

---

# Gradient Boosting Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Neural networks (NN) achieve remarkable results in various tasks, but lack key  
2 characteristics: interpretability, support for categorical features, and lightweight im-  
3 plementations suitable for edge devices. While ongoing efforts aim to address these  
4 challenges, Gradient Boosting Trees (GBT) inherently meet these requirements.  
5 As a result, GBTs have become the go-to method for supervised learning tasks  
6 in many real-world applications and competitions. However, their application in  
7 online learning scenarios, notably in reinforcement learning (RL), has been limited.  
8 In this work, we bridge this gap by introducing Gradient-Boosting RL (GBRL),  
9 a framework that extends the advantages of GBT to the RL domain. Using the  
10 GBRL framework, we implement various actor-critic algorithms and compare their  
11 performance with their NN counterparts. Inspired by shared backbones in NN  
12 we introduce a tree-sharing approach for policy and value functions with distinct  
13 learning rates, enhancing learning efficiency over millions of interactions. GBRL  
14 achieves competitive performance across a diverse array of tasks, excelling in  
15 domains with structured or categorical features. Additionally, we present a high-  
16 performance, GPU-accelerated implementation that integrates seamlessly with  
17 widely-used RL libraries. GBRL expands the toolkit for RL practitioners, demon-  
18 strating the viability and promise of GBT within the RL paradigm, particularly in  
19 domains characterized by structured or categorical features.

## 20 1 Introduction

21 Reinforcement Learning (RL) has shown great promise in various domains that involve sequential  
22 decision making. However, many real-world tasks, such as inventory management, traffic signal  
23 optimization, network optimization, resource allocation, and robotics, are represented by structured  
24 observations with categorical or mixed data types. These tasks can benefit significantly from  
25 deployment and training on edge devices due to resource constraints. Moreover, interpretability  
26 is crucial in these applications for regulatory reasons and for trust in the decision-making process.  
27 Current neural network (NN) based solutions struggle with interpretability, handling categorical data,  
28 and supporting light implementations suitable for low-compute devices.

29 Gradient Boosting Trees (GBT) is a powerful ensemble method extensively used in supervised  
30 learning due to its simplicity, accuracy, interpretability, and natural handling of structured and  
31 categorical data. Frameworks such as XGBoost [7], LightGBM [20], and CatBoost [36] have become  
32 integral in applications spanning finance [49], healthcare [54, 27, 43], and competitive data science  
33 [6]. Despite their successes, GBT has seen limited application in RL. This is primarily because  
34 traditional GBT libraries are designed for static datasets with predefined labels, contrasting with the  
35 dynamic nature of RL. The distribution shift in both input (state) and output (reward) poses significant  
36 challenges for the direct application of GBT in RL. Moreover, there is a notable lack of benchmarks  
37 or environments tailored for structured data, further hindering progress in this area.

38 In this paper, we introduce Gradient Boosting Reinforcement Learning (GBRL), a GBT framework  
 39 tailored for RL. Our contributions are:

- 40 1. **GBT for RL.** We demonstrate the viability and potential of GBT as function approximators  
 41 in RL. We present GBT-based implementations of PPO, A2C, and AWR, and show that  
 42 GBRL is competitive with NNs across a range of environments. In addition, similarly to  
 43 supervised learning, GBRL outperforms NNs on categorical tasks (see Figure 1).
- 44 2. **Tree-based Actor-Critic architecture.** Inspired by shared architectures in NN-based  
 45 actor-critic (AC), we introduce a GBT-based AC architecture. This reduces the memory and  
 46 computational requirements by sharing a common ensemble structure for both the policy and  
 47 value. This approach significantly reduces runtime compared to existing GBT frameworks,  
 48 thus removing the barrier to solving complex, high-dimensional RL tasks with millions of  
 49 interactions.
- 50 3. **Modern GBT-based RL library.** We provide a CUDA-based [33] hardware-accelerated  
 51 GBT framework optimized for RL. GBRL is designed to work as part of a broader system  
 52 and seamlessly integrates with popular repositories such as Stable-baselines3 [39]. This new  
 53 tool offers practitioners a powerful option for exploring GBT in RL settings.<sup>1</sup>

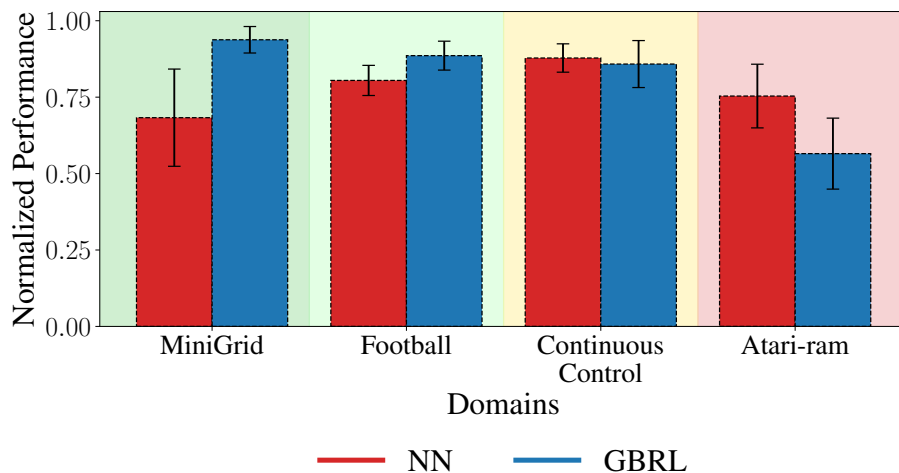


Figure 1: **PPO GBRL vs PPO NN.** Aggregated mean and standard deviation of the normalized average reward for the final 100 episodes. Rewards were normalized as:  $\text{reward}_{\text{norm}} = \frac{\text{reward}}{\text{reward}_{\text{max}}\{\text{NN}, \text{GBRL}\}}$  per environment and then aggregated across each domain.

## 54 2 Related Work

55 **Gradient boosted trees.** Recent advances have extended GBT’s capabilities beyond traditional  
 56 regression and classification. In ranking problems, GBT has been used to directly optimize ranking  
 57 metrics, as demonstrated by frameworks like StochasticRank [51] and recent advancements explored  
 58 in Lyzhin et al. [26]. Additionally, GBT offer probabilistic predictions through frameworks like  
 59 NGBoost [11], enabling uncertainty quantification [28]. The connection between GBT and Gaussian  
 60 Processes [52, 45] offers further possibilities for uncertainty-aware modeling. Recently, Ivanov and  
 61 Prokhorenkova [18] modeled graph-structured data by combining GBT with graph neural networks.

62 Despite their versatility, applying GBT in RL remains a relatively less explored area. Several works  
 63 have employed GBT as a function approximator within off-policy RL methods, including its use  
 64 in Q-learning [1] and in bandit settings to learn inverse propensity scores [24]. Recently, Brukhim  
 65 et al. [5] proposed a boosting framework for RL where a base policy class is incrementally enhanced  
 66 using linear combinations and nonlinear transformations. However, these previous works have not  
 67 yet demonstrated the scalability and effectiveness in complex, high-dimensional RL environments

<sup>1</sup>We attached the GBRL repository as supplementary material and will release it after the review process.

68 requiring extensive interactions. In this work, we show how to adapt the framework of GBT to  
69 successfully solve large-scale RL problems.

70 **Interpretability.** Due to the inherent non-linearities, NNs are challenging to interpret and require  
71 sophisticated methods to do so. Interpreting NNs often involves either approximation with simpler  
72 models such as decision trees or using gradient-based techniques, which require additional forward  
73 and backward passes [16, 9, 44, 3, 37]. On the other hand, interpretability methods for GBT can take  
74 advantage of the structure of a decision tree for high speed, efficiency, and accuracy [25, 10].

75 **Structured and categorical data.** Previous work in RL has predominantly focused on using  
76 NNs due to their ability to capture complex patterns in high-dimensional data. Techniques such as  
77 Q-learning and AC methods have advanced significantly, demonstrating success in tasks involving  
78 raw sensory inputs like images, text, and audio. However, NNs that perform well on structured and  
79 categorical data typically have very specialized architectures [19, 46, 14, 2] and are not standard  
80 multi-layer perceptrons (MLPs) that are often used in many RL tasks and algorithms [34]. Even with  
81 these specialized architectures, Gradient Boosting Trees (GBT) often perform equally or better on  
82 structured and categorical datasets [19, 31, 14, 15].

83 **Policy optimization through functional gradient ascent.** In this approach, the policy is parameter-  
84 ized by a growing linear combination of functions [29]. Each linear addition represents the functional  
85 gradient with respect to current parameters. Kersting and Driessens [21] demonstrated the direct  
86 optimization of policies using the policy gradient theorem [48]. Similarly, Scherrer and Geist [41]  
87 proposed a functional gradient ascent approach as a local policy search algorithm. While these works  
88 lay theoretical groundwork, practical results on complex, high-dimensional RL environments have  
89 not been shown. To adapt GBT’s to RL, we leverage the framework of functional gradient ascent.  
90 This combination enables a seamless integration of GBRL directly into existing RL optimization  
91 packages, such as Stable-baselines [39].

## 92 3 Preliminaries

93 We begin by introducing Markov Decision Processes (MDPs) and the AC schema. Then, we introduce  
94 GBT. In the following section, we show how to combine both of these paradigms into GBRL.

### 95 3.1 Markov Decision Process

96 We consider a fully observable infinite-horizon Markov decision process (MDP) characterized by the  
97 tuple  $(\mathcal{S}, \mathcal{A}, P, \mathcal{R})$ . At each step, the agent observes a state  $\mathbf{s} \in \mathcal{S}$  and samples an action  $\mathbf{a} \in \mathcal{A}$  from  
98 its policy  $\pi(\mathbf{s}, \mathbf{a})$ . Performing the action causes the system to transition to a new state  $\mathbf{s}'$  based on  
99 the transition probabilities  $P(\mathbf{s}' | \mathbf{s}, \mathbf{a})$ , and the agent receives a reward  $\mathbf{r} \sim \mathcal{R}(\mathbf{s}, \mathbf{a})$ . The objective is  
100 to find an optimal policy  $\pi^*$  that maximizes the expected discounted reward  $J(\pi) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}_t]$ ,  
101 with a discount factor  $\gamma \in [0, 1)$ .

102 The action-value function  $Q_{\pi}(\mathbf{s}, \mathbf{a}) := \mathbb{E}_{\pi}[\sum_{t'=0}^{\infty} \gamma^{t'} \mathcal{R}(\mathbf{s}_{t+t'}, \mathbf{a}_{t+t'}) | \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a}]$  estimates the  
103 expected returns of performing action  $\mathbf{a}$  in state  $\mathbf{s}$  and then acting according to  $\pi$ . Additionally, the  
104 value function  $V_{\pi}(\mathbf{s}) := \mathbb{E}_{\pi}[\sum_{t'=0}^{\infty} \gamma^{t'} \mathcal{R}(\mathbf{s}_{t+t'}, \mathbf{a}_{t+t'}) | \mathbf{s}_t = \mathbf{s}]$ , predicts the expected return starting  
105 from state  $\mathbf{s}$  and acting according to  $\pi$ . Finally, the advantage function  $A_{\pi}(\mathbf{s}, \mathbf{a}) := Q_{\pi}(\mathbf{s}, \mathbf{a}) - V_{\pi}(\mathbf{s})$ ,  
106 indicates the expected relative benefit of performing action  $\mathbf{a}$  over acting according to  $\pi$ .

### 107 3.2 Actor-Critic Reinforcement Learning

108 Actor-critic methods are a common method to solve the objective  $J(\pi)$ . They learn both the policy  
109 and value. In the GBRL framework, we extend three common AC algorithms to support GBT-based  
110 function approximators.

111 **A2C** [32] is a synchronous, on-policy AC algorithm designed to improve learning stability. The critic  
112 learns a value function,  $V(\mathbf{s})$ , used to estimate the advantage. This advantage is incorporated into the  
113 policy gradient updates, reducing variance and leading to smoother learning. The policy is updated  
114 using the following gradient:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}[\nabla_{\theta} \log \pi_{\theta}(\mathbf{a} | \mathbf{s}) A(\mathbf{s}, \mathbf{a})]$ .

115 **PPO** [42] extends A2C by improving stability. This is achieved through constrained policy  
 116 update steps using a clipped surrogate objective. This prevents drastic policy changes and  
 117 leads to smoother learning. To achieve this, PPO solves the following objective:  $\nabla_{\theta} J(\pi_{\theta}) =$   
 118  $\mathbb{E}[\nabla_{\theta} \text{clip}(\frac{\log \pi_{\theta}(\mathbf{a} | \mathbf{s})}{\log \pi_{\theta_{\text{old}}}(\mathbf{a} | \mathbf{s})}, 1 - \epsilon, 1 + \epsilon) A(\mathbf{s}, \mathbf{a})]$ . Additionally, PPO enhances sample efficiency by  
 119 performing multiple optimization steps on each collected rollout.

