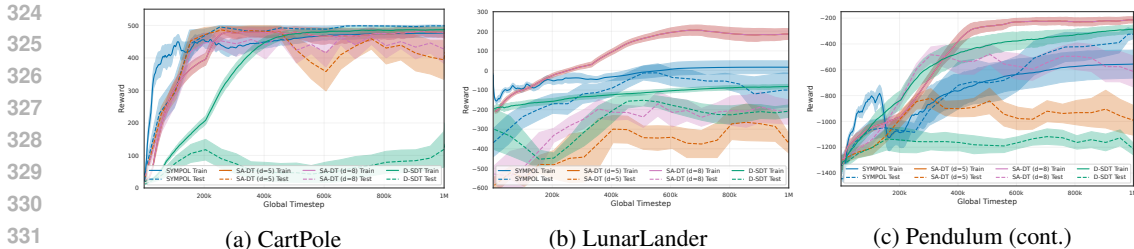


Figure 3: **Selected Training Curves.** Shows the training reward of the full-complexity policy (e.g. MLP in the case of SA-DT) as solid line and the test reward of the interpretable policy as dashed line for three control environments. Additional, more detailed results are in Appendix A.7.

Methods. The goal of this evaluation is to compare SYMPOL to alternative methods that allow an interpretation of RL policies as a symbolic, axis-aligned DTs. Therefore, we build on previous work (Silva et al., 2020) and use two methods grounded in the interpretable RL literature, as follows:

- **State-Action DTs (SA-DT):** SA-DTs are the most common method to generate interpretable policies post-hoc. Hereby, we first train an MLP policy, which is then distilled into a DT as a post-processing step after the training. SA-DT can be considered as a version of DAGGER (Ross et al., 2011) and therefore a simplified version of VIPER (Bastani et al., 2018). In a comparative experiment (see Appendix A.3), we showed that for the case of learning small, interpretable DTs the performance of SA-DT is similar to those of VIPER, which is in-line with results reported e.g. by Kohler et al. (2024).
- **Discretized Soft DTs (D-SDT):** SDTs allow gradient computation by assigning probabilities to each node. While SDTs exhibit a hierarchical structure, they are usually considered as less interpretable, since multiple features are considered in a single split and the whole tree is traversed simultaneously (Marton et al., 2024a). Therefore, Silva et al. (2020) use SDTs as policies which are discretized post-hoc to allow an easy interpretation.

We further included an MLP and SDT, providing an orientation to state-of-the-art results.

Evaluation procedure. We report the average undiscounted cumulative reward over 5 random trainings with 5 random evaluation episodes each (=25 evaluations for each method). We trained each method for 1mio timesteps. For SYMPOL, SDT and MLP, we optimized the hyperparameters based on the validation reward with optuna (Akiba et al., 2019) for 60 trials using a predefined grid. For D-SDT we discretized the SDT and for SA-DT, we distilled the MLP with the highest performance. More details on the hyperparameters can be found in Appendix C.

5.2 RESULTS

SYMPOL does not exhibit information loss. Existing methods for learning DT policies usually involve post-processing to obtain the interpretable model. Therefore, they introduce a mismatch between the optimized and interpreted policy, which can result in information loss. The main advantage of SYMPOL is the direct optimization of a DT policy, which guarantees that there is no information loss between the optimized and interpreted policy. To show this, we calculated Cohens’s D to measure the effect size comparing the validation reward of the trained model with the test reward of the applied, optionally post-processed model (Table 1). We can observe very large effects for SA-DT and D-SDT and only a very small effect for SYMPOL, similar to full-complexity models MLP and SDT. This discrepancy can also be observed in the training curves in Figure 3.

Table 1: Information Loss. We calculated Cohens’s D to measure effect size between the validation reward of the trained and the test reward of the applied model. Values > 0.8 are considered as a large effect. Detailed results are in Appendix A.2

	Cohen’s D ↓
SYMPOL (ours)	-0.019
SA-DT (d=5)	3.449
SA-DT (d=8)	2.527
D-SDT	3.126
MLP	0.306
SDT	0.040

SYMPOL learns accurate DT policies. We evaluated our approach against existing methods on control environments in Table 2. SYMPOL is consistently among the best interpretable models and achieves significantly higher rewards compared to alternative methods for learning DT policies on several environments, especially on LL and PD-C. Further, SYMPOL consistently solves CP and AB and is competitive to full-complexity models on most environments.

DT policies offer a good inductive bias for categorical environments.

While SYMPOL achieves great results in control benchmarks, it may not be an ideal method for environments modeling physical relationships. As recently also noted by Fuhrer et al. (2024), tree-based models are best suited for categorical environments due to their effective use of axis-aligned splits. In our experiments on MiniGrid (Table 3), SYMPOL achieves comparable or superior results to full-complexity models (e.g. on LG-7). The performance gap between SA-DT and SYMPOL is smaller in certain MiniGrid environments due to less complex environment transition functions and missing randomness, making the distillation easy. Considering more complex environments with randomness or lava like E-R or LG-7, SYMPOL outperforms alternative methods by a substantial margin.

DT policies learned with SYMPOL are small and interpretable.

While we trained SYMPOL with a depth of 7 and therefore 255 possible nodes, the effective tree size after pruning is significantly smaller with only 50.5 nodes (internal and leaf combined) on average. This can be attributed to a self-pruning mechanism that is inherently applied by SYMPOL in learning redundant paths during the training and therefore only optimizing relevant parts. Furthermore, DTs learned with SYMPOL are smaller than SA-DTs (d=5) with an average of 60.3 nodes and significantly smaller than SA-DTs (d=8) averaging 291.6 nodes. The pruned D-SDTs are significantly smaller with only 16.5, but also have a very poor performance, as shown in the previous experiment. An exemplary DT learned by SYMPOL is visualized in Figure 4. Extended results, including a comparison with SDTs are in Appendix A.4 and A.5.

SYMPOL is efficient. In RL, the actor-environment interaction frequently constitutes a significant portion of the total runtime. For smaller policies, in particular, the runtime is mainly determined by the time required to execute actions within the environment to obtain the next observation, while the time required to execute the policy itself having a comparatively minimal impact on runtime. Therefore, recent research put much effort into optimizing this interaction through environment

Table 2: **Control Performance.** We report the average undiscounted cumulative test reward over 25 random trials. The best interpretable method, and methods not statistically different, are marked bold.

	CP	AB	LL	MC-C	PD-C
SYMPOL (ours)	500	- 80	- 57	94	- 323
D-SDT	128	-205	-221	-10	-1343
SA-DT (d=5)	446	-97	-197	97	-1251
SA-DT (d=8)	476	- 75	-150	96	- 854
MLP	500	- 72	241	95	- 191
SDT	500	- 77	-124	- 4	- 310

Table 3: **MiniGrid Performance.** We report the average undiscounted cumulative test reward over 25 random trials. The best interpretable method, and methods not statistically different, are marked bold.

	E-R	DK	LG-5	LG-7	DS
SYMPOL (ours)	0.964	0.959	0.951	0.953	0.939
D-SDT	0.662	0.654	0.262	0.381	0.932
SA-DT (d=5)	0.583	0.958	0.951	0.458	0.952
SA-DT (d=8)	0.845	0.961	0.951	0.799	0.954
MLP	0.963	0.963	0.951	0.760	0.951
SDT	0.966	0.959	0.839	0.953	0.954

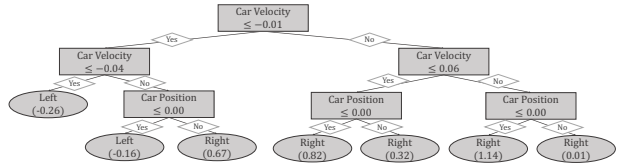


Figure 4: **SYMPOL Policy for MC-C.** The main rule encoded by this DT is that the car should accelerate to the left, if its velocity is negative and to the right if it is positive. This essentially increases the speed of the car over time, making it possible to reach the goal at the top of the hill. The magnitude of acceleration is mainly determined by the current position, reducing the action cost.

432
433
434
435
436
437
438
439
440
441
442

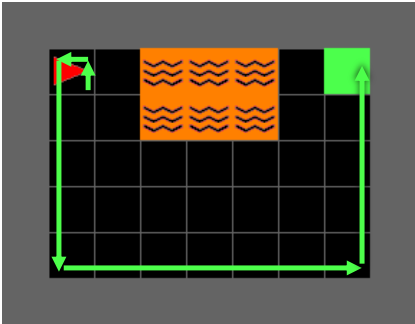


Figure 6: **DistShift**. We show the training environment for the agent along with the starting position and goal. The path taken by SYMPOL (see Figure 7) is marked by green arrows and solves the environment.

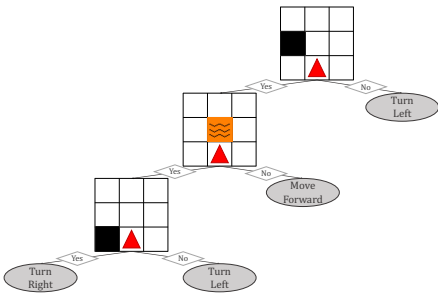


Figure 7: **SYMPOL Policy**. This image shows the DT policy of SYMPOL. Split nodes are visualized as the 3x3 view grid of the agent with one square marking the considered object and position. If the visualized object is present at this position, the true path (left) is taken.

