

---

# Policy Gradient with Tree Expansion

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Policy gradient methods are notorious for having a large variance and high sample  
2 complexity. To mitigate this, we introduce SoftTreeMax—a generalization of  
3 softmax that employs planning. In SoftTreeMax, we extend the traditional logits  
4 with the multi-step discounted cumulative reward, topped with the logits of future  
5 states. We analyze SoftTreeMax and explain how tree expansion helps to reduce  
6 its gradient variance. We prove that the variance decays exponentially with the  
7 planning horizon as a function of the chosen tree-expansion policy. Specifically,  
8 we show that the closer the induced transitions are to being state-independent,  
9 the stronger the decay. With approximate forward models, we prove that the  
10 resulting gradient bias diminishes with the approximation error while retaining  
11 the same variance reduction. Ours is the first result to bound the gradient bias for  
12 an approximate model. In a practical implementation of SoftTreeMax, we utilize  
13 a parallel GPU-based simulator for fast and efficient tree expansion. Using this  
14 implementation in Atari, we show that SoftTreeMax reduces the gradient variance  
15 by three orders of magnitude. This leads to better sample complexity and improved  
16 performance compared to distributed PPO.

## 17 1 Introduction

18 Policy Gradient (PG) methods [Sutton et al., 1999] for Reinforcement Learning (RL) are often the  
19 first choice for environments that allow numerous interactions at a fast pace [Schulman et al., 2017].  
20 Their success is attributed to several factors: they are easy to distribute to multiple workers, require  
21 no assumptions on the underlying value function, and have both on-policy and off-policy variants.

22 Despite these positive features, PG algorithms are also notoriously unstable due to the high variance  
23 of the gradients computed over entire trajectories [Liu et al., 2020, Xu et al., 2020]. As a result, PG  
24 algorithms tend to be highly inefficient in terms of sample complexity. Several solutions have been  
25 proposed to mitigate the high variance issue, including baseline subtraction [Greensmith et al., 2004,  
26 Thomas and Brunskill, 2017, Wu et al., 2018], anchor-point averaging [Papini et al., 2018], and other  
27 variance reduction techniques [Zhang et al., 2021, Shen et al., 2019, Pham et al., 2020].

28 A second family of algorithms that achieved state-of-the-art results in several domains is based on  
29 planning. Planning is exercised primarily in the context of value-based RL and is usually implemented  
30 using a Tree Search (TS) [Silver et al., 2016, Schrittwieser et al., 2020]. In this work, we combine  
31 PG with TS by introducing a parameterized differentiable policy that incorporates tree expansion.  
32 Namely, our SoftTreeMax policy replaces the standard policy logits of a state and action, with the  
33 expected value of trajectories that originate from these state and action. We consider two variants of  
34 SoftTreeMax, one for cumulative reward and one for exponentiated reward.

35 Combining TS and PG should be done with care given the biggest downside of PG—its high gradient  
36 variance. This raises questions that were ignored until this work: (i) How to design a PG method based  
37 on tree-expansion that is stable and performs well in practice? and (ii) How does the tree-expansion

38 policy affect the PG variance? Here, we analyze SoftTreeMax, and provide a practical methodology  
 39 to choose the expansion policy to minimize the resulting variance. Our main result shows that a  
 40 desirable expansion policy is one, under which the induced transition probabilities are similar for  
 41 each starting state. More generally, we show that the gradient variance of SoftTreeMax decays at  
 42 a rate of  $|\lambda_2|^d$ , where  $d$  is the depth of the tree and  $\lambda_2$  is the second eigenvalue of the transition  
 43 matrix induced by the tree expansion policy. This work is the first to prove such a relation between  
 44 PG variance and tree expansion policy. In addition, we prove that the with an approximate forward  
 45 model, the bias of the gradient is bounded proportionally to the approximation error of the model.

46 To verify our results, we implemented a practical version of SoftTreeMax that exhaustively searches  
 47 the entire tree and applies a neural network on its leaves. We test our algorithm on a parallelized  
 48 Atari GPU simulator [Dalton et al., 2020]. To enable a tractable deep search, up to depth eight, we  
 49 also introduce a pruning technique that limits the width of the tree. We do so by sampling only the  
 50 most promising nodes at each level. We integrate our SoftTreeMax GPU implementation into the  
 51 popular PPO [Schulman et al., 2017] and compare it to the flat distributed variant of PPO. This allows  
 52 us to demonstrate the potential benefit of utilizing learned models while isolating the fundamental  
 53 properties of TS without added noise. In all tested Atari games, our results outperform the baseline  
 54 and obtain up to 5x more reward. We further show in Section 6 that the associated gradient variance  
 55 is smaller by three orders of magnitude in all games, demonstrating the relation between low gradient  
 56 variance and high reward.