120 **AWR** [35] is an off-policy AC algorithm. Provided a dataset  $\mathcal{D}$ , AWR updates both the policy and the  
 121 value through supervised learning. This dataset can be pre-defined and fixed (offline), or continually  
 122 updated using the agents experience (replay buffer). At each training iteration  $k$ , AWR solves the  
 123 following two regression problems:

$$V_k = \arg \min_V \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [\|G(\mathbf{s}, \mathbf{a}) - V(\mathbf{s})\|_2^2], \pi_{k+1} = \arg \max_{\pi} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [\log \pi(\mathbf{a} | \mathbf{s}) \exp(\frac{1}{\beta} A_k(\mathbf{s}, \mathbf{a}))],$$

124 where  $G(\mathbf{s}, \mathbf{a})$  represents the monte-carlo estimate or TD( $\lambda$ ) of the expected return [47].

### 125 3.3 Gradient Boosting Trees as Functional Gradient Descent

126 Gradient boosting trees (GBT) [12] are a non-parametric machine learning technique that combines  
 127 decision tree ensembles with functional gradient descent [30]. GBT iteratively minimizes the expected  
 128 loss  $L(F(\mathbf{x})) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [L(\mathbf{y}, F(\mathbf{x}))]$  over a dataset  $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ . A GBT model,  $F_K$ , predicts  
 129 outputs using  $K$  additive trees as follows:

$$F_K(\mathbf{x}_i) = F_0 + \sum_{k=1}^K \epsilon h_k(\mathbf{x}_i), \quad (1)$$

130 where  $\epsilon$  is the learning rate,  $F_0$  is the base learner, and each  $h_k$  is an independent regression tree  
 131 partitioning the feature space.

132 In the context of functional gradient descent, the objective is to minimize the expected loss  $L(F(\mathbf{x})) =$   
 133  $\mathbb{E}_{\mathbf{x}, \mathbf{y}} [L(\mathbf{y}, F(\mathbf{x}))]$  with respect to the functional  $F$ . Here, a functional  $F : \mathcal{H} \rightarrow \mathbb{R}$  maps a function  
 134 space to real numbers. A GBT model can be viewed as a functional  $F$  that maps a linear combination  
 135 of binary decision trees to outputs:  $F : \text{lin}(\mathcal{H}) \rightarrow \mathbb{R}^D$ , where  $\mathcal{H}$  is the decision tree function class.

136 We start with an initial model,  $F_0$ , and iteratively add trees to  $F$  to minimize the expected loss.  
 137 Similar to parametric gradient descent, at each iteration  $k$ , we minimize the loss by taking a step in  
 138 the direction of the functional gradient  $g_k^i := \nabla_{F_{k-1}} L(\mathbf{y}_i, F_{k-1}(\mathbf{x}_i))$ . However, we are constrained  
 139 to gradient directions within  $\mathcal{H}$ . Thus, we project the gradient  $g_k$  into a decision tree by solving:

$$h_k = \arg \min_h \|\epsilon g_k - h(\mathbf{x})\|_2^2. \quad (2)$$

## 140 4 Gradient Boosting Reinforcement Learning

141 In this work, we extend the framework of GBT to support AC algorithms in the task of RL. The  
 142 objective in RL is to optimize the return  $J$ , the cumulative reward an agent receives. Unlike in  
 143 supervised learning, the target predictions are unknown a priori. RL agents learn through trial  
 144 and error. Good actions are reinforced by taking a step in the direction of the gradient  $\nabla_{\pi} J$ . This  
 145 formulation aligns perfectly with functional gradient ascent; thus, in GBRL, we optimize the objective  
 146 directly over the decision tree function class. This is achieved by iteratively growing the ensemble of  
 147 trees  $\{h_i\}$ . The ensemble outputs  $\theta$ , representing AC parameters such as the policy  $\pi$  and the value  
 148 function. For example,  $\theta = [\mu(\mathbf{s}), \sigma(\mathbf{s}), V(\mathbf{s})]$  for a Gaussian policy. At each iteration, a new tree  
 149  $h_k$ , constructed to minimize the distance to  $\nabla_{\theta_{k-1}} J$ , is added to the ensemble. Here, The resulting  
 150 method is an application of GBT as a functional gradient optimizer  $\theta_k \approx \theta_0 + \epsilon \sum_{m=0}^{k-1} \nabla_{\theta_m} J$ .

151 However, RL presents unique challenges for GBT. RL involves a nonstationary state distribution and  
 152 inherent online learning, causing gradients to vary in magnitude and direction. Large gradients in  
 153 unfavorable directions risk destabilizing training or leading to catastrophic forgetting. Moreover,  
 154 feedback in RL is provided through interactions with the environment and is not available a priori.  
 155 This contrasts with supervised learning settings, where gradients decrease with boosting iterations,  
 156 and targets are predefined. As a result, many of the key features that traditional GBT libraries rely on

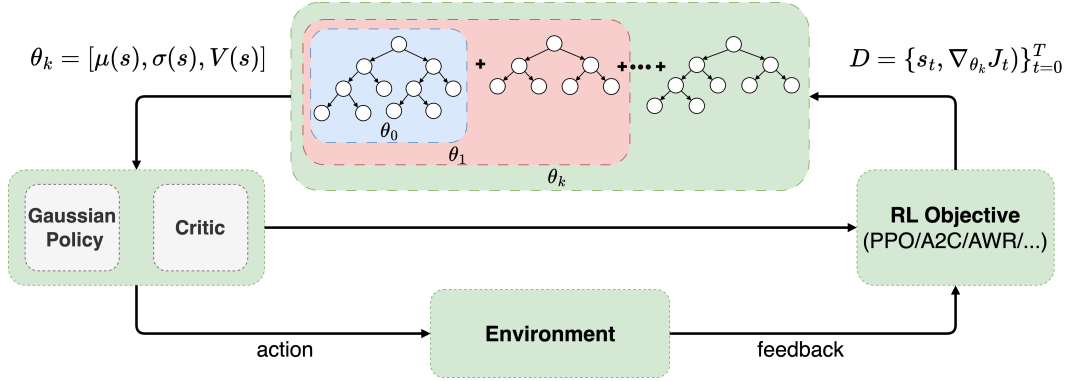


Figure 2: **The GBRL framework.** The actor’s policy and critic’s value function are parameterized by  $\theta_k$ . For example,  $\theta_k = [\mu(s), \sigma(s), V(s)]$  for a Gaussian policy.  $\theta_k$  is calculated by summing all the outputs of trees in the ensemble. Starting from  $\theta_0$ , at each training iteration, GBRL collects a rollout and computes the gradient  $\nabla_{\theta_0} J$ . This gradient is then used to fit the next tree added to the ensemble, which is updated to  $\theta_1$ . This process repeats with each iteration fitting a new tree, refining the parameterization, and expanding the ensemble towards  $\theta_k \approx \theta_0 + \epsilon \sum_{m=0}^{k-1} \nabla_{\theta_m} J$ , an approximated scaled sum of gradients with respect to past parametrizations.

157 are not suitable. For example, GOSS [20], categorical feature encoding [36], early-stopping signals,  
 158 pruning methods [53], and strategies to tackle online learning [55].

159 To address these challenges, we employ appropriate tools from the NN and GBT literature, such as  
 160 batch learning [13, 40] to update the ensemble. At each boosting iteration, we fit a decision tree on  
 161 a random batch sampled with replacement from the experience buffer. This approach helps handle  
 162 non-stationary distributions and improve stability by focusing on different parts of the state space,  
 163 allowing beneficial gradient directions to accumulate and minimizing the impact of detrimental ones.  
 164 Additionally, GBRL fits gradients directly to optimize objectives, whereas traditional GBT methods  
 165 require targets and need workarounds to utilize gradients effectively.

166 A common theme in AC algorithms is to utilize a shared approximation for the actor and the critic.  
 167 We adopt this approach in GBRL, constructing trees where each leaf provides two predictions. GBRL  
 168 predicts both the policy (distribution over actions) and the value estimate. The internal structure of the  
 169 tree is shared, providing a single feature representation for both objectives and significantly reducing  
 170 memory and computational bottlenecks. Accordingly, in GBRL we apply differentiated learning rates  
 171 to the policy and value outputs during prediction, effectively optimizing distinct objectives within  
 172 this shared structure. We present the full algorithm in Algorithm 1 and diagram in Figure 2.

## 173 5 Experiments

174 Our experiments aim to answer the following questions:

- 175 1. **GBT as RL Function Approximator:** Can GBT-based AC algorithms effectively solve  
 176 complex high-dimensional RL tasks?
- 177 2. **Comparison to NNs:** How does GBRL compare with NN-based training in various RL  
 178 algorithms?
- 179 3. **Benefits in Categorical Domains:** Do the benefits of GBT in supervised learning transfer  
 180 to the realm of RL?
- 181 4. **Comparison to Traditional GBT libraries:** Can we use traditional GBT libraries instead  
 182 of GBRL for RL tasks?
- 183 5. **Evaluating the shared AC architecture:** How does sharing the tree structure between the  
 184 actor and the critic impact performance?

185 We implemented GBT-based versions of A2C, PPO, and AWR within Stable Baselines3. We refer to  
 186 our implementations as PPO GBRL, A2C GBRL, and AWR GBRL. We evaluated GBRL against the

---

**Algorithm 1** Gradient Boosting for Reinforcement Learning (GBRL)

---

```
1: Initialize:  $\theta_0, \epsilon_{\text{actor}}, \epsilon_{\text{critic}}$ , experience buffer  $\mathcal{B}$ , total training iterations  $K$ , number of updates  $U$ ,  
   batch size  $N$ ,  $k \leftarrow 1$   
2: while  $k < K$  do  
3:   Collect trajectory  $\tau^{(k)} = (\mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_T, \mathbf{a}_T)^{(k)}$  and rewards  $(\mathbf{r}_0, \dots, \mathbf{r}_T)^{(k)}$  using  $\pi_{\theta_{k-1}}$   
4:   Add trajectory  $\tau^{(k)}$  and rewards to the experience buffer  $\mathcal{B}$   
5:   for each update  $u = 1$  to  $U$  do  
6:     Sample a batch from the experience buffer  $\mathcal{B}$   
7:     Compute gradients  $g$  according to AC algorithm (e.g., PPO, A2C, AWR)  
8:     Construct dataset  $D = \{(\mathbf{s}_n, g_n)\}_{n=0}^N$  and fit a decision tree  $h_k$   
9:     for each dimension  $d = 0$  to  $D$  do  
10:      if  $0 \leq d < D$  then  
11:        Update  $\theta_k^{(d)} = \theta_{k-1}^{(d)} + \epsilon_{\text{actor}} h_k(\mathbf{s})$   
12:      else  
13:        Update  $\theta_k^{(d)} = \theta_{k-1}^{(d)} + \epsilon_{\text{critic}} h_k(\mathbf{s})$   
14:       $k \leftarrow k + 1$   
15: Output: AC parameters  $\theta_K^{(d)}(\mathbf{s})$  for  $d = 0, 1, \dots, D$ 
```

---

187 equivalent NN implementations. Where available, we utilize hyperparameters from RL Baselines3  
188 Zoo [38]; otherwise, we optimize the hyperparameters for specific environments. The AWR NN  
189 implementation is based on the original paper [35].

190 We conducted experiments on a range of RL domains. We test classic control tasks, high-dimensional  
191 vectorized problems, and finally categorical tasks. We use 5 random seeds per experiment on a  
192 single NVIDIA V100-32GB GPU. We present the cumulative non-discounted reward, averaged  
193 across the last 100 episodes. We normalize the plots for simple visual comparison between GBRL  
194 and the corresponding NN implementations. The normalized score is computed as  $\text{score}_{\text{norm}} =$   
195  $\frac{\text{score}_{\text{GBRL}} - \text{score}_{\text{NN}}}{\text{score}_{\text{max}\{\text{NN}, \text{GBRL}\}} - \text{score}_{\text{min}\{\text{NN}, \text{GBRL}\}}}$ . We provide the full learning curves, implementation details, compute  
196 resources, un-normalized numerical results, and hyperparameters in the supplementary material.

197 **Classic Environments.** We evaluate GBRL’s ability to solve classic RL tasks using Continuous-  
198 Control and Box2D environments, provided via Gym [50]. We trained agents for 1M steps (1.5M for  
199 LunarLander-v2) and provide the results in Figure 3. For exact values, refer to Table 2.

200 Considering the algorithmic objective, we observe that GBRL and NN present similar performance  
201 when optimized using PPO. In contrast, the other methods demonstrate inconclusive results. In  
202 certain environments, such as MountainCar, GBRL outperforms NN with all AC methods. On the  
203 other hand, in Pendulum NN is better.