443
444
445
446
447
448
449

vectorization. The design of SYMPOL, in contrast to existing methods for tree-based RL, allows a seamless integration with these highly efficient training frameworks. As a result, the runtime of SYMPOL is almost identical to using an MLP or SDT as policy, averaging less than 30 seconds for 1mio timesteps. Detailed results are in Appendix A.7.

450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465

Ablation study. In Section 4, we introduced several crucial components to facilitate an efficient and stable training of SYMPOL. To support the intuitive justifications for our modifications, we performed an ablation study (Figure 5) to evaluate the relevance of the individual components. Our results confirm that each component substantially contributes to the overall performance.

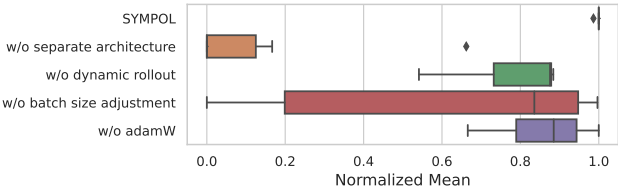


Figure 5: **Ablation Study**. We report the mean normalized reward over all control environments (details in Table 11).

6 CASE STUDY: DETECTING GOAL MISGENERALIZATION

466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

To demonstrate the benefits of SYMPOLs enhanced transparency, we present a case study on goal misgeneralization (Di Langosco et al., 2022). Good policy generalization is vital in RL, yet agents often exhibit poor out-of-distribution performance, even with minor environmental changes. Goal misgeneralization is a well-researched out-of-distribution robustness failure that occurs when an agent learns robust skills during training but follows unintended goals. This happens when the agent’s behavioral objective diverges from the intended objective, leading to high rewards during training but poor generalization during testing. For instance, NNs were shown to systematically misgeneralize on Atari environments (Farebrother et al., 2018; Delfosse et al., 2024a).

To demonstrate that SYMPOL can help in detecting misaligned behavior, let us consider the DistShift environment from MiniGrid, shown in Figure 6. The environment is designed to test for misgeneralization (Chevalier-Boisvert et al., 2023), as the goal is repeatedly placed in the top right corner and the lava remains at the same position. We can formulate the intended behavior according to the task description as avoiding the lava and reaching a specific goal location. SYMPOL, similar to other methods, solved the task consistently. The advantage of SYMPOL is the tree-based structure, which is easily interpretable. When inspecting the SYMPOL policy (Figure 7), we can immediately observe that the agent has not captured the actual task correctly. Essentially, it has only learned to keep an empty space on the left of the agent (which translates into following the wall) and not to step into lava (but not to get around it). While this is sufficient to solve this exact environment, it is evident, that the agent has not generalized to the overall goal.

REFERENCES

- 540
541
542 Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:
543 A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM*
544 *SIGKDD international conference on knowledge discovery & data mining*, pp. 2623–2631, 2019.
- 545 Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy
546 extraction. *Advances in Neural Information Processing Systems*, 31, 2018.
- 547 Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients
548 through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- 549
550 Clément Bonnet, Daniel Luo, Donal Byrne, Shikha Surana, Sasha Abramowitz, Paul Duckworth,
551 Vincent Coyette, Laurence I. Midgley, Elshadai Tegegn, Tristan Kalloniatis, Omayma Mahjoub,
552 Matthew Macfarlane, Andries P. Smit, Nathan Grinsztajn, Raphael Boige, Cemlyn N. Waters,
553 Mohamed A. Mimouni, Ulrich A. Mbou Sob, Ruan de Kock, Siddarth Singh, Daniel Furelos-
554 Blanco, Victor Le, Arnú Pretorius, and Alexandre Laterre. Jumanji: a diverse suite of scalable
555 reinforcement learning environments in jax, 2024. URL [https://arxiv.org/abs/2306.](https://arxiv.org/abs/2306.09884)
556 09884.
- 557 Yushi Cao, Zhiming Li, Tianpei Yang, Hao Zhang, Yan Zheng, Yi Li, Jianye Hao, and Yang Liu.
558 Galois: boosting deep reinforcement learning via generalizable logic synthesis. *Advances in*
559 *Neural Information Processing Systems*, 35:19930–19943, 2022.
- 560 Chun-Hao Chang, Rich Caruana, and Anna Goldenberg. Node-gam: Neural generalized additive
561 model for interpretable deep learning. *arXiv preprint arXiv:2106.01613*, 2021.
- 562
563 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems,
564 Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld:
565 Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*,
566 abs/2306.13831, 2023.
- 567
568 Youri Coppens, Kyriakos Efthymiadis, Tom Lenaerts, Ann Nowé, Tim Miller, Rosina Weber, and
569 Daniele Magazzeni. Distilling deep reinforcement learning policies in soft decision trees. In
570 *Proceedings of the IJCAI 2019 workshop on explainable artificial intelligence*, pp. 1–6, 2019.
- 571 Vinícius G Costa, Jorge Pérez-Aracil, Sancho Salcedo-Sanz, and Carlos E Pedreira. Evolving inter-
572 pretable decision trees for reinforcement learning. *Artificial Intelligence*, 327:104057, 2024.
- 573
574 Leonardo L Custode and Giovanni Iacca. Evolutionary learning of interpretable decision trees. *IEEE*
575 *Access*, 11:6169–6184, 2023.
- 576
577 Quentin Delfosse, Jannis Blüml, Bjarne Gregori, and Kristian Kersting. Hacktari: Atari learning
578 environments for robust and continual reinforcement learning. *arXiv preprint arXiv:2406.03997*,
579 2024a.
- 580
581 Quentin Delfosse, Hikaru Shindo, Devendra Dhama, and Kristian Kersting. Interpretable and ex-
582 plainable logical policies via neurally guided symbolic abstraction. *Advances in Neural Informa-*
tion Processing Systems, 36, 2024b.
- 583
584 Quentin Delfosse, Sebastian Sztwiertnia, Mark Rothermel, Wolfgang Stammer, and Kristian Ker-
585 sting. Interpretable concept bottlenecks to align reinforcement learning agents. *arXiv preprint*
586 *arXiv:2401.05821*, 2024c.
- 587
588 Lauro Langosco Di Langosco, Jack Koch, Lee D Sharkey, Jacob Pfau, and David Krueger. Goal mis-
589 generalization in deep reinforcement learning. In *International Conference on Machine Learning*,
pp. 12004–12019. PMLR, 2022.
- 590
591 Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning.
592 *Journal of Machine Learning Research*, 6:503–556, 2005.
- 593
Jesse Farebrother, Marlos C Machado, and Michael Bowling. Generalization and regularization in
dqn. *arXiv preprint arXiv:1810.00123*, 2018.

- 594 Gregory Farquhar, Tim Rocktäschel, Maximilian Igl, and Shimon Whiteson. Treeqn and
595 atreec: Differentiable tree-structured models for deep reinforcement learning. *arXiv preprint*
596 *arXiv:1710.11417*, 2017.
- 597 Benjamin Fuhrer, Chen Tessler, and Gal Dalal. Gradient boosting reinforcement learning. *arXiv*
598 *preprint arXiv:2407.08250*, 2024.
- 600 Jiaming Guo, Rui Zhang, Shaohui Peng, Qi Yi, Xing Hu, Ruizhi Chen, Zidong Du, Ling Li, Qi Guo,
601 Yunji Chen, et al. Efficient symbolic policy learning with differentiable symbolic expression.
602 *Advances in Neural Information Processing Systems*, 36, 2024.
- 603 Ujjwal Das Gupta, Erik Talvitie, and Michael Bowling. Policy tree: Adaptive representation for
604 policy gradient. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29,
605 2015.
- 607 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
608 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*
609 *ence on machine learning*, pp. 1861–1870. PMLR, 2018.
- 611 Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wil-
612 son. Averaging weights leads to wider optima and better generalization. *arXiv preprint*
613 *arXiv:1803.05407*, 2018.
- 614 Aman Jhunjunwala, Jaeyoung Lee, Sean Sedwards, Vahdat Abdelzad, and Krzysztof Czarnecki.
615 Improved policy extraction via online q-value distillation. In *2020 International Joint Conference*
616 *on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2020.
- 618 Zhengyao Jiang and Shan Luo. Neural logic reinforcement learning. In *International conference on*
619 *machine learning*, pp. 3110–3119. PMLR, 2019.
- 621 Akansha Kalra and Daniel S Brown. Interpretable reward learning via differentiable decision trees.
622 In *NeurIPS ML Safety Workshop*, 2022.
- 623 Akansha Kalra and Daniel S Brown. Can differentiable decision trees learn interpretable reward
624 functions? *arXiv preprint arXiv:2306.13004*, 2023.
- 626 Kentaro Kanamori, Takuya Takagi, Ken Kobayashi, and Yuichi Ike. Counterfactual explanation
627 trees: Transparent and consistent actionable recourse with decision trees. In *International Con-*
628 *ference on Artificial Intelligence and Statistics*, pp. 1846–1870. PMLR, 2022.
- 629 Eoin M. Kenny, Mycal Tucker, and Julie Shah. Towards interpretable deep reinforcement learn-
630 ing with human-friendly prototypes. In *International Conference on Learning Representations*
631 *(ICLR)*, 2023. URL https://openreview.net/forum?id=hWwY_Jq0xsN.
- 632 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
633 *arXiv:1412.6980*, 2014.
- 634 Hector Kohler, Quentin Delfosse, Riad Akrouf, Kristian Kersting, and Philippe Preux. Inter-
635 pretable and editable programmatic tree policies for reinforcement learning. *arXiv preprint*
636 *arXiv:2405.14956*, 2024.
- 637 Mikel Landajuela, Brenden K Petersen, Sookyung Kim, Claudio P Santiago, Ruben Glatt, T Nathan
638 Mundhenk, Jacob F Pettit, and Daniel M Faissol. Discovering symbolic policies with deep rein-
639 forcement learning. In *Proceedings of the 38th International Conference on Machine Learning*,
640 pp. 5979–5989. PMLR, 2021.
- 641 Robert Tjarko Lange. gymnax: A JAX-based reinforcement learning environment library, 2022.
642 URL <http://github.com/RobertTLange/gymnax>.
- 643 Zhao-Hua Li, Yang Yu, Yingfeng Chen, Ke Chen, Zhipeng Hu, and Changjie Fan. Neural-to-tree
644 policy distillation with policy improvement criterion. *arXiv preprint arXiv:2108.06898*, 2021.