57 We summarize our key contributions. (i) We show how to combine two families of SoTA approaches:  
 58 PG and TS by **introducing SoftTreeMax**: a novel parametric policy that generalizes softmax to  
 59 planning. Specifically, we propose two variants based on cumulative and exponentiated rewards. (ii)  
 60 **We prove that the gradient variance of SoftTreeMax in its two variants decays exponentially**  
 61 with its tree depth. Our analysis sheds new light on the choice of tree expansion policy. It raises  
 62 the question of optimality in terms of variance versus the traditional regret; e.g., in UCT [Kocsis  
 63 and Szepesvári, 2006]. (iii) We prove that with an approximate forward model, the **gradient bias is**  
 64 **proportional to the approximation error**, while retaining the variance decay. This quantifies the  
 65 accuracy required from a learned forward model. (iv) We **implement a differentiable deep version**  
 66 **of SoftTreeMax** that employs a parallelized GPU tree expansion. We demonstrate how its gradient  
 67 variance is reduced by three orders of magnitude over PPO while obtaining up to 5x reward.

## 68 2 Preliminaries

69 Let  $\Delta_U$  denote simplex over the set  $U$ . Throughout, we consider a discounted Markov Decision  
 70 Process (MDP)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma, \nu)$ , where  $\mathcal{S}$  is a finite state space of size  $S$ ,  $\mathcal{A}$  is a finite action  
 71 space of size  $A$ ,  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is the reward function,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$  is the transition  
 72 function,  $\gamma \in (0, 1)$  is the discount factor, and  $\nu \in \mathbb{R}^S$  is the initial state distribution. We denote  
 73 the transition matrix starting from state  $s$  by  $P_s \in [0, 1]^{A \times S}$ , i.e.,  $[P_s]_{a,s'} = P(s'|a, s)$ . Similarly,  
 74 let  $R_s = r(s, \cdot) \in \mathbb{R}^A$  denote the corresponding reward vector. Separately, let  $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$  be a  
 75 stationary policy. Let  $P^\pi$  and  $R_\pi$  be the induced transition matrix and reward function, respectively,  
 76 i.e.,  $P^\pi(s'|s) = \sum_a \pi(a|s) \Pr(s'|s, a)$  and  $R_\pi(s) = \sum_a \pi(a|s) r(s, a)$ . Denote the stationary  
 77 distribution of  $P^\pi$  by  $\mu_\pi \in \mathbb{R}^S$  s.t.  $\mu_\pi^\top P^\pi = P^\pi$ , and the discounted state visitation frequency  
 78 by  $d_\pi$  so that  $d_\pi^\top = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \nu^\top (P^\pi)^t$ . Also, let  $V^\pi \in \mathbb{R}^S$  be the value function of  $\pi$   
 79 defined by  $V^\pi(s) = \mathbb{E}^\pi [\sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s]$ , and let  $Q^\pi \in \mathbb{R}^{S \times A}$  be the Q-function  
 80 such that  $Q^\pi(s, a) = \mathbb{E}^\pi [r(s, a) + \gamma V^\pi(s')]$ . Our goal is to find an optimal policy  $\pi^*$  such that  
 81  $V^*(s) \equiv V^{\pi^*}(s) = \max_\pi V^\pi(s), \forall s \in \mathcal{S}$ .

82 For the analysis in Section 4, we introduce the following notation. Denote by  $\Theta \in \mathbb{R}^S$  the vector  
 83 representation of  $\theta(s) \forall s \in \mathcal{S}$ . For a vector  $u$ , denote by  $\exp(u)$  the coordinate-wise exponent of  
 84  $u$  and by  $D(u)$  the diagonal square matrix with  $u$  in its diagonal. For a matrix  $A$ , denote its  $i$ -th  
 85 eigenvalue by  $\lambda_i(A)$ . Denote the  $k$ -dimensional identity matrix and all-ones vector by  $I_k$  and  $\mathbf{1}_k$ ,  
 86 respectively. Also, denote the trace operator by  $\text{Tr}$ . Finally, we treat all vectors as column vectors.

87 **2.1 Policy Gradient**

88 PG schemes seek to maximize the cumulative reward as a function of the policy  $\pi_\theta(a|s)$  by performing  
 89 gradient steps on  $\theta$ . The celebrated Policy Gradient Theorem [Sutton et al., 1999] states that

$$\frac{\partial}{\partial \theta} \nu^\top V^{\pi_\theta} = \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a)].$$

90 The variance of the gradient is thus

$$\text{Var}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} (\nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a)). \quad (1)$$

91 In the notation above, we denote the variance of a vector random variable  $X$  by

$$\text{Var}_x(X) = \text{Tr} \left[ \mathbb{E}_x \left[ (X - \mathbb{E}_x X)^\top (X - \mathbb{E}_x X) \right] \right],$$

92 similarly as in [Greensmith et al., 2004]. From now on, we drop the subscript from Var in (1)  
 93 for brevity. When the action space is discrete, a commonly used parameterized policy is softmax:  
 94  $\pi_\theta(a|s) \propto \exp(\theta(s, a))$ , where  $\theta : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is a state-action parameterization.