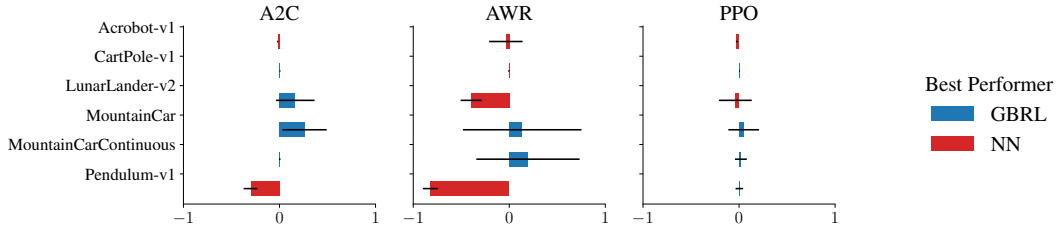


Figure 3: **Continuous-Control and Box2D environments.** Normalized comparison between GBRL and NN. PPO, the best performing method, shows similar performance with GBRL and NN function classes.

204 **High-Dimensional Vectorized Environments.** The decision-tree function class operates on in-  
205 dividual features at each step. Consequently, this function class is not well-suited for handling  
206 pixel-based representations, which require more complex feature interactions. Therefore, we evaluated  
207 GBRL in the Football [22] and Atari RAM [4] domains. These offer high-dimensional vectorized

208 representations. We trained agents in both environments for 10 million timesteps. The complete  
 209 results are reported in Tables 3 and 4 and illustrated in Figure 4.

210 The results portray the following phenomenon. While both tasks may seem similar, there is a distinct  
 211 difference. The features in the football domain are manually constructed and represent identifiable  
 212 information, such as the location of the ball and the players. However, the Atari RAM domain  
 213 provides a flattened view of the system RAM, which is unstructured.

214 At their core, binary decision trees are if-else clauses. This function class is naturally suited to work  
 215 with structured data. These insights are emphasized in the football domain. Here, PPO GBRL greatly  
 216 outperforms PPO NN across most environments and exhibits equivalent performance on the rest. In  
 217 addition, as Atari RAM is unstructured, we observe that, as can be expected, in most cases GBRL  
 218 underperforms NN, except for AWR. However, AWR NN underperformed considerably compared to  
 219 the other NN implementations.

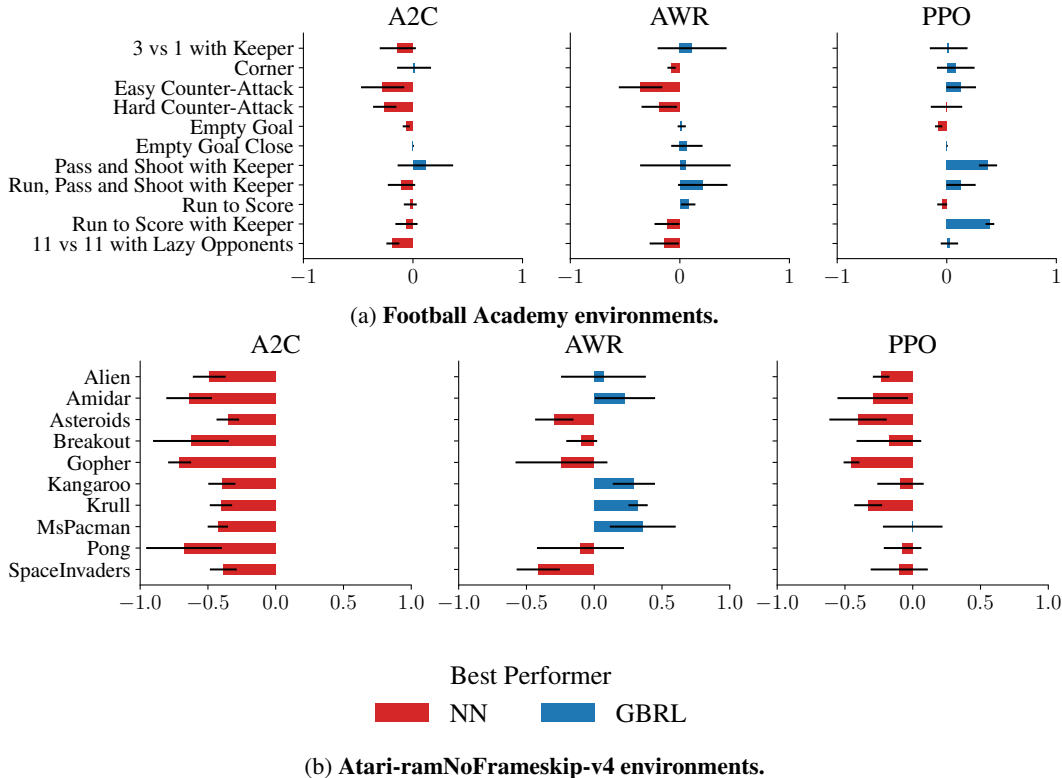


Figure 4: **High-Dimensional Vectorized Environments.** GBRL outperforms NN on the structured Football domain using PPO. NN outperforms on unstructured tasks, such as Atari RAM.

220 **Categorical Environments.** The football experiment suggests that GBRL outperforms when  
 221 assigned structured data. Here, we focus on categorical environments. This is a regime where GBT  
 222 excels in supervised learning. In these experiments, we evaluated the MiniGrid domain [8]. It consists  
 223 of 2D grid worlds with goal-oriented tasks that require object interaction. We trained in PutNear,  
 224 FourRooms, and Fetch tasks for 10M timesteps, matching the reported PPO NN in RL Baselines3  
 225 Zoo. We trained the remaining environments for 1M timesteps. We give the results in Figure 5. For  
 226 exact numbers, see Table 5.

227 In MiniGrid, GBRL outperforms or is on-par with NN in most tasks. Specifically, PPO GBRL is  
 228 significantly better than PPO NN. We observe the same trend when comparing between environments.  
 229 These results emphasize that GBRL is a strong candidate for problems characterized by structured  
 230 data, specifically when using PPO as the algorithmic backend.

231 **GBRL vs Traditional GBT Libraries.** Here, we compare GBRL with Catboost and XGBoost.  
 232 We focus on the PPO variant. When comparing to the standard libraries, we utilize their built-in

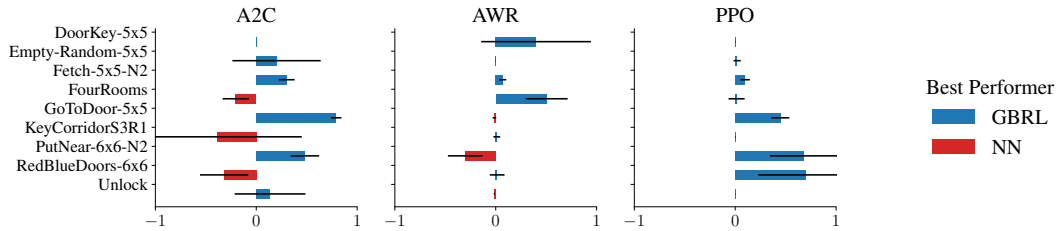


Figure 5: **MiniGrid environments.** GBRL combined with the PPO backend outperforms NN across a range of categorical environments.

233 options for incremental learning, vectorized leaves, and custom loss functions. As both CatBoost and  
 234 XGBoost do not support differential learning rates, we used separate ensembles for the actor and the  
 235 critic. For the comparison, we use the CartPole-v1 environment, training for 1M steps. The results  
 236 are shown in Figure 6.

237 As seen, standard GBT libraries are unable to solve RL tasks in a realistic timeframe. GBRL, however,  
 238 efficiently solves the task while also remaining competitive with NN across a range of environments.

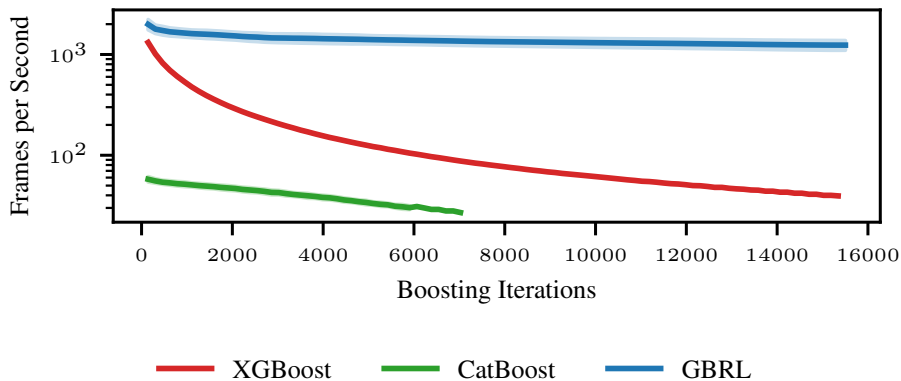


Figure 6: **Comparing to standard GBT libraries.** CatBoost and XGBoost are intractable in RL. Specifically, CatBoost lacks GPU support for custom losses, leading to low FPS and early termination.

239 **Evaluating the shared AC architecture.** Finally, we evaluated the benefits of using a shared AC  
 240 architecture by training PPO GBRL on three MiniGrid environments. We train agents with shared  
 241 and non-shared architectures for 10M timesteps and compare the score, GPU memory usage, and  
 242 FPS. We provide the aggregated results in Figure 7, and environment-specific breakdowns in the  
 243 supplementary.

244 The benefit of the shared structure is clear both in terms of GPU memory consumption and computa-  
 245 tion speed. By sharing the tree structure, GBRL requires less than half the memory and almost triples  
 246 the training FPS. This is achieved without any negative performance on the resulting policy, as seen  
 247 in the reward plot.

248 **Result summary.** The performance of GBRL varied across RL algorithms, but environments like  
 249 MiniGrid highlight the potential advantages of using GBT in RL. The results suggest that GBT’s  
 250 strengths in handling structured and categorical data from supervised learning can effectively transfer  
 251 to the RL domain. Conversely, GBRL underperformed in Atari-RAM environments, indicating that  
 252 certain environments, characterized by unstructured observations, are less suited for GBTs.

253 The results can be explained by the findings of Grinsztajn et al. [15], which suggest that NNs have  
 254 an inductive bias toward overly smooth solutions and that MLP-like architectures are not robust  
 255 to uninformative features. The optimal solutions for Atari-RAM might be smoother, which could  
 256 explain the better performance of NNs. On the other hand, McElfresh et al. [31] argue that GBT  
 257 outperforms NNs on ‘irregular’ datasets. Tree-based models excel in handling irregular patterns and  
 258 categorical data, aligning with GBRL’s success in environments like MiniGrid.



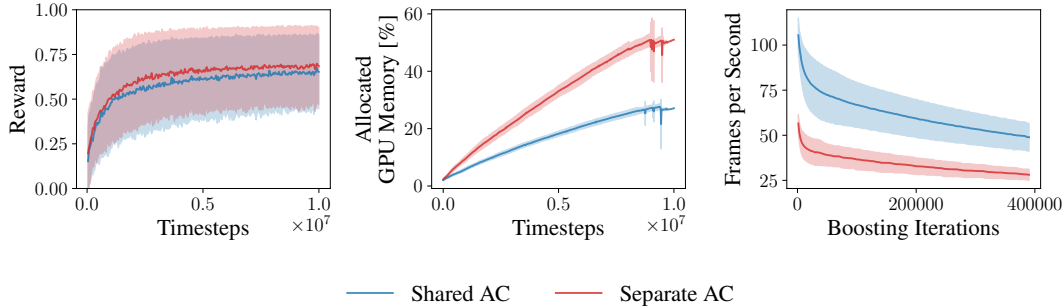


Figure 7: **Shared Actor-Critic**. Sharing the tree structure significantly increases training efficiency and memory, without impacting on the score.

259 Comparing different algorithmic backbones, we find PPO to be the strongest. PPO GBRL excelled  
 260 in the MiniGrid and Football domains, and performed comparably with NN in classic control  
 261 tasks. PPO GBRL’s success can be attributed to its alignment with GBRL’s incremental learning  
 262 strategy. On the other hand, A2C’s single gradient update per rollout may limit its effectiveness and  
 263 contribute to its underwhelming performance in many environments. Similarly, AWR’s design for  
 264 multiple sample updates results in very large ensembles, creating a trade-off between large, slow, and  
 265 memory-intensive ensembles, and lighter, less performant versions.

## 266 6 Conclusion

267 Historically, RL practitioners have relied on tabular, linear, and NN-based function approximators.  
 268 But, GBT, a widely successful tool in supervised learning, has been absent from this toolbox. We  
 269 present a method for effectively integrating it into RL and demonstrate domains where it excels  
 270 compared to NNs. GBRL is a step toward solutions that are more interpretable, well suited for  
 271 real-world tasks with structured data, or capable of deployment on low-compute devices.

272 The choice of an RL method depends on the task characteristics: tabular and linear approaches  
 273 are suitable for small state spaces or simple mappings, while NNs handle complex relationships in  
 274 unstructured data. GBT thrives in complex, yet structured environments. In such cases, we observe  
 275 the advantage of GBRL over NNs, reflecting its already known benefits in supervised learning.