- 648 Guiliang Liu, Oliver Schulte, Wang Zhu, and Qingcan Li. Toward interpretable deep reinforce-
649 ment learning with linear model u-trees. In *Machine Learning and Knowledge Discovery in*
650 *Databases: European Conference, ECML PKDD 2018, Dublin, Ireland, September 10–14, 2018,*
651 *Proceedings, Part II 18*, pp. 414–429. Springer, 2019.
- 652 Xiao Liu, Shuyang Liu, Bo An, Yang Gao, Shangdong Yang, and Wenbin Li. Effective interpretable
653 policy distillation via critical experience point identification. *IEEE Intelligent Systems*, 38(5):
654 28–36, 2023.
- 655 Zichuan Liu, Yuanyang Zhu, Zhi Wang, Yang Gao, and Chunlin Chen. Mixrts: Toward inter-
656 interpretable multi-agent reinforcement learning via mixing recurrent soft decision trees. *arXiv*
657 *preprint arXiv:2209.07225*, 2022.
- 658 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
659 *arXiv:1711.05101*, 2017.
- 660 Lirui Luo, Guoxi Zhang, Hongming Xu, Yaodong Yang, Cong Fang, and Qing Li. End-to-end
661 neuro-symbolic reinforcement learning with textual explanations. In *Forty-first International*
662 *Conference on Machine Learning*.
- 663 Lirui Luo et al. End-to-end neuro-symbolic reinforcement learning with textual explanations. *arXiv*
664 *preprint arXiv:2403.12451*, 2024.
- 665 Sascha Marton, Stefan Lüdtkke, Christian Bartelt, and Heiner Stuckenschmidt. Gradtree: Learning
666 axis-aligned decision trees with gradient descent. In *Proceedings of the AAAI Conference on*
667 *Artificial Intelligence*, volume 38, pp. 14323–14331, 2024a.
- 668 Sascha Marton, Stefan Lüdtkke, Christian Bartelt, and Heiner Stuckenschmidt. Grande: Gradient-
669 based decision tree ensembles for tabular data. In *The Twelfth International Conference on Learn-*
670 *ing Representations*, 2024b.
- 671 Stephanie Milani, Zhicheng Zhang, Nicholay Topin, Zheyuan Ryan Shi, Charles Kamhoua, Evan-
672 gelos E Papalexakis, and Fei Fang. Maviper: Learning decision tree policies for interpretable
673 multi-agent reinforcement learning. In *Joint European Conference on Machine Learning and*
674 *Knowledge Discovery in Databases*, pp. 251–266. Springer, 2022.
- 675 Joosung Min and Lloyd T Elliott. Q-learning with online random forests. *arXiv preprint*
676 *arXiv:2204.03771*, 2022.
- 677 Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim
678 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement
679 learning, 2016. URL <https://arxiv.org/abs/1602.01783>.
- 680 Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:
681 Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 682 Svetlin Penkov and Subramanian Ramamoorthy. Learning programmatically structured representa-
683 tions with perceptor gradients. *arXiv preprint arXiv:1905.00956*, 2019.
- 684 Ben Peters, Vlad Niculae, and André FT Martins. Sparse sequence-to-sequence models. *arXiv*
685 *preprint arXiv:1905.05702*, 2019.
- 686 Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision ensembles for
687 deep learning on tabular data. *arXiv preprint arXiv:1909.06312*, 2019.
- 688 Stéphane Ross, Geoffrey Gordon, and J Andrew Bagnell. A reduction of imitation learning and
689 structured prediction to no-regret online learning. In *Proceedings of the 14th International Con-*
690 *ference on Artificial Intelligence and Statistics (AISTATS)*, pp. 627–635. JMLR Workshop and
691 Conference Proceedings, 2011.
- 692 Aaron M Roth, Nicholay Topin, Pooyan Jamshidi, and Manuela Veloso. Conservative q-
693 improvement: Reinforcement learning for an interpretable decision-tree policy. *arXiv preprint*
694 *arXiv:1907.01180*, 2019.

- 702 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
703 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
704
- 705 Andrew Silva and Matthew Gombolay. Encoding human domain knowledge to warm start rein-
706 forcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35,
707 pp. 5042–5050, 2021.
708
- 709 Andrew Silva, Matthew Gombolay, Taylor Killian, Ivan Jimenez, and Sung-Hyun Son. Optimiza-
710 tion methods for interpretable differentiable decision trees applied to reinforcement learning. In
711 *International conference on artificial intelligence and statistics*, pp. 1855–1865. PMLR, 2020.
712
- 713 Pradyumna Tambwekar, Andrew Silva, Nakul Gopalan, and Matthew Gombolay. Natural lan-
714 guage specification of reinforcement learning policies through differentiable decision trees. *IEEE*
715 *Robotics and Automation Letters*, 8(6):3621–3628, 2023.
716
- 717 Nicholay Topin, Stephanie Milani, Fei Fang, and Manuela Veloso. Iterative bounding mdps: Learn-
718 ing interpretable policies via non-interpretable methods. In *Proceedings of the AAAI Conference*
719 *on Artificial Intelligence*, volume 35, pp. 9923–9931, 2021.
720
- 721 Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu,
722 Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard
723 interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
724
- 725 Dweep Trivedi, Jesse Zhang, Shao-Hua Sun, and Joseph J Lim. Learning to synthesize programs
726 as interpretable and generalizable policies. *Advances in Neural Information Processing Systems*,
727 34:25146–25163, 2021.
728
- 729 Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri.
730 Programmatically interpretable reinforcement learning. In *International Conference on Machine*
731 *Learning*, pp. 5045–5054. PMLR, 2018.
732
- 733 Daniël Vos and Sicco Verwer. Optimal decision tree policies for markov decision processes. *arXiv*
734 *preprint arXiv:2301.13185*, 2023.
735
- 736 Daniël Vos and Sicco Verwer. Optimizing interpretable decision tree policies for reinforcement
737 learning. *arXiv preprint arXiv:2408.11632*, 2024. URL <https://arxiv.org/abs/2408.11632>.
738
- 739 Maxime Wabartha and Joelle Pineau. Piecewise linear parametrization of policies: Towards in-
740 terpretable deep reinforcement learning. In *The Twelfth International Conference on Learning*
741 *Representations*, 2024. URL <https://openreview.net/forum?id=hOMVq57Ce0>.
742
- 743 Jiayi Weng, Min Lin, Shengyi Huang, Bo Liu, Denys Makoviichuk, Viktor Makoviychuk,
744 Zichen Liu, Yufan Song, Ting Luo, Yukun Jiang, Zhongwen Xu, and Shuicheng Yan. En-
745 vPool: A highly parallel reinforcement learning environment execution engine. In S. Koyejo,
746 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural In-*
747 *formation Processing Systems*, volume 35, pp. 22409–22421. Curran Associates, Inc., 2022.
748 URL https://proceedings.neurips.cc/paper_files/paper/2022/file/8caaf08e49ddb6694fae067442ee21-Paper-Datasets_and_Benchmarks.pdf.
749
- 750 Kai Xu et al. Symbolic-model-based reinforcement learning. *arXiv preprint arXiv:2203.12346*,
751 2022.
752
- 753 Zhuorui Ye, Stephanie Milani, Fei Fang, and Geoff Gordon. Concept-based interpretable rein-
754 forcement learning with limited to no human labels. In *Workshop on Interpretable Policies in*
755 *Reinforcement Learning@ RLC-2024*, 2024.

A ADDITIONAL RESULTS

In this section, we present additional results to support the claims made in the main paper, along with extended results for the summarizing tables. We focus on the control environments because they offer a diverse suite of benchmarks that cover different tasks and include both continuous and discrete action spaces. We chose not to include the MiniGrid environments here because their inclusion could distort the results, particularly the averages calculated in the main paper, as all MiniGrid environments involve similar tasks and feature a discrete action and observation space. The primary reason for including MiniGrid in the main paper is to provide additional experimental results that confirm the robustness and applicability of our method across different domains, as well as to highlight that tree-based methods offer a beneficial inductive bias for these categorical environments.

A.1 EVALUATION WITH ALTERNATIVE RL ALGORITHMS: ADVANTAGE ACTOR CRITIC (A2C)

To support our claim that SYMPOL can be integrated into arbitrary on-policy frameworks, we provide results on A2C in the following. The reported results use optimized hyperparameters for each method. In general, A2C is considered as less stable compared to PPO, which is an additional challenge for SYMPOL, as training stability is especially crucial for DT policies. As A2C does not update the policy in minibatches over multiple epochs, we did not include a dynamic batch size here, but update SYMPOL with a single update over the rollout to stay consistent with the A2C algorithm.

Table 4: **A2C Control Performance.** We report the average undiscounted cumulative test reward over 25 random trials. The best interpretable method, and methods not statistically different, are marked bold.

	CP	AB	LL	MC-C	PD-C
SYMPOL (ours)	500	- 84	- 85	58	- 502
D-SDT	11	-427	-396	-0	-1137
SA-DT (d=5)	295	-102	-348	0	-1467
SA-DT (d=8)	223	- 99	-367	2	-1526
MLP	500	- 78	208	0	- 202
SDT	500	- 85	-159	0	- 201

Table 5: **A2C Control Performance Comparison.** We report the average undiscounted cumulative test reward over 25 random trials, comparing A2C with PPO using optimized hyperparameters. A number is marked bold if the performance achieved with the underlying RL algorithm (PPO or A2C) is significantly better or not statistically different from the best result.

	MLP		SDT		SYMPOL (ours)		SA-DT (d=5)		SA-DT (d=8)		D-SDT	
	A2C	PPO	A2C	PPO	A2C	PPO	A2C	PPO	A2C	PPO	A2C	PPO
CP	500	500	500	500	500	500	295	446	223	476	11	128
AB	-78	-72	-85	-77	-84	-80	-102	-97	-99	-75	-427	-205
LL	208	241	-159	-124	-85	-57	-348	-197	-367	-150	-396	-221
MC-C	0	95	0	-4	58	94	0	97	-2	96	0	-10
PD-C	-202	-191	-201	-310	-502	-323	-1467	-1251	-1526	-854	-1137	-1343

We compared the performance of all methods using A2C on control environments in Table 4, and additionally provided a direct comparison between PPO and A2C in Table 5. When using A2C, SYMPOL consistently outperforms other interpretable models. The performance gap becomes even more pronounced when using A2C instead of PPO, as SYMPOL achieves substantially higher performance than all other interpretable models in each environment. On MC-C, SYMPOL is the only method that achieves a positive reward, whereas even full-complexity models were unable to solve the task. This can be attributed to the lower training stability of A2C compared to PPO. This could also explain the poor results of distillation methods, as the policy learned by full-complexity models, even when achieving a high test reward, is potentially less consolidated, making it harder to distill.

810 However, to confirm this assumption, further experiments would be required. Based on these results,
 811 we can confirm that SYMPOL can seamless be integrated into other RL algorithms, demonstrating
 812 the high flexibility of our proposed method. Additionally, our method can benefit from advances in
 813 RL, as it can be seamlessly integrated into novel frameworks.

815 A.2 INFORMATION LOSS

816 We provide detailed results on the information loss which can result as a consequence of discretiza-
 817 tion (for D-SDT) or distillation (for SA-DT). In Table 6, we report the validation reward of the
 818 trained model along with the test reward of the discretized model. We can clearly observe that there
 819 are major differences for SA-DT and D-SDT on several datasets, indicating information loss. In
 820 Table 7, we report Cohen’s D to measure the effect size comparing the validation reward of the
 821 trained model with the test reward of the applied, optionally post-processed model. Again, we can
 822 clearly see large effects for SA-DT and D-SDT on several datasets, especially for PD-C and LL, but
 823 also CP. Furthermore, the training curves in Figure 14 visually show the information loss during the
 824 training.

826 Table 6: **Information Loss (Comparison)**. We report the validation reward of the trained model
 827 and the test reward of the applied model.

	MLP		SDT		SYMPOL (ours)		SA-DT (d=5)		SA-DT (d=8)		D-SDT	
	valid	test	valid	test	valid	test	valid	test	valid	test	valid	test
831 CP	500	500	500	500	500	500	500	446	500	476	500	128
832 AB	-71	-72	-89	-77	-79	-80	-71	-97	-71	-75	-89	-205
833 LL	256	241	-91	-124	-9	-57	256	-197	256	-150	-91	-221
834 MC-C	95	95	-4	-4	87	94	95	97	95	96	-4	-10
835 PD-C	-169	-191	-295	-310	-305	-323	-169	-1251	-169	-854	-295	-1343

837 Table 7: **Information Loss (Cohen’s D)**. We calculated Cohens’s D to measure effect size between
 838 the validation reward of the trained model and the test reward of the applied model. Typically,
 839 values > 0.5 are considered a medium and values > 0.8 a large effect. positive effects that are at
 840 least medium are marked as bold.

	MLP	SDT	SYMPOL (ours)	SA-DT (d=5)	SA-DT (d=8)	D-SDT
842 CP	0.000	0.000	0.000	0.632	1.214	4.075
843 AB	0.035	-0.630	0.104	0.728	0.338	0.982
844 LL	0.341	0.750	0.370	4.776	8.155	2.254
845 MC-C	-0.042	-0.002	-1.035	-2.011	-1.172	0.745
846 PD-C	1.195	0.081	0.468	13.120	4.101	7.573
847 Mean ↓	0.306	0.040	-0.019	3.449	2.527	3.126

851 A.3 COMPARISON OF SYMPOL, SA-DT (DAGGER) AND VIPER (Q-DAGGER)

852 In this section, we provide a direct comparison of SYMPOL with SA-DT and VIPER for control
 853 environments (see Table 8) and MiniGrid environments (see Table 9). SA-DT (Silva et al., 2020)
 854 can be considered as a version of DAGGER (Ross et al., 2011) and is conceptually similar to VIPER
 855 (Q-DAGGER) (Bastani et al., 2018) which improves data collection including additional weighting.
 856 Our results remain consistent with our original claims, demonstrating that SYMPOL outperforms
 857 alternative approaches. This is also in-line with the results reported by Kohler et al. (2024), where
 858 the authors show, that the sampling in VIPER does not yield a better performance compared do
 859 DAGGER/SA-DT for interpretable DTs. The results reported in the original VIPER paper (Bastani
 860 et al., 2018) stating to achieve a perfect reward for CP are on a different version of the environment
 861 (CartPole-v0) with only 200 opposed to 500 time steps and less randomness (CartPole-v1), making
 862 the underlying task easier. Also, we want to note, that the reported results are in-line with related
 863 work, reporting comparable or worse results than ours. For instance, Vos & Verwer (2024) report
 a mean reward of only 367 for VIPER on CartPole-v1. Also, Kenny et al. (2023) showcase a poor

Table 8: **Control Performance.** We report the average undiscounted cumulative test reward over 25 random trials. The best interpretable method, and methods not statistically different, are marked bold. Please note that VIPER cannot be applied to continuous environments.

	CP	AB	LL	MC-C	PD-C
SYMPOL (ours)	500	- 80	- 57	94	- 323
D-SDT	128	-205	-221	-10	-1343
SA-DT (d=5)	446	-97	-197	97	-1251
SA-DT (d=8)	476	- 75	-150	96	- 854
VIPER (d=5)	457	- 77	-200	-	-
VIPER (d=8)	480	- 75	-169	-	-
MLP	500	- 72	241	95	- 191
SDT	500	- 77	-124	- 4	- 310

Table 9: **MiniGrid Performance.** We report the average undiscounted cumulative test reward over 25 random trials. The best interpretable method, and methods not statistically different, are marked bold.

	E-R	DK	LG-5	LG-7	DS
SYMPOL (ours)	0.964	0.959	0.951	0.953	0.939
D-SDT	0.662	0.654	0.262	0.381	0.932
SA-DT (d=5)	0.583	0.958	0.951	0.458	0.952
SA-DT (d=8)	0.845	0.961	0.951	0.799	0.954
VIPER (d=5)	0.651	0.958	0.948	0.456	0.954
VIPER (d=8)	0.845	0.963	0.948	0.801	0.954
MLP	0.963	0.963	0.951	0.760	0.951
SDT	0.966	0.959	0.839	0.953	0.954

performance of VIPER in general and specifically for LL the performance is worse than what we reported. Our findings align with these, suggesting that differences in performance may reflect randomness and missing generalizability in the evaluation.

A.4 TREE SIZE

We report the average tree sizes over 25 trials for each environment. The DTs for SYMPOL and D-SDT are automatically pruned by removing redundant paths. There are mainly two identifiers, making a path redundant:

- The split threshold of a split is outside the range specified by the environment. For instance, if $x_1 \in [0.0, 1.0]$ the decision $x_1 \leq -0.1$ will always be false as $-0.1 \leq 0.0$.
- A decision at a higher level of the tree already predefines the current decision. For instance, if the split at the root node is $x_1 \leq 0.5$ and the subsequent node following the true path is $x_1 \leq 0.6$ we know that this node will always be evaluated to true as $0.5 \leq 0.6$.

We excluded the MiniGrid environments here, as they require a more sophisticated, automated pruning as there exist more requirements making a path redundant. For instance, if for the decision whether there is a wall in front of the agent is true, the decision for all other objects at the same position has to be always false.

Table 10: **Tree Size.** We report the average size of the learned DT for each environment.