276 A crucial component of GBRL is our efficient adaptation of GBT for AC methods, which allows  
 277 the simultaneous optimization of distinct objectives. We optimized this approach for large-scale  
 278 ensembles using GPU acceleration (CUDA). Furthermore, GBRL integrates seamlessly with existing  
 279 RL libraries, promoting ease of use and adoption.

## 280 7 Limitations and Future Directions

281 In this work, we integrated the highly popular GBT, typically used in supervised learning, into RL.  
 282 Our results show that GBT is competitive across a range of problems. However, we identified several  
 283 limitations and compelling areas for further research. First, a significant challenge lies in the continu-  
 284 ous generation of trees. As the policy improves through numerous updates, the size of the ensemble  
 285 increases. This unbounded growth has implications for memory usage, computational efficiency, and  
 286 the feasibility of online real-time adaptation. The problem is exacerbated by off-policy methods that  
 287 build many trees per sample. Moreover, the redundancy of trees, especially those from early stages,  
 288 suggests that the final policy could be represented with a much smaller ensemble. Consequently,  
 289 developing strategies for tree pruning, ensemble compression, or dynamically managing ensemble  
 290 size could offer crucial optimizations without compromising performance.

291 Another key challenge lies in effectively integrating GBT with additional state-of-the-art RL algo-  
 292 rithms such as DDPG [23] or SAC [17]. These require differentiable Q-functions to update the policy.  
 293 Since GBTs are not differentiable, new solutions are needed to incorporate them into these algorithms.  
 294 One such possible direction can be probabilistic trees, where each node represents the probability of  
 295 traversing the graph.

## References

- 296  
297 [1] D. Abel, A. Agarwal, F. Diaz, A. Krishnamurthy, and R. E. Schapire. Exploratory gradient  
298 boosting for reinforcement learning in complex domains, 2016.
- 299 [2] S. O. Arık and T. Pfister. Tabnet: Attentive interpretable tabular learning, 2020.
- 300 [3] O. Bastani, Y. Pu, and A. Solar-Lezama. Verifiable reinforcement learning via policy extraction,  
301 2019.
- 302 [4] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An  
303 evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279,  
304 June 2013. ISSN 1076-9757. doi: 10.1613/jair.3912. URL [http://dx.doi.org/10.1613/  
305 jair.3912](http://dx.doi.org/10.1613/jair.3912).
- 306 [5] N. Brukhim, E. Hazan, and K. Singh. A boosting approach to reinforcement learning, 2023.
- 307 [6] T. Chen. Machine learning challenge winning solutions, 2023. [https://github.com/dmlc/  
308 xgboost/tree/master/demo#machine-learning-challenge-winning-solutions](https://github.com/dmlc/xgboost/tree/master/demo#machine-learning-challenge-winning-solutions).
- 309 [7] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd  
310 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16.  
311 ACM, Aug. 2016. doi: 10.1145/2939672.2939785. URL [http://dx.doi.org/10.1145/  
312 2939672.2939785](http://dx.doi.org/10.1145/2939672.2939785).
- 313 [8] M. Chevalier-Boisvert, B. Dai, M. Towers, R. de Lazcano, L. Willems, S. Lahlou, S. Pal, P. S.  
314 Castro, and J. Terry. Minigrid & miniworld: Modular & customizable reinforcement learning  
315 environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.
- 316 [9] Q. Delfosse, S. Sztwiertnia, M. Rothermel, W. Stammer, and K. Kersting. Interpretable concept  
317 bottlenecks to align reinforcement learning agents, 2024.
- 318 [10] A. Delgado-Panadero, B. Hernandez-Lorca, M. T. Garcia-Ordas, and J. A. Benitez-Andrades.  
319 Implementing local-explainability in gradient boosting trees: Feature contribution. *Information  
320 Sciences*, 589:199–212, Apr. 2022. ISSN 0020-0255. doi: 10.1016/j.ins.2021.12.111. URL  
321 <http://dx.doi.org/10.1016/j.ins.2021.12.111>.
- 322 [11] T. Duan, A. Avati, D. Y. Ding, K. K. Thai, S. Basu, A. Y. Ng, and A. Schuler. Ngboost: Natural  
323 gradient boosting for probabilistic prediction, 2020.
- 324 [12] J. H. Friedman. Greedy function approximation: A gradient boosting machine. *The Annals of  
325 Statistics*, 29(5):1189 – 1232, 2001. doi: 10.1214/aos/1013203451. URL [https://doi.org/  
326 10.1214/aos/1013203451](https://doi.org/10.1214/aos/1013203451).
- 327 [13] J. H. Friedman. Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38  
328 (4):367–378, 2002. ISSN 0167-9473. doi: [https://doi.org/10.1016/S0167-9473\(01\)00065-2](https://doi.org/10.1016/S0167-9473(01)00065-2).  
329 URL <https://www.sciencedirect.com/science/article/pii/S0167947301000652>.  
330 Nonlinear Methods and Data Mining.
- 331 [14] Y. Gorishniy, I. Rubachev, V. Khrulkov, and A. Babenko. Revisiting deep learning models for  
332 tabular data, 2023.
- 333 [15] L. Grinsztajn, E. Oyallon, and G. Varoquaux. Why do tree-based models still outperform deep  
334 learning on tabular data?, 2022.
- 335 [16] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, D. Pedreschi, and F. Giannotti. A survey of  
336 methods for explaining black box models, 2018.
- 337 [17] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy  
338 deep reinforcement learning with a stochastic actor, 2018.
- 339 [18] S. Ivanov and L. Prokhorenkova. Boost then convolve: Gradient boosting meets graph neural  
340 networks, 2021.

- 341 [19] L. Katzir, G. Elidan, and R. El-Yaniv. Net-dnf: Effective deep modeling of tabular data. In  
342 *International Conference on Learning Representations*, 2021. URL [https://openreview.](https://openreview.net/forum?id=73WTGs96kho)  
343 [net/forum?id=73WTGs96kho](https://openreview.net/forum?id=73WTGs96kho).
- 344 [20] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu.  
345 Lightgbm: A highly efficient gradient boosting decision tree. In I. Guyon, U. V.  
346 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, edi-  
347 tors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,  
348 Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/](https://proceedings.neurips.cc/paper_files/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf)  
349 [6449f44a102fde848669bdd9eb6b76fa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf).
- 350 [21] K. Kersting and K. Driessens. Non-parametric policy gradients: a unified treatment of  
351 propositional and relational domains. In *Proceedings of the 25th International Conference*  
352 *on Machine Learning*, ICML '08, page 456–463, New York, NY, USA, 2008. Association  
353 for Computing Machinery. ISBN 9781605582054. doi: 10.1145/1390156.1390214. URL  
354 <https://doi.org/10.1145/1390156.1390214>.
- 355 [22] K. Kurach, A. Raichuk, P. Stańczyk, M. Zajac, O. Bachem, L. Espeholt, C. Riquelme, D. Vincent,  
356 M. Michalski, O. Bousquet, and S. Gelly. Google research football: A novel reinforcement  
357 learning environment, 2020.
- 358 [23] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.  
359 Continuous control with deep reinforcement learning, 2019.
- 360 [24] B. London, L. Lu, T. Sandler, and T. Joachims. Boosted off-policy learning. In F. Ruiz, J. Dy,  
361 and J.-W. van de Meent, editors, *Proceedings of The 26th International Conference on Artificial*  
362 *Intelligence and Statistics*, volume 206, pages 5614–5640, 2023.
- 363 [25] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb,  
364 N. Bansal, and S.-I. Lee. From local explanations to global understanding with explainable ai  
365 for trees. *Nature Machine Intelligence*, 2(1):2522–5839, 2020.
- 366 [26] I. Lyzhin, A. Ustimenko, A. Gulin, and L. Prokhorenkova. Which tricks are important for  
367 learning to rank?, 2023.
- 368 [27] H. Ma, J. Cao, Y. Fang, W. Zhang, W. Sheng, S. Zhang, and Y. Yu. Retrieval-based gradient  
369 boosting decision trees for disease risk assessment. In *Proceedings of the 28th ACM SIGKDD*  
370 *Conference on Knowledge Discovery and Data Mining*, KDD '22, page 3468–3476, New  
371 York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393850. doi:  
372 [10.1145/3534678.3539052](https://doi.org/10.1145/3534678.3539052). URL <https://doi.org/10.1145/3534678.3539052>.
- 373 [28] A. Malinin, L. Prokhorenkova, and A. Ustimenko. Uncertainty in gradient boosting via  
374 ensembles, 2021.
- 375 [29] L. Mason, J. Baxter, P. Bartlett, and M. Frean. Boosting algorithms as gradient descent. In  
376 S. Solla, T. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems*,  
377 volume 12. MIT Press, 1999. URL [https://proceedings.neurips.cc/paper/1999/](https://proceedings.neurips.cc/paper/1999/file/96a93ba89a5b5c6c226e49b88973f46e-Paper.pdf)  
378 [file/96a93ba89a5b5c6c226e49b88973f46e-Paper.pdf](https://proceedings.neurips.cc/paper/1999/file/96a93ba89a5b5c6c226e49b88973f46e-Paper.pdf).
- 379 [30] L. Mason, J. Baxter, P. Bartlett, and M. Frean. Boosting algorithms as gradient descent. In  
380 S. Solla, T. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems*,  
381 volume 12. MIT Press, 1999. URL [https://proceedings.neurips.cc/paper\\_files/](https://proceedings.neurips.cc/paper_files/paper/1999/file/96a93ba89a5b5c6c226e49b88973f46e-Paper.pdf)  
382 [paper/1999/file/96a93ba89a5b5c6c226e49b88973f46e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1999/file/96a93ba89a5b5c6c226e49b88973f46e-Paper.pdf).
- 383 [31] D. McElfresh, S. Khandagale, J. Valverde, V. Prasad C, G. Ramakrishnan, M. Goldblum,  
384 and C. White. When do neural nets outperform boosted trees on tabular data? In A. Oh,  
385 T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neu-*  
386 *ral Information Processing Systems*, volume 36, pages 76336–76369. Curran Associates,  
387 Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/f06d5ebd4ff40b40dd97e30cee632123-Paper-Datasets_and_Benchmarks.pdf)  
388 [f06d5ebd4ff40b40dd97e30cee632123-Paper-Datasets\\_and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/f06d5ebd4ff40b40dd97e30cee632123-Paper-Datasets_and_Benchmarks.pdf).
- 389 [32] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and  
390 K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. 2016. doi:  
391 [10.48550/ARXIV.1602.01783](https://arxiv.org/abs/1602.01783). URL <https://arxiv.org/abs/1602.01783>.

- 392 [33] NVIDIA, P. Vingelmann, and F. H. Fitzek. Cuda, release: 10.2.89, 2020. URL <https://developer.nvidia.com/cuda-toolkit>.  
393
- 394 [34] K. Ota, D. K. Jha, and A. Kanezaki. Training larger networks for deep reinforcement learning,  
395 2021.
- 396 [35] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and  
397 scalable off-policy reinforcement learning, 2019.
- 398 [36] L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Drogush, and A. Gulin. Catboost: unbiased  
399 boosting with categorical features, 2019.
- 400 [37] Y. Qing, S. Liu, J. Song, H. Wang, and M. Song. A survey on explainable reinforcement  
401 learning: Concepts, algorithms, challenges, 2023.
- 402 [38] A. Raffin. RL baselines3 zoo. <https://github.com/DLR-RM/rl-baselines3-zoo>, 2020.
- 403 [39] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann. Stable-baselines3:  
404 Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22  
405 (268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.
- 406 [40] S. Ruder. An overview of gradient descent optimization algorithms, 2017.
- 407 [41] B. Scherrer and M. Geist. Local policy search in a convex space and conservative policy  
408 iteration as boosted policy search. In T. Calders, F. Esposito, E. Hüllermeier, and R. Meo,  
409 editors, *Machine Learning and Knowledge Discovery in Databases*, pages 35–50, Berlin,  
410 Heidelberg, 2014. Springer Berlin Heidelberg. ISBN 978-3-662-44845-8.
- 411 [42] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
412 algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- 413 [43] H. Seto, A. Oyama, S. Kitora, H. Toki, R. Yamamoto, J. Kotoku, A. Haga, M. Shinzawa,  
414 M. Yamakawa, S. Fukui, and T. Moriyama. Gradient boosting decision tree becomes more  
415 reliable than logistic regression in predicting probability for diabetes with big data. *Scientific*  
416 *Reports*, 12(1):15889, Oct. 2022.
- 417 [44] A. Shrikumar, P. Greenside, and A. Kundaje. Learning important features through propagating  
418 activation differences, 2019.
- 419 [45] F. Sigrist. Gaussian process boosting, 2022.
- 420 [46] G. Somepalli, M. Goldblum, A. Schwarzschild, C. B. Bruss, and T. Goldstein. Saint: Improved  
421 neural networks for tabular data via row attention and contrastive pre-training, 2021.
- 422 [47] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second  
423 edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>.
- 424 [48] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods  
425 for reinforcement learning with function approximation. In S. Solla, T. Leen, and  
426 K. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12. MIT  
427 Press, 1999. URL [https://proceedings.neurips.cc/paper\\_files/paper/1999/  
428 file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf).
- 429 [49] Z. Tian, J. Xiao, H. Feng, and Y. Wei. Credit risk assessment based on gradient boosting decision  
430 tree. *Procedia Computer Science*, 174:150–160, 2020. ISSN 1877-0509. doi: [https://doi.org/10.  
431 1016/j.procs.2020.06.070](https://doi.org/10.1016/j.procs.2020.06.070). URL [https://www.sciencedirect.com/science/article/  
432 pii/S1877050920315842](https://www.sciencedirect.com/science/article/pii/S1877050920315842). 2019 International Conference on Identification, Information and  
433 Knowledge in the Internet of Things.
- 434 [50] M. Towers, J. K. Terry, A. Kwiatkowski, J. U. Balis, G. d. Cola, T. Deleu, M. Goulão,  
435 A. Kallinteris, A. KG, M. Krimmel, R. Perez-Vicente, A. Pierré, S. Schulhoff, J. J. Tai, A. T. J.  
436 Shen, and O. G. Younis. Gymnasium, Mar. 2023. URL [https://zenodo.org/record/  
437 8127025](https://zenodo.org/record/8127025).