	SYMPOL (ours)	D-SDT	SA-DT (d=5)	SA-DT (d=8)
CP	39.4	14.2	61.8	315.0
AB	78.6	17.0	56.5	173.0
LL	55.0	19.8	59.8	270.2
MC-C	23.4	3.0	61.0	311.8
PD-C	56.2	28.6	62.2	388.2
Mean ↓	50.5	16.5	60.3	291.6

A.5 INTERPRETABILITY COMPARISON BETWEEN HARD, AXIS-ALIGNED AND SOFT DTs

In the main paper, we showed a visualization of a hard, axis-aligned DT learned by SYMPOL on the MC-C environment. While this was a comparatively small tree, we provide another example of a comparatively large tree with 59 nodes in Figure 10. While the tree is comparatively large, we can observe that the main logic is contained in the nodes at higher levels, focusing on the pole angle and the pole angular velocity. The less important features are in the lower levels where splits are often made on the cart position, which is not required to solve the task perfectly. This also highlights the potential for advanced post-hoc pruning methods to increase interpretability and potentially even generalization. When comparing the hard, axis-aligned decision trees (DTs) (Figure 10 and Figure 4) learned by SYMPOL with corresponding soft decision trees (SDTs) learned by existing direct optimization methods (Figure 11 and Figure 12), the superiority of axis-aligned splits over oblique, multidimensional boundaries becomes evident.

The DTs learned by SYMPOL are substantially more interpretable, both at the level of individual splits and the overall tree structure. For example, when examining the root node of the CartPole task, the standard DT (Figure 10) makes a straightforward comparison, " $s_3(\text{Pole Velocity}) \leq -0.800$ ", whereas the corresponding SDT (Figure 11) expresses the decision as " $\sigma(-1.14s_0 - 0.30s_1 + 0.94s_2 + 0.11s_3 + 0.99)$ ", which is significantly more complex and harder to interpret. This disparity becomes even more pronounced when considering complete paths or the entire tree. In SYMPOL's DTs, decisions at nodes are binary (yes/no), whereas SDTs employ probabilistic routing.

Probabilistic routing introduces two key disadvantages in interpretability compared to axis-aligned DTs: (1) the need to consider multiple paths simultaneously, and (2) the inability to directly interpret the leaf outputs, as they are weighted by path probabilities. We also want to note that in order to allow a visualization of the SDT, we reduced the tree size to only 4, while the tree size used during the evaluation was 7.

A.6 ABLATION STUDY

Our ablation study was designed to support our intuitive justifications for the modifications made to the RL framework and our method. Therefore, we disabled individual component of our method and evaluated the performance without the specific component. This includes the following modifications introduced in Section 4:

- w/o separate architecture:** Instead of using separate architectures for actor and critic, we use the same architecture and hyperparameters for the actor and critic.
- w/o dynamic rollout:** We proposed a dynamic rollout buffer that increases with a stepwise exponential rate during training to increase stability while maintaining exploration early on. Here we used a standard, static rollout buffer.
- w/o batch size adjustment:** Similar to the dynamic rollout buffer, we proposed using a dynamic batch size to increase gradient stability in later stages of the training. Here, we used standard, static batch size.
- w/o adamW:** We introduced an Adam optimizer with weight decay to SYMPOL to support the adjustment of the features to split on and the class predicted. Here, we use a standard Adam optimizer without weight decay.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

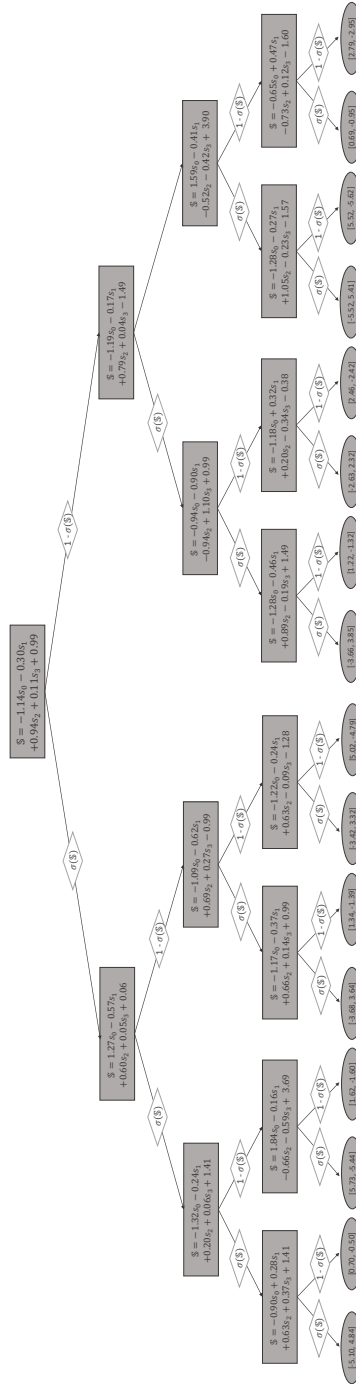


Figure 11: **Exemplary SDT for CartPole.** This figure visualizes a soft / differentiable DT learned on CartPole. The tree involves oblique decisions involving multiple features at each split. Additionally, there is no hard decision on which path is selected, but multiple paths are taken with an associated probability. The final prediction is obtained by weighting the leaf outputs with the corresponding path probability.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

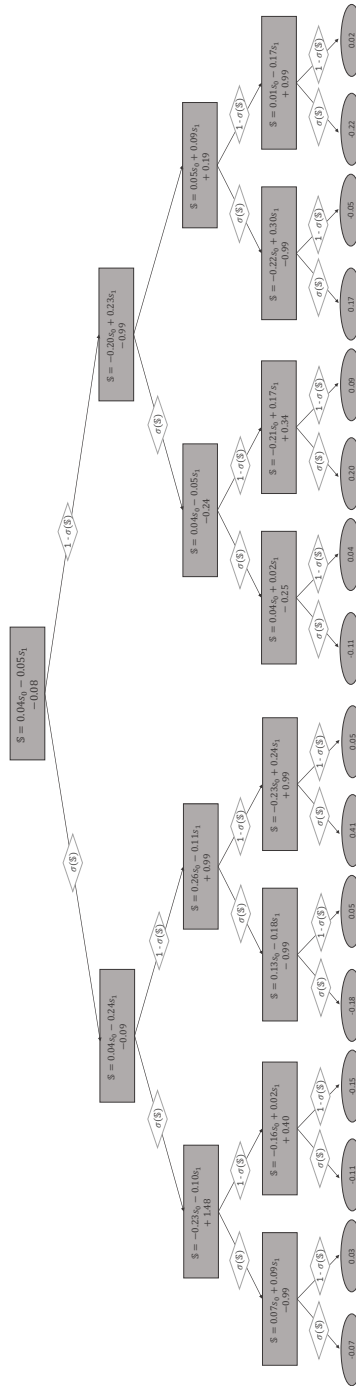


Figure 12: **Exemplary SDT for MountainCar.** This figure visualizes a soft / differentiable DT learned on MountainCar. The tree involves oblique decisions involving multiple features at each split. Additionally, there is no hard decision on which path is selected, but multiple paths are taken with an associated probability. The final prediction is obtained by weighting the leaf outputs with the corresponding path probability.

Detailed results for each of the control datasets are reported in Table 11. The results clearly confirm our intuitive justifications, as each adjustment has a crucial impact on the performance of SYMPOL.

Table 11: **Ablation Study.** We report the average test performance over a total of 25 random trials. This normalized performance consists in normalizing each reward between 0 and 1 via an affine renormalization between the top- and worse-performing models. Instead of the worse-performing model, we use the 20% test reward quantile to account for outliers.

Agent Type	CP	AB	LL	MC-C	PD-C	Normalized Mean (\uparrow)
SYMPOL	500.0	- 79.9	- 57.4	94.3	- 323.3	0.988
w/o separate architecture	135.6	-196.4	-276.8	-552.4	-1219.4	0.080
w/o dynamic rollout	456.1	- 92.0	-178.2	-144.1	- 434.4	0.598
w/o batch size adjustment	498.8	- 81.7	-320.7	-1818.8	- 434.4	0.372
w/o adamW	416.1	- 78.3	- 97.3	0.0	- 393.5	0.865

A.7 RUNTIME AND TRAINING CURVES

The experiments were conducted on a single NVIDIA RTX A6000. The environments for CP, AB, MC-C and PD-C were vectorized (Lange, 2022) and therefore the training is highly efficient, taking only 30 seconds for 1mio timesteps on average (excluding the sequential evaluation which cannot be vectorized). The remaining environments are not vectorized, and we used the standard Gymnasium (Towers et al., 2024) implementation. In Table 12 can clearly see the impact of environment vectorization, as the runtime for LL, which is not vectorized, is more than 10 times higher with over 400 seconds.

Table 12: **Runtime.** We report the average runtime over 25 trials. One trial spans 1mio timesteps for each environment. We excluded LL from the mean runtime calculation, as this is the only non-vectorized environment. To provide a fair comparison of different methods, we aligned the step and batch size.

	SYMPOL (ours)	SDT	MLP
CP	28.8	23.9	25.2
AB	35.5	37.7	33.8
MC-C	23.4	19.4	18.4
PD-C	28.7	28.2	18.5
Mean \downarrow	29.1	27.3	24.0
LL	402.3	394.0	405.6

In addition to the training times, we report detailed training curves for each method. Figure 13 compares the training reward and the test reward of SYMPOL with the full-complexity models MLP and SDT. SYMPOL shows a similar convergence compared to full-complexity models on most environments. For AB, SYMPOL converges even faster than an MLP which can be attributed to the dynamic rollout buffer and batch size. For MC-C we can see that the training of SYMPOL is very unstable at the beginning. We believe that this can be attributed to the sparse reward function of this certain environment and the fact that as a result, minor changes in the policy can result in a severe drop in the reward. Combined with the small rollout buffer and batch size early in the training of SYMPOL, this can result in an unstable training. However, we can see that the training stabilizes later on, which again confirms the design of our dynamic buffer size increasing over time.