- 438 [51] A. Ustimenko and L. Prokhorenkova. Stochasticrank: Global optimization of scale-free discrete  
439 functions, 2020.
- 440 [52] A. Ustimenko, A. Beliakov, and L. Prokhorenkova. Gradient boosting performs gaussian  
441 process inference, 2023.
- 442 [53] K. Wang, J. Lu, A. Liu, G. Zhang, and L. Xiong. Evolving gradient boost: A pruning scheme  
443 based on loss improvement ratio for learning under concept drift. *IEEE Trans. Cybern.*, 53(4):  
444 2110–2123, Apr. 2023.
- 445 [54] S. Wassan, B. Suhail, R. Mubeen, B. Raj, U. Agarwal, E. Khatri, S. Gopinathan, and G. Dhi-  
446 man. Gradient boosting for health iot federated learning. *Sustainability*, 14(24), 2022. ISSN  
447 2071-1050. doi: 10.3390/su142416842. URL [https://www.mdpi.com/2071-1050/14/24/](https://www.mdpi.com/2071-1050/14/24/16842)  
448 16842.
- 449 [55] C. Zhang, Y. Zhang, X. Shi, G. Almpanidis, G. Fan, and X. Shen. On incremental learning  
450 for gradient boosting decision trees. *Neural Processing Letters*, 50(1):957–987, Aug 2019.  
451 ISSN 1573-773X. doi: 10.1007/s11063-019-09999-3. URL [https://doi.org/10.1007/](https://doi.org/10.1007/s11063-019-09999-3)  
452 s11063-019-09999-3.

## 453 Appendix

454 This appendix provides supplementary materials that support the findings and methodologies dis-  
455 cussed in the main text. It is organized into four sections to present the full experiment results,  
456 implementation details, hyperparameters used during the experiments, training progression plots, and  
457 experimental plots, respectively. These materials offer detailed insights into the research process and  
458 outcomes, facilitating a deeper understanding and replication of the study.

## 459 A Implementaion Details and Hyperparameters

460 Included in this section are implementation details, information regarding compute resources, and  
461 tables containing the hyperparameters used in our experiments enabling the reproducibility of our  
462 results. Table 1 lists GBRL hyperparameters for all experiments.

### 463 A.1 Environments

464 The Football domain consists of a vectorized 115-dimensional observation space that summarizes the  
465 main aspects of the game and 19 discrete actions. We focus on its academy scenarios, which present  
466 situational tasks involving scoring a single goal. A standard reward of +1 is granted for scoring, and  
467 we employed the "Checkpoints" shaped reward structure. This structure provides additional points as  
468 the agent moves closer towards the goal, with a maximum reward of 2 per scenario. The Atari-ram  
469 environment consists of a vectorized 128-dimensional observational space representing the 128 byte  
470 RAM state and up to 18 discrete actions. We trained agents in both domains for 10M timesteps.

471 The MiniGrid environment [8] is a 2D grid world with goal-oriented tasks requiring object interaction.  
472 The observation space consists of a 7x7 image representing the grid, a mission string, and the agent's  
473 direction. Each tile in the observed image contains a 3D tuple dictating an object's color, type, and  
474 state. All MiniGrid tasks emit a reward of +1 for successful completion and 0 otherwise.

475 We trained our NN-based agents on a flattened observation space using the built-in one-hot wrapper.  
476 For GBRL agents, we generated a 51-dimensional categorical observational space by encoding each  
477 unique tile tuple as a categorical string to represent the observed image. Categorical features were  
478 added for the agent's direction (up, left, right, down) and missions. All agents were trained for 1M  
479 timesteps, except for PutNear, FourRooms, and Fetch tasks, which were trained for 10M based on the  
480 reported values for PPO NN in RL Baselines3 Zoo.

### 481 A.2 Compute Resources

482 All experiments were done on the NVIDIA NGC platform on a single NVIDIA V100-32GB GPU  
483 per experiment. Training time and compute requirements vary between algorithms and according  
484 to hyperparameters. The number of boosting iterations has the largest impact on both runtime and  
485 memory. GBRL experimental runs required from 1GB to 24GB of GPU memory. Moreover, runtime  
486 varied from 20 minutes for 1M timesteps training on classic environments and up to 5 days for 10M  
487 timesteps on Atari-ram. NN experimental runs required up to 3GB of GPU memory and runtime  
488 ranged from 10 minutes and up to 3 days. The total compute time for all experiments combined was  
489 approximately 1800 GPU hours. Additionally, the research project involved preliminary experiments  
490 and hyperparameter tuning, which required an estimated additional 168 GPU hours.

## 491 B Detailed Results Tables

492 This section contains tables presenting the mean and standard deviation of the average episode reward  
493 for the final 100 episodes within each experiment. More specifically, Table 2 presents results for  
494 Continuous Control & Block2D environments, Tables 3 and 4 present results for the high-dimensional  
495 vectorized environments, and Table 5 presents results for the categorical environments.

	batch size	clip range	ent coef	gae lambda	gamma	num epochs	num steps	num envs	policy lr	value lr
Acrobot	512	0.2	0.0	0.94	0.99	20	128	16	0.16	0.034
CartPole	64	0.2	0.0	0.8	0.98	1	128	8	0.029	0.015
LunarLander	256	0.2	0.0033	0.98	0.999	20	512	16	0.031	0.003
MountainCar	256	0.2	0.033	0.98	0.999	20	512	16	0.031	0.003
MountainCar Continuous	256	0.2	0.033	0.98	0.999	20	512	16	0.031	0.003
Pendulum	512	0.2	0.0	0.93	0.91	20	256	16	0.031	0.013
Football	512	0.2	0.0	0.95	0.998	10	256	16	0.033	0.006
Atari-Ram	64	0.92	8e-5	0.95	0.99	4	512	16	0.05	0.002
MiniGrid	512	0.2	0.0	0.95	0.99	20	256	16	0.17	0.01

(a) PPO. For continuous action spaces we used log std init = -2 and log std lr = lin\_0.0017. We utilized gradient norm clipping for Gym environments. Specifically, 10 for the value gradients and 150 for the policy gradients.

	ent coef	gae lambda	gamma	num steps	num envs	policy lr	value lr	log std init	log std lr
Acrobot	0.0	1	0.99	8	4	0.79	0.031	-	-
CartPole	0.0	1	0.99	8	16	0.13	0.047	-	-
LunarLander	0.0	1	0.995	5	32	0.16	0.04	-	-
MountainCar	0.0	1	0.99	8	16	0.64	0.032	-	-
MountainCar Continuous	0.0	1	0.995	128	16	0.0008	2.8e-6	0	0.0004
Pendulum	0.0	0.9	0.9	10	32	0.003	0.056	-2	0.00018
Football	0.0004	0.95	0.998	128	8	0.87	0.017	-	-
Atari-Ram	0.0009	0.95	0.993	128	8	0.17	0.013	-	-
MiniGrid	0.0	0.95	0.99	10	128	0.34	0.039	-	-

(b) A2C

	batch size	ent coef	gae lambda	gamma	train freq	gradient steps	num envs	policy lr	value lr	log std init	log std lr
Acrobot	1024	0.0	0.95	0.99	2000	150	1	0.05	0.1	-	-
CartPole	1024	0.0	0.95	0.99	2000	150	1	0.05	0.1	-	-
LunarLander	1024	0.0	0.95	0.99	2000	150	1	0.05	0.1	-	-
MountainCar	64	0.0	0.95	0.99	2000	150	1	0.64	0.032	-	-
MountainCar Continuous	64	0.0	0.95	0.99	2000	150	1	0.089	0.083	-2	lin_0.0017
Pendulum	1024	0.0	0.9	0.9	1000	50	1	0.003	0.07	-2	0.0005
Football	512	0.03	0.95	0.99	750	10	1	0.09	0.00048	-	-
Atari-Ram	1024	0.0	0.95	0.993	2000	50	1	0.0779	0.0048	-	-
MiniGrid	1024	0.0	0.95	0.99	1500	25/100*	1	0.0075	0.005	-	-

(c) AWR. For all envs, buffer size = 50,000,  $\beta = 0.05$ . \*MiniGrid environments used 100 gradient steps for tasks trained for 1M steps, and 25 gradient steps for tasks trained for 10M steps, for a reduced tree size.

Table 1: **GBRL hyperparameters** - NN represented by an MLP with two hidden layers.

Table 2: Continuous-Control and Box2D environments: Average episode reward for the final 100 episodes.

	Acrobot	CartPole	LunarLander	MountainCar	MountainCar Continuous	Pendulum-v1
NN: A2C	-82.27 ± 3.29	500.00 ± 0.0	-43.01 ± 106.26	-148.90 ± 24.10	92.66 ± 0.32	-183.64 ± 22.32
GBRL: A2C	-90.73 ± 2.98	500.00 ± 0.0	<b>47.93 ± 41.00</b>	-124.42 ± 5.74	93.15 ± 1.19	-538.83 ± 66.25
NN: AWR	-102.53 ± 57.25	500.00 ± 0.0	<b>282.48 ± 1.96</b>	-160.65 ± 53.97	18.93 ± 42.34	-159.64 ± 9.42
GBRL: AWR	-118.12 ± 33.54	497.54 ± 3.11	76.03 ± 56.62	-146.68 ± 24.53	44.38 ± 45.94	-1257.61 ± 98.10
NN: PPO	-74.83 ± 1.22	500.00 ± 0.0	261.73 ± 6.93	-115.53 ± 1.39	85.81 ± 7.51	-249.31 ± 60.00
GBRL: PPO	-87.82 ± 2.16	500.00 ± 0.0	248.72 ± 59.10	-110.55 ± 15.60	89.42 ± 5.73	-246.89 ± 20.61

## 496 C Training Plots

497 This section presents learning curves depicting model performance throughout the training phase.  
 498 Figures 8 to 11 show the training reward as a function of environment steps of the agents trained in  
 499 the experiments. The column order is: A2C, AWR, and PPO.

Table 3: Football Academy environments: Average episode reward for the final 100 episodes.

	3 vs 1 with keeper	Corner	Counterattack Easy	Counterattack Hard	Empty Goal	Empty Goal Close
NN: A2C	1.78 ± 0.10	1.00 ± 0.17	<b>1.58 ± 0.35</b>	<b>1.43 ± 0.17</b>	<b>1.93 ± 0.05</b>	2.0 ± 0.0
GBRL: A2C	1.59 ± 0.17	1.01 ± 0.07	1.11 ± 0.14	1.00 ± 0.05	1.81 ± 0.03	2.00 ± 0.00
NN: AWR	1.50 ± 0.37	1.01 ± 0.04	<b>1.59 ± 0.36</b>	1.18 ± 0.21	1.90 ± 0.08	1.92 ± 0.17
GBRL: AWR	1.66 ± 0.34	0.92 ± 0.05	0.95 ± 0.05	0.92 ± 0.05	1.93 ± 0.07	2.0 ± 0.0
NN: PPO	1.61 ± 0.05	0.95 ± 0.02	1.43 ± 0.15	1.23 ± 0.18	<b>1.98 ± 0.01</b>	1.99 ± 0.00
GBRL: PPO	1.63 ± 0.19	1.05 ± 0.20	1.64 ± 0.09	1.23 ± 0.07	1.84 ± 0.06	2.0 ± 0.0

	Pass & Shoot keeper	Run Pass & Shoot keeper	Run to Score	Run to score w/ keeper	Single Goal vs Lazy
NN: A2C	1.41 ± 0.37	1.77 ± 0.08	1.87 ± 0.12	1.25 ± 0.23	<b>1.65 ± 0.04</b>
GBRL: A2C	1.60 ± 0.21	1.60 ± 0.14	1.82 ± 0.10	1.15 ± 0.08	1.31 ± 0.11
NN: AWR	1.26 ± 0.46	1.15 ± 0.14	1.81 ± 0.14	1.25 ± 0.34	1.28 ± 0.27
GBRL: AWR	1.35 ± 0.37	1.53 ± 0.40	<b>1.98 ± 0.01</b>	0.99 ± 0.16	1.03 ± 0.12
NN: PPO	1.31 ± 0.13	1.64 ± 0.16	1.91 ± 0.09	1.13 ± 0.06	1.68 ± 0.09
GBRL: PPO	<b>1.87 ± 0.09</b>	1.85 ± 0.08	1.83 ± 0.04	<b>1.95 ± 0.02</b>	1.73 ± 0.06

Table 4: Atari-ramNoFrameskip-v4 environments: Average episode reward for the final 100 episodes.