Furthermore, we provide a pairwise comparison of SYMPOL with SA-DT and D-SDT in Figure 14. Here, we can again observe the severe information loss for D-SDT and SA-DT by comparing the training curve with the test reward.

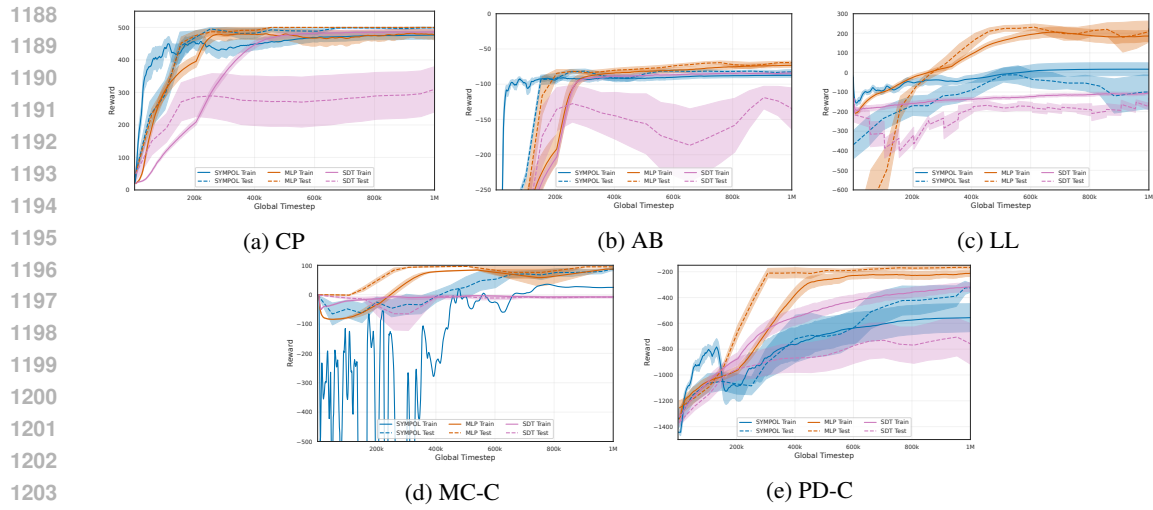


Figure 13: **Training Curves (Full-Complexity)**. Shows the training reward as solid line and the test reward as dashed line for SYMPOL (blue), MLP (orange) and SDT (green).

1241

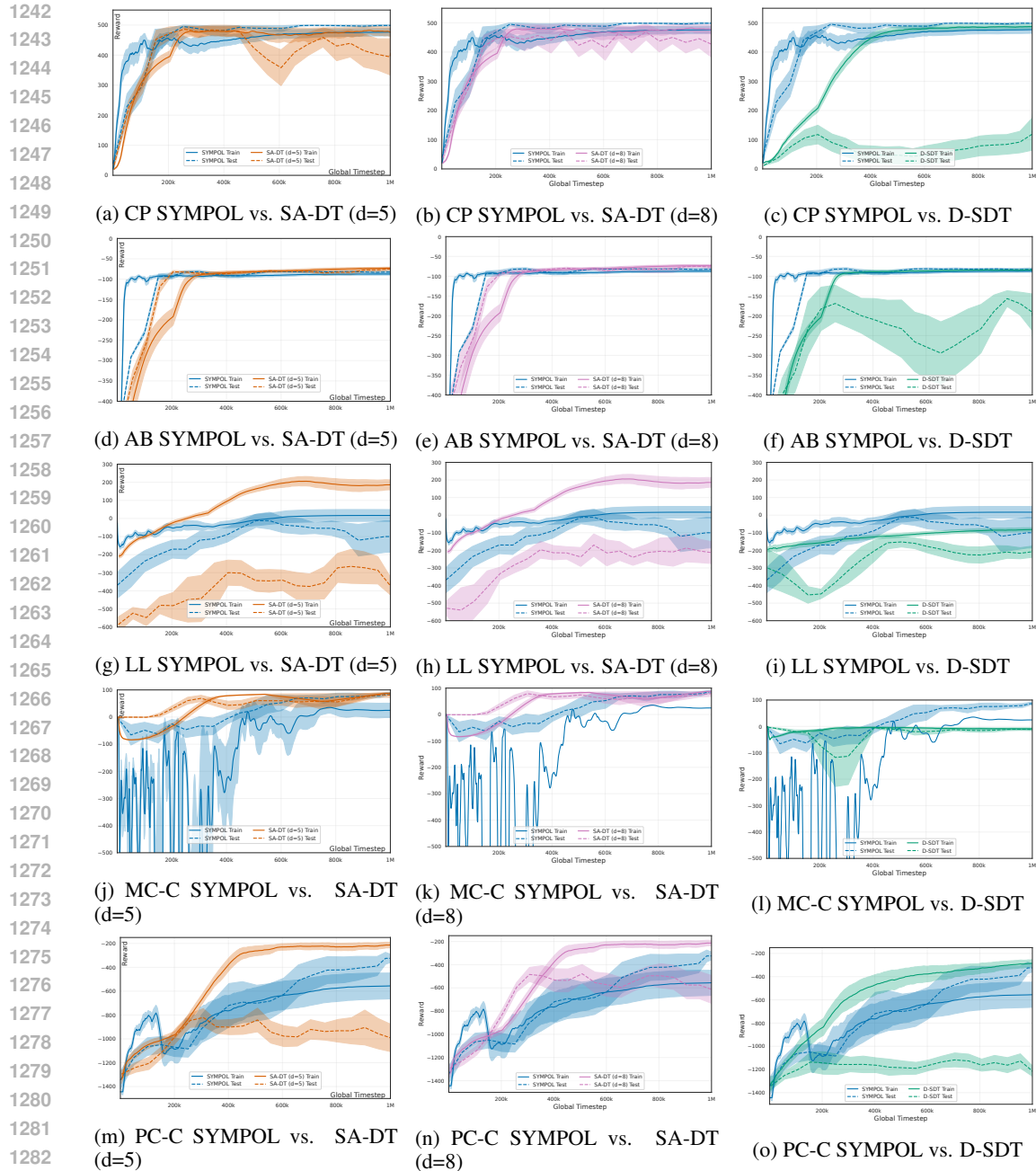


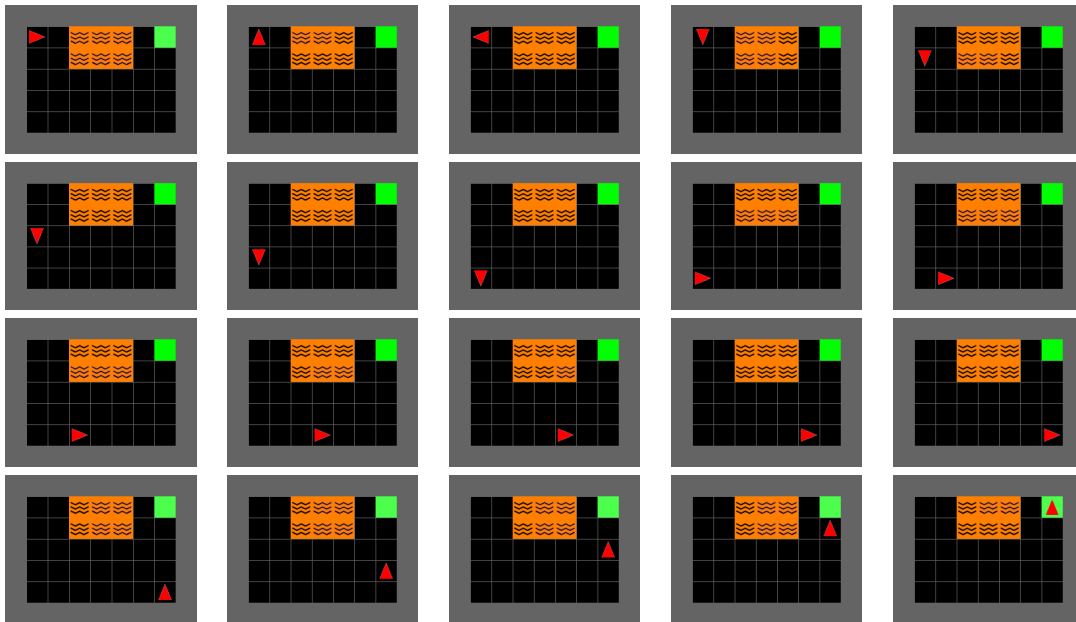
Figure 14: **Training Curves.** Shows the training reward as solid line and the test reward as dashed line for SYMPOL (blue), SA-DT-5 (orange) SA-DT-8 (green) and D-SDT (red). Thereby, the test reward is calculated with the discretized/distilled policy for SA-DT and D-SDT. For several datasets, we can again observe the severe information loss introduced with the post-processing (e.g. for PD-C and LL).