	Alien	Amidar	Asteroids	Breakout	Gopher
NN: A2C	<b>1802.24 ± 323.12</b>	<b>304.62 ± 55.61</b>	<b>2770.46 ± 271.97</b>	<b>76.69 ± 30.08</b>	<b>3533.84 ± 118.50</b>
GBRL: A2C	595.08 ± 43.51	48.71 ± 14.65	1402.66 ± 161.67	11.52 ± 2.34	502.20 ± 341.88
NN: AWR	739.82 ± 303.06	86.32 ± 40.16	<b>2308.68 ± 257.72</b>	26.57 ± 9.91	1471.93 ± 716.65
GBRL: AWR	829.99 ± 166.48	125.53 ± 25.25	1592.63 ± 109.96	17.32 ± 1.89	913.06 ± 79.95
NN: PPO	<b>1555.32 ± 107.59</b>	<b>310.93 ± 80.13</b>	<b>2309.46 ± 145.66</b>	<b>32.88 ± 15.74</b>	<b>2507.84 ± 108.37</b>
GBRL: PPO	1163.86 ± 76.54	186.32 ± 50.63	1514.34 ± 317.46	19.96 ± 1.93	1215.04 ± 81.01

	Kangaroo	Krull	MsPacman	Pong	SpaceInvaders
NN: A2C	<b>2137.6 ± 425.64</b>	<b>9325.38 ± 777.12</b>	<b>2007.64 ± 116.52</b>	<b>15.39 ± 4.26</b>	<b>462.30 ± 35.56</b>
GBRL: A2C	948.8 ± 483.80	5291.4 ± 433.35	989.68 ± 100.02	-12.80 ± 11.10	265.36 ± 44.64
NN: AWR	1214.8 ± 313.42	4519.78 ± 522.11	892.31 ± 289.36	-10.25 ± 2.11	842.00 ± 130.51
GBRL: AWR	<b>1809.26 ± 37.51</b>	<b>6419.26 ± 387.76</b>	<b>1641.84 ± 284.19</b>	-11.68 ± 3.79	397.85 ± 566.38
NN: PPO	2487.4 ± 829.65	9167.3 ± 294.30	2069.22 ± 202.48	18.50 ± 1.60	479.77 ± 65.07
GBRL: PPO	2160.8 ± 826.92	6888.66 ± 756.18	2069.22 ± 538.62	15.40 ± 6.55	434.84 ± 31.83

Table 5: MiniGrid environments: Average episode reward for the final 100 episodes.

	DoorKey-5x5	Empty-Random-5x5	Fetch-5x5-N2	FourRooms	GoToDoor-5x5
NN: A2C	0.96 ± 0.00	0.77 ± 0.42	0.43 ± 0.03	0.62 ± 0.19	0.05 ± 0.04
GBRL: A2C	0.96 ± 0.00	0.96 ± 0.00	0.62 ± 0.02	0.51 ± 0.07	0.78 ± 0.02
NN: AWR	0.57 ± 0.52	0.96 ± 0.00	0.90 ± 0.26	0.19 ± 0.12	0.95 ± 0.01
GBRL: AWR	<b>0.96 ± 0.00</b>	0.97 ± 0.00	<b>0.95 ± 0.01</b>	<b>0.54 ± 0.05</b>	0.94 ± 0.01
NN: PPO	0.78 ± 0.40	0.96 ± 0.00	0.89 ± 0.03	0.53 ± 0.03	0.60 ± 0.06
GBRL: PPO	0.96 ± 0.00	0.96 ± 0.00	<b>0.96 ± 0.01</b>	0.56 ± 0.04	<b>0.96 ± 0.00</b>

	KeyCorridorS3R1	PutNear-6x6-N2	RedBlueDoors-6x6	Unlock
NN: A2C	0.75 ± 0.42	0.01 ± 0.00	<b>0.30 ± 0.22</b>	0.77 ± 0.43
GBRL: A2C	0.39 ± 0.48	<b>0.18 ± 0.018</b>	0.0 ± 0.0	0.90 ± 0.09
NN: AWR	0.93 ± 0.00	<b>0.60 ± 0.13</b>	0.83 ± 0.00	0.96 ± 0.00
GBRL: AWR	0.94 ± 0.00	0.36 ± 0.01	0.84 ± 0.03	0.95 ± 0.00
NN: PPO	0.76 ± 0.42	0.001 ± 0.00	0.17 ± 0.40	0.97 ± 0.00
GBRL: PPO	<b>0.95 ± 0.00</b>	<b>0.44 ± 0.19</b>	<b>0.88 ± 0.02</b>	0.97 ± 0.00



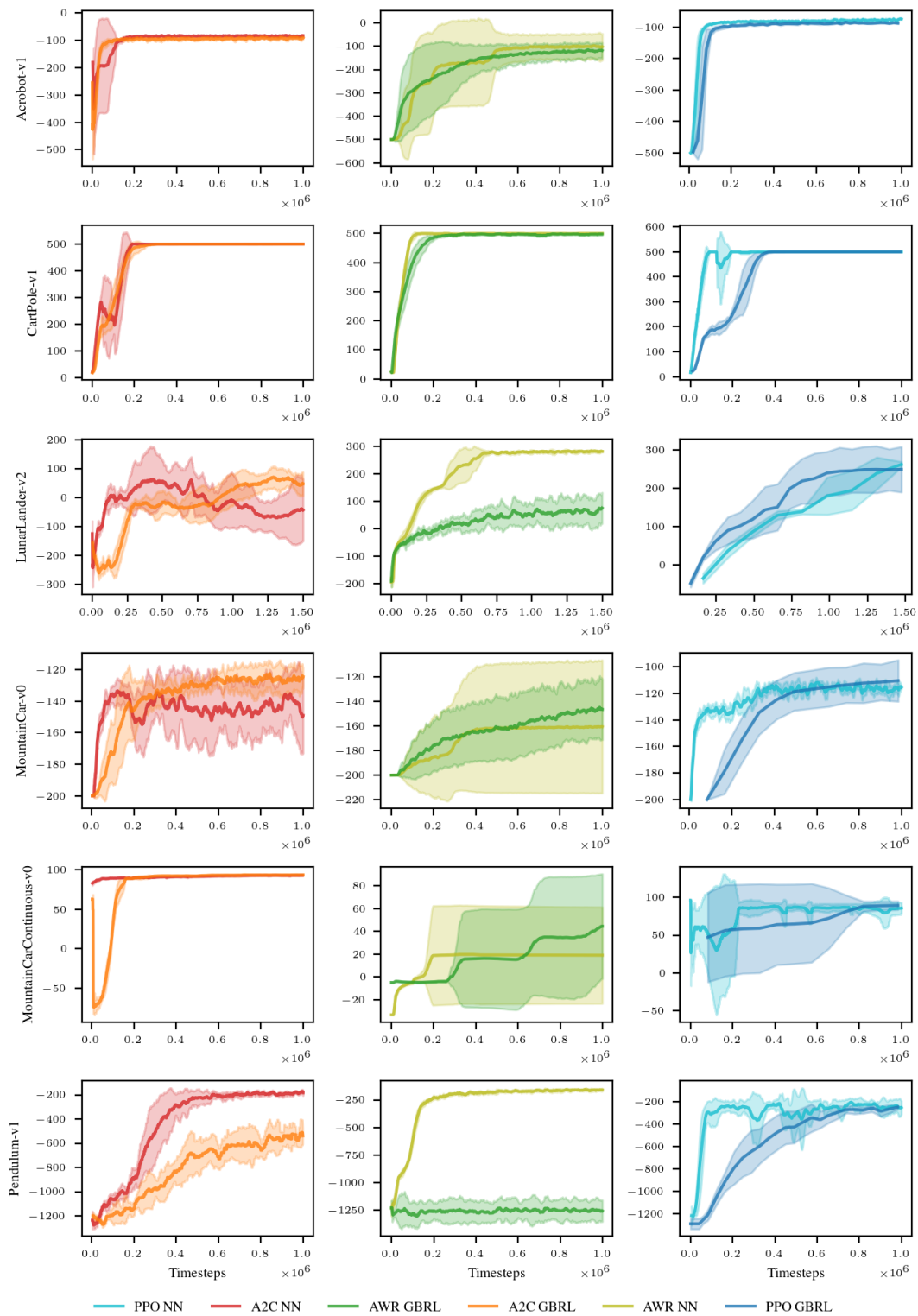


Figure 8: Classic Control and Box2D environments: Training reward as a function of environment steps.

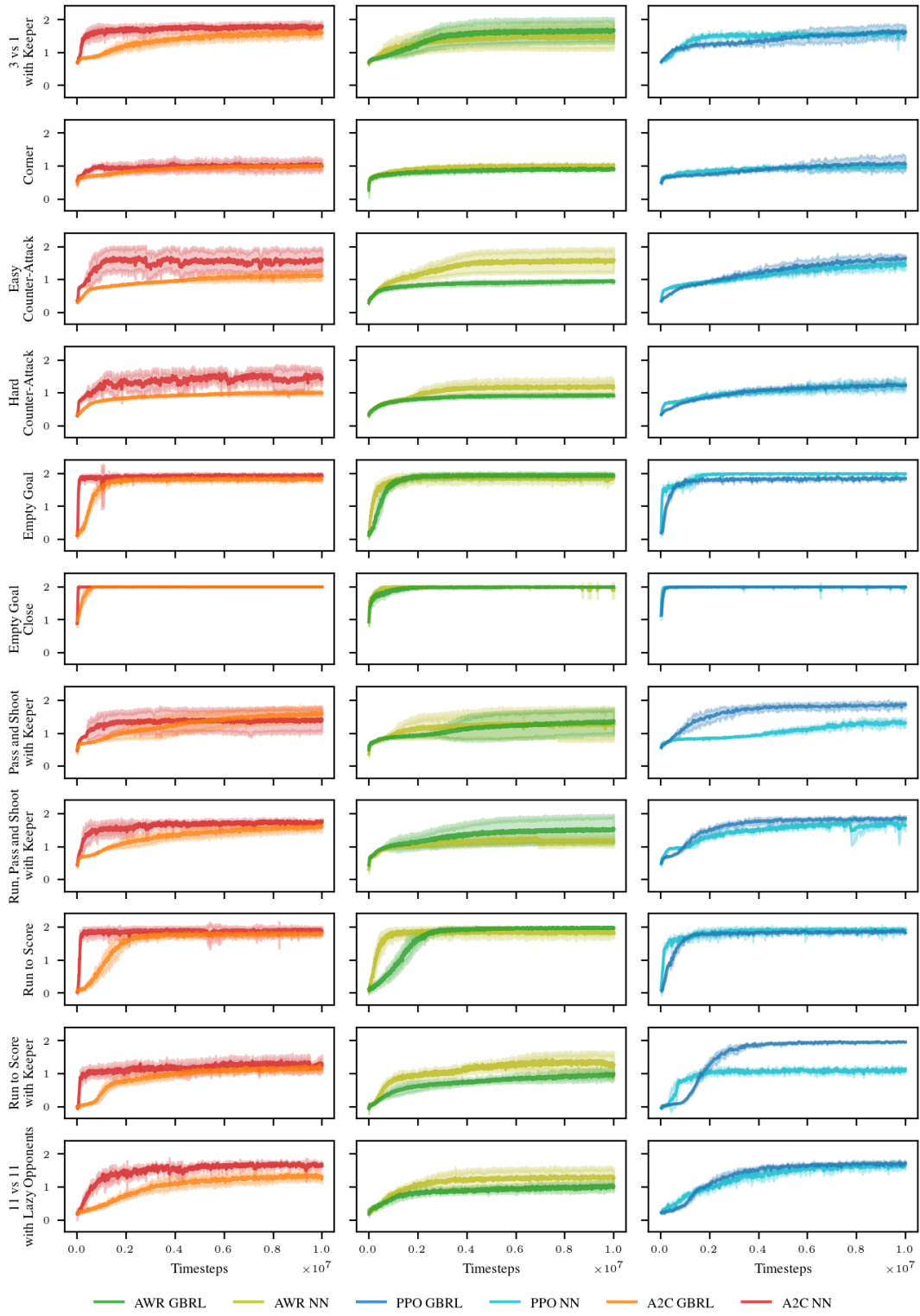


Figure 9: Football Academy environments: Training reward as a function of environment step.

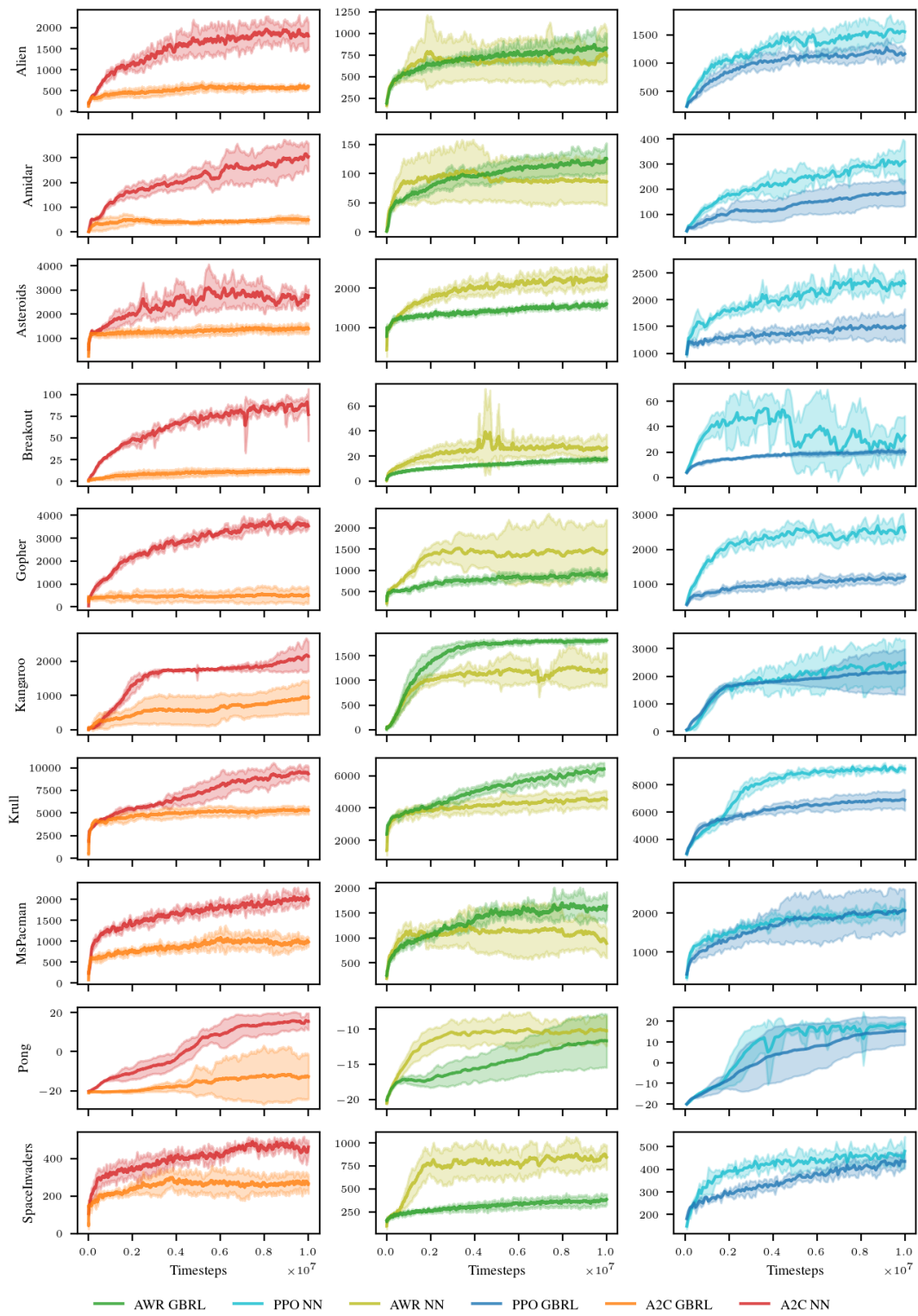


Figure 10: Atari-ramNoFrameskip-v4 environments: Training reward as a function of environment step.

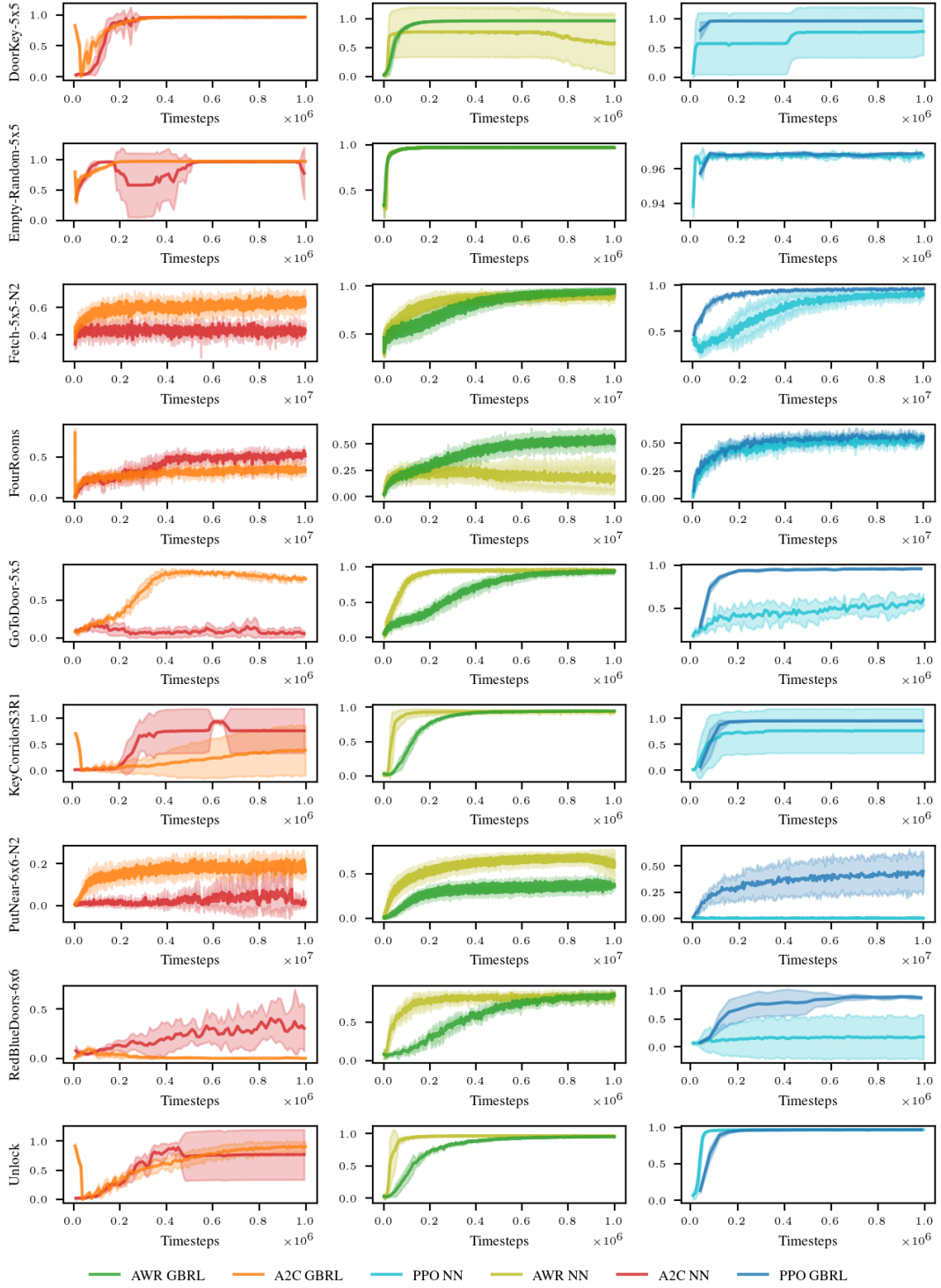


Figure 11: MiniGrid environments: Training reward as a function of environment step.

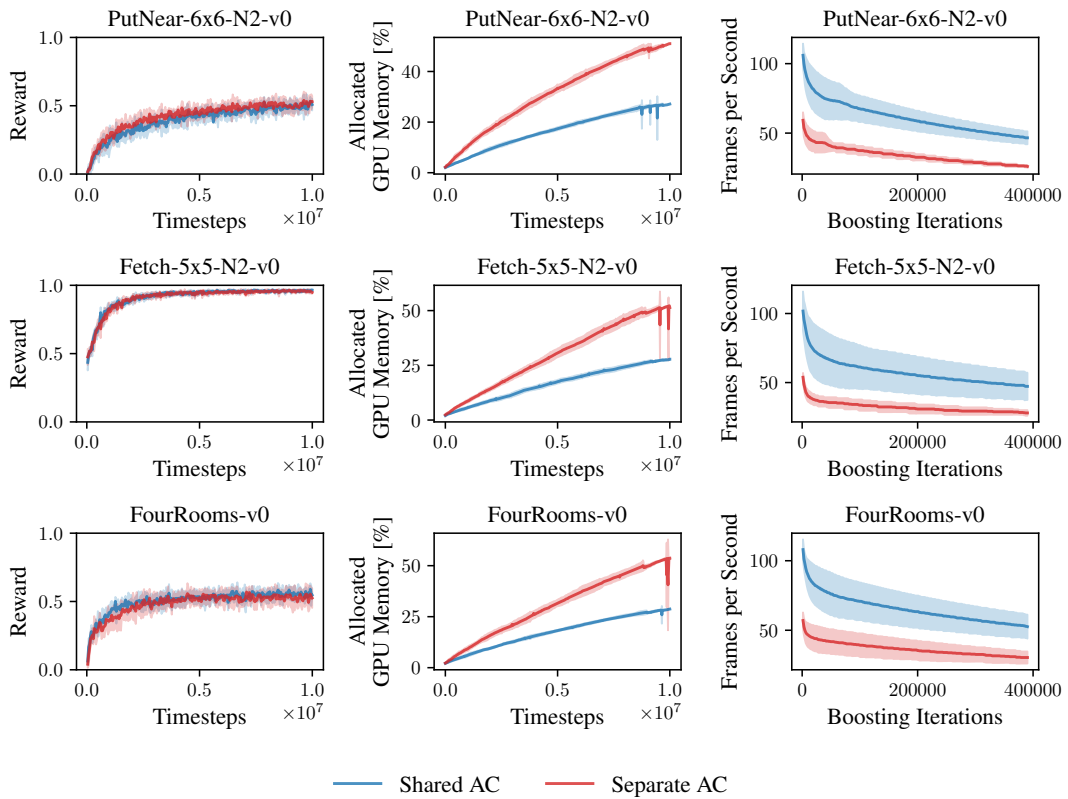


Figure 12: **Sharing actor critic tree structure significantly increases efficiency while retraining similar performance.** Training reward, GPU memory usage, and FPS, are compared across 10M environment (5 seeds, 3 MiniGrid environments)

## 500 **NeurIPS Paper Checklist**

### 501 **1. Claims**

502 Question: Do the main claims made in the abstract and introduction accurately reflect the  
503 paper's contributions and scope?

504 Answer: [\[Yes\]](#)

505 Justification: The main claims made in the abstract and introduction accurately reflect the  
506 paper's contributions and scope. The abstract provides a concise overview of the problem,  
507 the proposed solution, and the key contributions, including the introduction of GBRL, which  
508 extends the advantages of GBT to the RL domain. The paper demonstrates competitive  
509 performance with NN-based methods, especially in domains with structured or categorical  
510 features, and presents a high-performance, GPU-accelerated implementation that can be  
511 integrated with RL libraries. While the paper includes some aspirational goals related to  
512 future applications, these are clearly distinguished from the results shown and serve as  
513 motivation for further research..

514 Guidelines:

- 515 • The answer NA means that the abstract and introduction do not include the claims  
516 made in the paper.
- 517 • The abstract and/or introduction should clearly state the claims made, including the  
518 contributions made in the paper and important assumptions and limitations. A No or  
519 NA answer to this question will not be perceived well by the reviewers.
- 520 • The claims made should match theoretical and experimental results, and reflect how  
521 much the results can be expected to generalize to other settings.
- 522 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
523 are not attained by the paper.

### 524 **2. Limitations**

525 Question: Does the paper discuss the limitations of the work performed by the authors?

526 Answer: [\[Yes\]](#)

527 Justification: We have included a limitations section, in which we discuss the computational  
528 efficiency and the memory limits of our current method. Additionally, we discuss the  
529 challenges in implementing our methods to other popular RL algorithms.