B MINIGRID

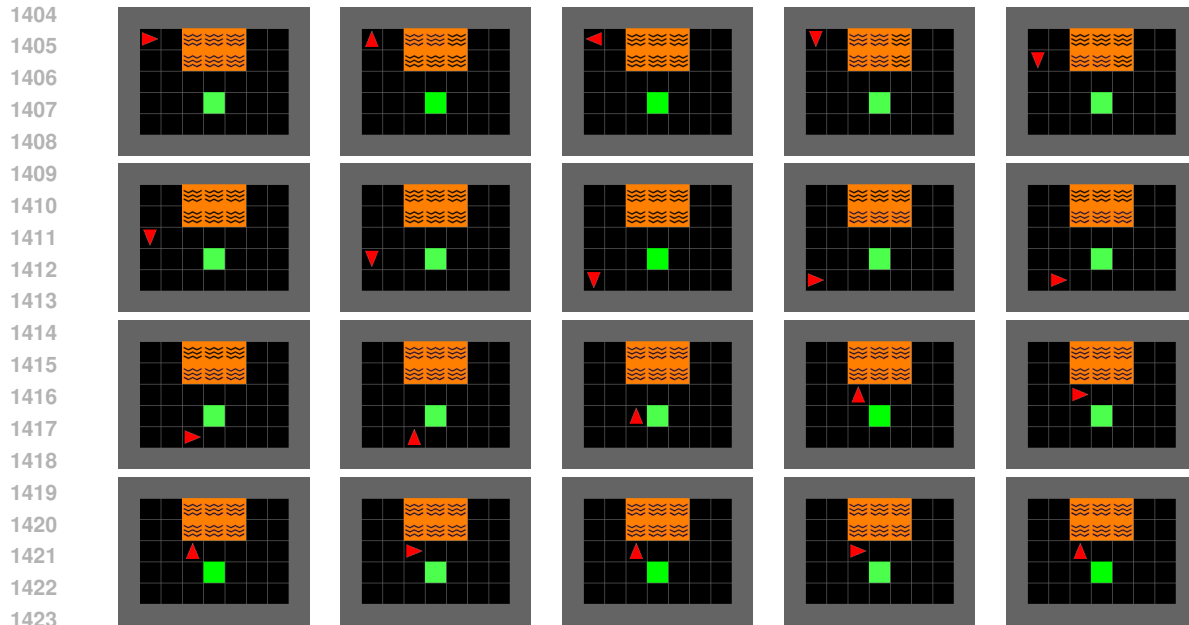
We used the MiniGrid implementation from (Chevalier-Boisvert et al., 2023). For each environment, we limited to observations and the actions to the available ones according to the documentation. Furthermore, we decided to use a view size of 3 to allow a good visualization of the results. In the following, we provide more examples for our MiniGrid Use-Case, along with more detailed visual-

1296 izations. In the following, we visualized the SYMPOL agent sequentially acting in the environment
 1297 as one image for one step from left to right and top to bottom. Figure 15 shows how SYMPOL (see
 1298 image in the main paper or `tree_function(obs)` defined below) solves the environment. Fig-
 1299 ure 16 and Figure 18 show the same agent failing on the environment with domain randomization,
 1300 proving that the agent did not generalize, as we could already observe by inspecting the symbolic,
 1301 tree-based policy. Retraining the agent with domain randomization (see image in the main paper or
 1302 `tree_function_retrained(obs)` defined below), SYMPOL is able to solve the environment
 1303 (see Figure 17 and Figure 19), maintaining interpretability.

1318 B.1 VISUALIZATIONS ENVIRONMENT

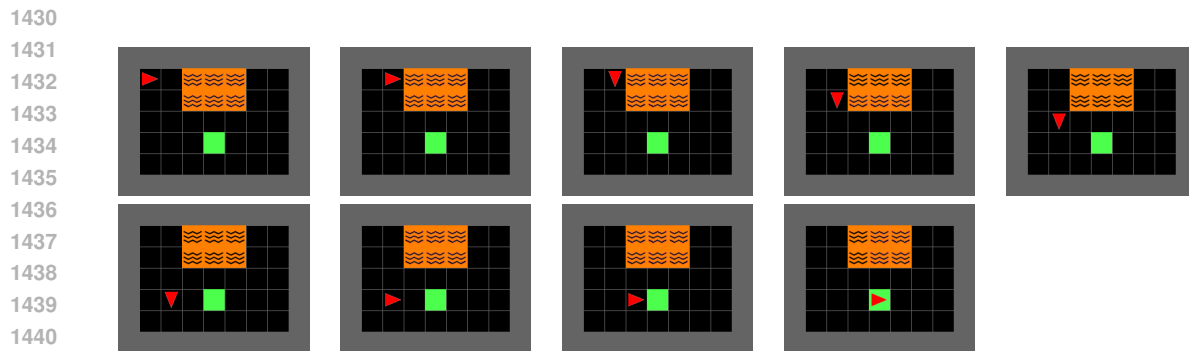


1346
1347 Figure 15: **DistShift SYMPOL**. This figure visualizes the path taken by the SYMPOL agent trained
 1348 on the basic DistShift environment (see image in the main paper or `tree_function(obs)` de-
 1349 fined below) from left to right and top to bottom. The agent follows the wall and reaches the goal at
 the top right corner.



1424
1425
1426
1427
1428
1429

Figure 18: **DistShift (Domain Randomization) SYMPOL Example 2.** This figure visualizes the path taken by the SYMPOL agent trained on the basic DistShift environment (see image in the main paper or `tree_function(obs)` defined below) from left to right and top to bottom. The agent follows the wall until there is no empty space on the left. Instead of an empty space there is the goal, but instead of walking into the goal, the agent surpasses it and again gets stuck at the lava.



1442
1443
1444
1445
1446

Figure 19: **DistShift (Domain Randomization) SYMPOL (retrained) Example 2.** This figure visualizes the path taken by the SYMPOL agent trained on the randomized DistShift environment (see image in the main paper or `tree_function_retrained(obs)` defined below) from left to right and top to bottom. The agent avoids the lava, identifies the goal, and walks into the goal.

1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

1458 B.2 SYMPOL ALGORITHMIC PRESENTATION

```

1459
1460 1 def tree_function(obs):
1461 2     if obs[field one to front and one to left] is 'empty':
1462 3         if obs[field one to front] is 'lava':
1463 4             if obs[field one to left] is 'empty':
1464 5                 action = 'turn right'
1465 6             else:
1466 7                 action = 'turn left'
1467 8         else:
1468 9             action = 'move forward'
1469 10    else:
1470 11        action = 'turn left'
1471 12    return action

```

```

1472
1473 1 def tree_function_retrained(obs):
1474 2     if obs[field one to left] is 'goal':
1475 3         action = 'turn left'
1476 4     else:
1477 5         if obs[field one to right] is 'goal':
1478 6             action = 'turn right'
1479 7         else:
1480 8             if obs[field two to front] is 'wall':
1481 9                 if obs[field one to front and one to right] is 'goal':
1482 10                    if obs[field one to right] is 'lava':
1483 11                        action = 'turn left'
1484 12                    else:
1485 13                        action = 'move forward'
1486 14                else:
1487 15                    if obs[field one to front] is 'goal':
1488 16                        action = 'move forward'
1489 17                    else:
1490 18                        action = 'turn left'
1491 19            else:
1492 20                if obs[field one to front] is 'lava':
1493 21                    action = 'turn right'
1494 22                else:
1495 23                    action = 'move forward'
1496 24    return action

```

1493 C METHODS AND HYPERPARAMETERS

1494
1495 The main methods we compared SYMPOL against are behavioral cloning state-action DTs (SA-DT)
1496 and discretized soft decision trees (D-SDT). In addition to the information given in the paper, we
1497 want to provide some more detailed results of the implementation and refer to our source code for
1498 the exact definition.
1499
1500

- 1501 • **State-Action DTs (SA-DT)** Behavioral cloning SA-DTs are the most common method to
1502 generate interpretable policies post-hoc. Hereby, we first train an MLP policy, which is then
1503 distilled into a DT as a post-processing step after the training. Specifically, we train the DT
1504 on a dataset of expert trajectories generated with the MLP policy. The number of expert
1505 trajectories was set to 25 which we experienced as a good trade-off between dataset size
1506 for the distillation and model complexity during preliminary experiments. The 25 expert
1507 trajectories result in a total of approximately 12500 state-action pairs, varying based on the
1508 environment specification.
- 1509 • **Discretized Soft Decision Trees (D-SDT)** SDTs allow gradient computation by assigning
1510 probabilities to each node. While SDTs exhibit a hierarchical structure, they are usually
1511 considered as less interpretable, since multiple features are considered in a single split and
the whole tree is traversed simultaneously (Marton et al., 2024a). Therefore, Silva et al.

(2020) use SDTs as policies which are discretized post-hoc to allow an easy interpretation considering only a single feature at each split. Discretization is achieved by employing an argmax to obtain the feature index and normalizing the split threshold based on the feature vector. We improved their method by replacing the scaled sigmoid and softmax, with an entmoid and entmax transformation (Peters et al., 2019), resulting in sparse feature selectors with more responsive gradients, as it is common practice Popov et al. (2019); Chang et al. (2021).

In the following, we list the parameter grids used during the hyperparameter optimization (HPO) as well as the optimal parameters selected for each environment. For SYMPOL, SDT and MLP, we optimized the hyperparameters based on the validation reward with optuna Akiba et al., 2019 for 60 trials. Thereby, we ensured that the environments evaluated during the HPO were distinct to the environments used for reporting the test performance in the rest of the paper. Additionally, we decrease the learning rate if no improvement in validation reward is observed for five consecutive iterations, allowing for finer model adjustments in later training stages.