530 Guidelines:

- 531 • The answer NA means that the paper has no limitation while the answer No means that  
532 the paper has limitations, but those are not discussed in the paper.
- 533 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 534 • The paper should point out any strong assumptions and how robust the results are to  
535 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
536 model well-specification, asymptotic approximations only holding locally). The authors  
537 should reflect on how these assumptions might be violated in practice and what the  
538 implications would be.
- 539 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
540 only tested on a few datasets or with a few runs. In general, empirical results often  
541 depend on implicit assumptions, which should be articulated.
- 542 • The authors should reflect on the factors that influence the performance of the approach.  
543 For example, a facial recognition algorithm may perform poorly when image resolution  
544 is low or images are taken in low lighting. Or a speech-to-text system might not be  
545 used reliably to provide closed captions for online lectures because it fails to handle  
546 technical jargon.
- 547 • The authors should discuss the computational efficiency of the proposed algorithms  
548 and how they scale with dataset size.
- 549 • If applicable, the authors should discuss possible limitations of their approach to  
550 address problems of privacy and fairness.

551 • While the authors might fear that complete honesty about limitations might be used by  
552 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
553 limitations that aren't acknowledged in the paper. The authors should use their best  
554 judgment and recognize that individual actions in favor of transparency play an impor-  
555 tant role in developing norms that preserve the integrity of the community. Reviewers  
556 will be specifically instructed to not penalize honesty concerning limitations.

### 557 3. Theory Assumptions and Proofs

558 Question: For each theoretical result, does the paper provide the full set of assumptions and  
559 a complete (and correct) proof?

560 Answer: [NA]

561 Justification: We do not include theoretical results.

562 Guidelines:

- 563 • The answer NA means that the paper does not include theoretical results.
- 564 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
565 referenced.
- 566 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 567 • The proofs can either appear in the main paper or the supplemental material, but if  
568 they appear in the supplemental material, the authors are encouraged to provide a short  
569 proof sketch to provide intuition.
- 570 • Inversely, any informal proof provided in the core of the paper should be complemented  
571 by formal proofs provided in appendix or supplemental material.
- 572 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 573 4. Experimental Result Reproducibility

574 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
575 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
576 of the paper (regardless of whether the code and data are provided or not)?

577 Answer: [Yes]

578 Justification: The paper fully discloses the architecture, the GBRL method, hyperparameters  
579 used, and the code will be released publicly.

580 Guidelines:

- 581 • The answer NA means that the paper does not include experiments.
- 582 • If the paper includes experiments, a No answer to this question will not be perceived  
583 well by the reviewers: Making the paper reproducible is important, regardless of  
584 whether the code and data are provided or not.
- 585 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
586 to make their results reproducible or verifiable.
- 587 • Depending on the contribution, reproducibility can be accomplished in various ways.  
588 For example, if the contribution is a novel architecture, describing the architecture fully  
589 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
590 be necessary to either make it possible for others to replicate the model with the same  
591 dataset, or provide access to the model. In general, releasing code and data is often  
592 one good way to accomplish this, but reproducibility can also be provided via detailed  
593 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
594 of a large language model), releasing of a model checkpoint, or other means that are  
595 appropriate to the research performed.
- 596 • While NeurIPS does not require releasing code, the conference does require all submis-  
597 sions to provide some reasonable avenue for reproducibility, which may depend on the  
598 nature of the contribution. For example
  - 599 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
600 to reproduce that algorithm.
  - 601 (b) If the contribution is primarily a new model architecture, the paper should describe  
602 the architecture clearly and fully.

- 603 (c) If the contribution is a new model (e.g., a large language model), then there should  
604 either be a way to access this model for reproducing the results or a way to reproduce  
605 the model (e.g., with an open-source dataset or instructions for how to construct  
606 the dataset).
- 607 (d) We recognize that reproducibility may be tricky in some cases, in which case  
608 authors are welcome to describe the particular way they provide for reproducibility.  
609 In the case of closed-source models, it may be that access to the model is limited in  
610 some way (e.g., to registered users), but it should be possible for other researchers  
611 to have some path to reproducing or verifying the results.

## 612 5. Open access to data and code

613 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
614 tions to faithfully reproduce the main experimental results, as described in supplemental  
615 material?

616 Answer: [Yes]

617 Justification: The code is attached as supplementary material and will be released publicly.

618 Guidelines:

- 619 • The answer NA means that paper does not include experiments requiring code.
- 620 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
621 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 622 • While we encourage the release of code and data, we understand that this might not be  
623 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
624 including code, unless this is central to the contribution (e.g., for a new open-source  
625 benchmark).
- 626 • The instructions should contain the exact command and environment needed to run to  
627 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
628 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 629 • The authors should provide instructions on data access and preparation, including how  
630 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 631 • The authors should provide scripts to reproduce all experimental results for the new  
632 proposed method and baselines. If only a subset of experiments are reproducible, they  
633 should state which ones are omitted from the script and why.
- 634 • At submission time, to preserve anonymity, the authors should release anonymized  
635 versions (if applicable).
- 636 • Providing as much information as possible in supplemental material (appended to the  
637 paper) is recommended, but including URLs to data and code is permitted.

## 638 6. Experimental Setting/Details

639 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
640 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
641 results?

642 Answer: [Yes]

643 Justification: We share our hyperparameters and code required to reproduce our results in  
644 the appendix and as supplemental material.

645 Guidelines:

- 646 • The answer NA means that the paper does not include experiments.
- 647 • The experimental setting should be presented in the core of the paper to a level of detail  
648 that is necessary to appreciate the results and make sense of them.
- 649 • The full details can be provided either with the code, in appendix, or as supplemental  
650 material.

## 651 7. Experiment Statistical Significance

652 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
653 information about the statistical significance of the experiments?

654 Answer: [Yes]



655 Justification: We report mean and standard deviation for the average episodic reward over the  
656 last 100 training episodes across five different agents trained random seeds per environment  
657 for all our experiments.

658 Guidelines:

- 659 • The answer NA means that the paper does not include experiments.
- 660 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
661 dence intervals, or statistical significance tests, at least for the experiments that support  
662 the main claims of the paper.
- 663 • The factors of variability that the error bars are capturing should be clearly stated (for  
664 example, train/test split, initialization, random drawing of some parameter, or overall  
665 run with given experimental conditions).
- 666 • The method for calculating the error bars should be explained (closed form formula,  
667 call to a library function, bootstrap, etc.)
- 668 • The assumptions made should be given (e.g., Normally distributed errors).
- 669 • It should be clear whether the error bar is the standard deviation or the standard error  
670 of the mean.
- 671 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
672 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
673 of Normality of errors is not verified.
- 674 • For asymmetric distributions, the authors should be careful not to show in tables or  
675 figures symmetric error bars that would yield results that are out of range (e.g. negative  
676 error rates).
- 677 • If error bars are reported in tables or plots, The authors should explain in the text how  
678 they were calculated and reference the corresponding figures or tables in the text.

## 679 8. Experiments Compute Resources

680 Question: For each experiment, does the paper provide sufficient information on the com-  
681 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
682 the experiments?

683 Answer: [Yes]

684 Justification: We provide information on computer resources in the appendix.

685 Guidelines:

- 686 • The answer NA means that the paper does not include experiments.
- 687 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
688 or cloud provider, including relevant memory and storage.
- 689 • The paper should provide the amount of compute required for each of the individual  
690 experimental runs as well as estimate the total compute.
- 691 • The paper should disclose whether the full research project required more compute  
692 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
693 didn't make it into the paper).

## 694 9. Code Of Ethics

695 Question: Does the research conducted in the paper conform, in every respect, with the  
696 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

697 Answer: [Yes]

698 Justification: Yes, our work conforms with the NeurIPS Code of Ethics.

699 Guidelines:

- 700 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 701 • If the authors answer No, they should explain the special circumstances that require a  
702 deviation from the Code of Ethics.
- 703 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
704 eration due to laws or regulations in their jurisdiction).

## 705 10. Broader Impacts

706 Question: Does the paper discuss both potential positive societal impacts and negative  
707 societal impacts of the work performed?

708 Answer: [No]

709 Justification: The primary focus of this paper is on the algorithmic development and  
710 performance evaluation of Gradient Boosting Trees in Reinforcement Learning (GBRL).  
711 While the paper does not explicitly discuss societal impacts, GBRL has the potential  
712 to positively influence various domains, such as inventory management, traffic signal  
713 optimization, network optimization, resource allocation, and robotics. These domains have  
714 direct implications for the day-to-day lives of many people. The enhanced performance and  
715 the capability of GBRL to be deployed on edge devices could bring AI to new applications,  
716 potentially leading to significant societal benefits. However, as this work is foundational  
717 research, it does not address specific societal impacts or applications directly.

718 Guidelines:

- 719 • The answer NA means that there is no societal impact of the work performed.
- 720 • If the authors answer NA or No, they should explain why their work has no societal  
721 impact or why the paper does not address societal impact.
- 722 • Examples of negative societal impacts include potential malicious or unintended uses  
723 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
724 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
725 groups), privacy considerations, and security considerations.
- 726 • The conference expects that many papers will be foundational research and not tied  
727 to particular applications, let alone deployments. However, if there is a direct path to  
728 any negative applications, the authors should point it out. For example, it is legitimate  
729 to point out that an improvement in the quality of generative models could be used to  
730 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
731 that a generic algorithm for optimizing neural networks could enable people to train  
732 models that generate Deepfakes faster.
- 733 • The authors should consider possible harms that could arise when the technology is  
734 being used as intended and functioning correctly, harms that could arise when the  
735 technology is being used as intended but gives incorrect results, and harms following  
736 from (intentional or unintentional) misuse of the technology.
- 737 • If there are negative societal impacts, the authors could also discuss possible mitigation  
738 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
739 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
740 feedback over time, improving the efficiency and accessibility of ML).

## 741 11. Safeguards

742 Question: Does the paper describe safeguards that have been put in place for responsible  
743 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
744 image generators, or scraped datasets)?

745 Answer: [NA]

746 Justification: The paper poses no such risks.

747 Guidelines:

- 748 • The answer NA means that the paper poses no such risks.
- 749 • Released models that have a high risk for misuse or dual-use should be released with  
750 necessary safeguards to allow for controlled use of the model, for example by requiring  
751 that users adhere to usage guidelines or restrictions to access the model or implementing  
752 safety filters.
- 753 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
754 should describe how they avoided releasing unsafe images.
- 755 • We recognize that providing effective safeguards is challenging, and many papers do  
756 not require this, but we encourage authors to take this into account and make a best  
757 faith effort.

## 758 12. Licenses for existing assets

759 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
760 the paper, properly credited and are the license and terms of use explicitly mentioned and  
761 properly respected?

762 Answer: [Yes]

763 Justification: We are the original owners of the models and algorithm code. We credit the  
764 repository and authors of all datasets/models we based our implementations on.

765 Guidelines:

- 766 • The answer NA means that the paper does not use existing assets.
- 767 • The authors should cite the original paper that produced the code package or dataset.
- 768 • The authors should state which version of the asset is used and, if possible, include a  
769 URL.
- 770 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 771 • For scraped data from a particular source (e.g., website), the copyright and terms of  
772 service of that source should be provided.
- 773 • If assets are released, the license, copyright information, and terms of use in the  
774 package should be provided. For popular datasets, `paperswithcode.com/datasets`  
775 has curated licenses for some datasets. Their licensing guide can help determine the  
776 license of a dataset.
- 777 • For existing datasets that are re-packaged, both the original license and the license of  
778 the derived asset (if it has changed) should be provided.
- 779 • If this information is not available online, the authors are encouraged to reach out to  
780 the asset's creators.

### 781 13. New Assets

782 Question: Are new assets introduced in the paper well documented and is the documentation  
783 provided alongside the assets?

784 Answer: [Yes]

785 Justification: We will release the GBRL code publicly. The code is documented, with  
786 instructions provided on how to install, use, and incorporate within RL libraries. Additionally,  
787 details regarding the training process, license, limitations, and other relevant information are  
788 included in the documentation. The documentation will be made available alongside the  
789 code upon release.

790 Guidelines:

- 791 • The answer NA means that the paper does not release new assets.
- 792 • Researchers should communicate the details of the dataset/code/model as part of their  
793 submissions via structured templates. This includes details about training, license,  
794 limitations, etc.
- 795 • The paper should discuss whether and how consent was obtained from people whose  
796 asset is used.
- 797 • At submission time, remember to anonymize your assets (if applicable). You can either  
798 create an anonymized URL or include an anonymized zip file.

### 799 14. Crowdsourcing and Research with Human Subjects

800 Question: For crowdsourcing experiments and research with human subjects, does the paper  
801 include the full text of instructions given to participants and screenshots, if applicable, as  
802 well as details about compensation (if any)?

803 Answer: [NA]

804 Justification: The paper does not involve crowdsourcing nor research with human subjects.

805 Guidelines:

- 806 • The answer NA means that the paper does not involve crowdsourcing nor research with  
807 human subjects.
- 808 • Including this information in the supplemental material is fine, but if the main contribu-  
809 tion of the paper involves human subjects, then as much detail as possible should be  
810 included in the main paper.

811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832

- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: the paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.