C.1 HPO GRIDS

Table 13: **HPO Grid SYMPOL**

hyperparameter	values
learning_rate_actor_weights	[0.0001, 0.1]
learning_rate_actor_split_values	[0.0001, 0.05]
learning_rate_actor_split_idx_array	[0.0001, 0.1]
learning_rate_actor_leaf_array	[0.0001, 0.05]
learning_rate_actor_log_std	[0.0001, 0.1]
learning_rate_critic	[0.0001, 0.01]
n_update_epochs	[0, 10]
reduce_lr	{True, False}
n_steps	{128, 512}
n_envs	[4, 16]
norm_adv	{True, False}
ent_coef	{0.0, 0.1, 0.2, 0.5}
gae_lambda	{0.8, 0.9, 0.95, 0.99}
gamma	{0.9, 0.95, 0.99, 0.999}
vf_coef	{0.25, 0.50, 0.75}
max_grad_norm	[None]
SWA	{True}
adamW	{True}
depth	{7}
minibatch_size	{64}

Table 14: **HPO Grid MLP**

hyperparameter	values
neurons_per_layer	[16, 256]
num_layers	[1, 3]
learning_rate_actor	[0.0001, 0.01]
learning_rate_critic	[0.0001, 0.01]
minibatch_size	{64, 128, 256, 512}
n_update_epochs	[1, 10]
n_steps	{128, 512}
n_envs	[4, 16]
norm_adv	{True, False}
ent_coef	{0.0, 0.1, 0.2, 0.5}
gae_lambda	{0.8, 0.9, 0.95, 0.99}
gamma	{0.9, 0.95, 0.99, 0.999}
vf_coef	{0.25, 0.50, 0.75}
max_grad_norm	{0.1, 0.5, 1.0, None}

Table 15: **HPO Grid SDT**

hyperparameter	values
critic	{'MLP', 'SDT'}
depth	[4, 8]
temperature	{0.01, 0.05, 0.1, 0.5, 1.0}
learning_rate_actor	[0.0001, 0.01]
learning_rate_critic	[0.0001, 0.01]
minibatch_size	{64, 128, 256, 512}
n_update_epochs	[1, 10]
n_steps	{128, 512}
n_envs	[4, 16]
norm_adv	{True, False}
ent_coef	{0.0, 0.1, 0.2, 0.5}
gae_lambda	{0.8, 0.9, 0.95, 0.99}
gamma	{0.9, 0.95, 0.99, 0.999}
vf_coef	{0.25, 0.50, 0.75}
max_grad_norm	{0.1, 0.5, 1.0, None}

C.2 BEST HYPERPARAMETERS

Table 16: Best Hyperparameters SYMPOL (Control)

	CP	AB	LL	MC-C	PD-C
ent_coef	0.200	0.000	0.000	0.500	0.100
gae_lambda	0.950	0.950	0.900	0.990	0.800
gamma	0.990	0.990	0.999	0.999	0.999
learning_rate_actor_weights	0.048	0.003	0.072	0.000	0.022
learning_rate_actor_split_values	0.000	0.000	0.001	0.000	0.000
learning_rate_actor_split_idx_array	0.026	0.052	0.010	0.000	0.010
learning_rate_actor_leaf_array	0.020	0.005	0.009	0.028	0.006
learning_rate_actor_log_std	0.001	0.002	0.021	0.094	0.000
learning_rate_critic	0.001	0.000	0.002	0.002	0.000
max_grad_norm	None	None	None	None	None
n_envs	7	8	6	5	15
n_steps	512	128	512	128	128
n_update_epochs	7	7	7	2	7
norm_adv	False	False	True	False	True
reduce_lr	True	True	True	True	False
vf_coef	0.500	0.250	0.500	0.500	0.750
SWA	True	True	True	True	True
adamW	True	True	True	True	True
dropout	0.000	0.000	0.000	0.000	0.000
depth	7	7	7	7	7
minibatch_size	64	64	64	64	64
n_estimators	1	1	1	1	1

Table 17: Best Hyperparameters SYMPOL (MiniGrid)

	E-R	DK	LG-5	LG-7	DS
ent_coef	0.100	0.200	0.100	0.100	0.500
gae_lambda	0.990	0.950	0.900	0.900	0.950
gamma	0.900	0.990	0.950	0.990	0.999
learning_rate_actor_weights	0.063	0.042	0.055	0.001	0.036
learning_rate_actor_split_values	0.001	0.001	0.006	0.001	0.000
learning_rate_actor_split_idx_array	0.001	0.001	0.012	0.001	0.009
learning_rate_actor_leaf_array	0.003	0.004	0.009	0.008	0.001
learning_rate_actor_log_std	0.043	0.021	0.005	0.002	0.038
learning_rate_critic	0.001	0.001	0.001	0.001	0.001
max_grad_norm	None	None	None	None	None
n_envs	14	14	16	7	10
n_steps	128	512	512	128	512
n_update_epochs	8	9	5	4	5
norm_adv	True	True	True	True	False
reduce_lr	False	True	True	True	True
vf_coef	0.500	0.500	0.250	0.500	0.250
SWA	True	True	True	True	True
adamW	True	True	True	True	True
dropout	0.000	0.000	0.000	0.000	0.000
depth	7	7	7	7	7
minibatch_size	64	64	64	64	64
n_estimators	1	1	1	1	1

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

Table 18: **Best Hyperparameters MLP (Control)**

	CP	AB	LL	MC-C	PD-C
adamW	False	False	False	False	False
ent_coef	0.200	0.000	0.100	0.100	0.100
gae_lambda	0.900	0.900	0.900	0.950	0.950
gamma	0.999	0.990	0.999	0.999	0.990
learning_rate_actor	0.001	0.000	0.001	0.005	0.000
learning_rate_critic	0.003	0.005	0.003	0.001	0.002
max_grad_norm	1.000	1.000	0.500	0.100	None
minibatch_size	256	256	128	512	128
n_envs	13	12	13	15	8
n_steps	128	512	512	512	512
n_update_epochs	7	9	8	2	2
neurons_per_layer	139	185	46	240	75
norm_adv	False	True	False	True	True
num_layers	2	2	3	2	2
reduce_lr	False	False	False	False	False
vf_coef	0.250	0.500	0.500	0.250	0.250

Table 19: **Best Hyperparameters MLP (MiniGrid)**

	E-R	DK	LG-5	LG-7	DS
adamW	False	False	False	False	False
ent_coef	0.100	0.100	0.100	0.100	0.100
gae_lambda	0.950	0.900	0.950	0.950	0.990
gamma	0.990	0.900	0.990	0.900	0.990
learning_rate_actor	0.000	0.000	0.002	0.000	0.000
learning_rate_critic	0.001	0.000	0.003	0.001	0.001
max_grad_norm	0.100	0.100	1	0.500	0.100
minibatch_size	64	256	128	512	256
n_envs	13	8	8	12	10
n_steps	512	256	512	128	128
n_update_epochs	5	7	9	8	7
neurons_per_layer	112	169	76	28	158
norm_adv	False	True	False	True	True
num_layers	3	1	1	1	2
reduce_lr	False	False	False	False	False
vf_coef	0.500	0.500	0.250	0.750	0.500

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

Table 20: **Best Hyperparameters SDT (Control)**

	CP	AB	LL	MC-C	PD-C
adamW	False	False	False	False	False
critic	mlp	mlp	mlp	mlp	mlp
depth	7	6	8	7	7
ent_coef	0.000	0.100	0.200	0.000	0.200
gae_lambda	0.950	0.950	0.990	0.900	0.900
gamma	0.990	0.990	0.999	0.990	0.900
learning_rate_actor	0.001	0.002	0.001	0.001	0.000
learning_rate_critic	0.000	0.000	0.001	0.007	0.000
max_grad_norm	0.100	0.100	1.000	0.500	0.100
minibatch_size	128	128	128	64	128
n_envs	15	6	7	14	7
n_steps	512	128	512	512	256
n_update_epochs	4	10	2	1	7
norm_adv	True	False	True	False	False
reduce_lr	False	False	False	False	False
temperature	1	0.500	1	1	0.100
vf_coef	0.500	0.500	0.750	0.250	0.500

Table 21: **Best Hyperparameters SDT (MiniGrid)**

	E-R	DK	LG-5	LG-7	DS
adamW	False	False	False	False	False
critic	sdt	mlp	sdt	sdt	sdt
depth	7	6	7	8	7
ent_coef	0.100	0.100	0.200	0.100	0.100
gae_lambda	0.900	0.950	0.990	0.950	0.900
gamma	0.990	0.900	0.999	0.950	0.950
learning_rate_actor	0.004	0.001	0.000	0.002	0.001
learning_rate_critic	0.000	0.002	0.000	0.005	0.002
max_grad_norm	0.100	0.100	0.500	0.100	None
minibatch_size	512	256	512	256	512
n_envs	10	10	10	13	5
n_steps	512	256	256	128	512
n_update_epochs	5	10	8	4	7
norm_adv	True	True	True	True	True
reduce_lr	False	False	False	False	False
temperature	1	1	1	1	1
vf_coef	0.750	0.750	0.750	0.250	0.750