95 **3 SoftTreeMax: Exponent of trajectories**

96 We introduce a new family of policies called SoftTreeMax, which are a model-based generalization  
 97 of the popular softmax. We propose two variants: Cumulative (C-SoftTreeMax) and Exponenti-  
 98 ated (E-SoftTreeMax). In both variants, we replace the generic softmax logits  $\theta(s, a)$  with the  
 99 score of a trajectory of horizon  $d$  starting from  $(s, a)$ , generated by applying a behavior policy  
 100  $\pi_b$ . In C-SoftTreeMax, we exponentiate the expectation of the logits. In E-SoftTreeMax, we first  
 101 exponentiate the logits and then only compute their expectation.

102 **Logits.** We define the SoftTreeMax logit  $\ell_{s,a}(d; \theta)$  to be the random variable depicting the score of a  
 103 trajectory of horizon  $d$  starting from  $(s, a)$  and following the policy  $\pi_b$ :

$$\ell_{s,a}(d; \theta) = \gamma^{-d} \left[ \sum_{t=0}^{d-1} \gamma^t r_t + \gamma^d \theta(s_d) \right]. \quad (2)$$

104 In the above expression, note that  $s_0 = s$ ,  $a_0 = a$ ,  $a_t \sim \pi_b(\cdot|s_t) \forall t \geq 1$ , and  $r_t \equiv r(s_t, a_t)$ .  
 105 For brevity of the analysis, we let the parametric score  $\theta$  in (2) be state-based, similarly to a value  
 106 function. Instead, one could use a state-action input analogous to a Q-function. Thus, SoftTreeMax  
 107 can be integrated into the two types of implementation of RL algorithms in standard packages. Lastly,  
 108 the preceding  $\gamma^{-d}$  scales the  $\theta$  parametrization to correspond to its softmax counterpart.

109 **C-SoftTreeMax.** Given an inverse temperature parameter  $\beta$ , we let C-SoftTreeMax be

$$\pi_{d,\theta}^C(a|s) \propto \exp[\beta \mathbb{E}^{\pi_b} \ell_{s,a}(d; \theta)]. \quad (3)$$

110 C-SoftTreeMax gives higher weight to actions that result in higher expected returns. While standard  
 111 softmax relies entirely on parametrization  $\theta$ , C-SoftTreeMax also interpolates a Monte-Carlo portion  
 112 of the reward.

113 **E-SoftTreeMax.** The second operator we propose is E-SoftTreeMax:

$$\pi_{d,\theta}^E(a|s) \propto \mathbb{E}^{\pi_b} \exp[(\beta \ell_{s,a}(d; \theta))]; \quad (4)$$

114 here, the expectation is taken outside the exponent. This objective corresponds to the exponentiated  
 115 reward objective which is often used for risk-sensitive RL [Howard and Matheson, 1972, Fei et al.,  
 116 2021, Noorani and Baras, 2021]. The common risk-sensitive objective is of the form  $\log \mathbb{E}[\exp(\delta R)]$ ,  
 117 where  $\delta$  is the risk parameter and  $R$  is the cumulative reward. Similarly to that literature, the exponent  
 118 in (4) emphasizes the most promising trajectories.

119 **SoftTreeMax properties.** SoftTreeMax is a natural model-based generalization of softmax. For  
 120  $d = 0$ , both variants above coincide since (2) becomes deterministic. In that case, for a state-action  
 121 parametrization, they reduce to standard softmax. When  $\beta \rightarrow 0$ , both variants again coincide and  
 122 sample actions uniformly (exploration). When  $\beta \rightarrow \infty$ , the policies become deterministic and

123 greedily optimize for the best trajectory (exploitation). For C-SoftTreeMax, the best trajectory is  
 124 defined in expectation, while for E-SoftTreeMax it is defined in terms of the best sample path.

125 **SoftTreeMax convergence.** Under regularity conditions, for any parametric policy, PG converges  
 126 to local optima [Bhatnagar et al., 2009], and thus also SoftTreeMax. For softmax PG, asymptotic  
 127 [Agarwal et al., 2021] and rate results [Mei et al., 2020b] were recently obtained, by showing that  
 128 the gradient is strictly positive everywhere [Mei et al., 2020b, Lemmas 8-9]. We conjecture that  
 129 SoftTreeMax satisfies the same property, being a generalization of softmax, but formally proving it is  
 130 subject to future work.

131 **SoftTreeMax gradient.** The two variants of SoftTreeMax involve an expectation taken over  $S^d$   
 132 many trajectories from the root state  $s$  and weighted according to their probability. Thus, during  
 133 the PG training process, the gradient  $\nabla_{\theta} \log \pi_{\theta}$  is calculated using a weighted sum of gradients over  
 134 all reachable states starting from  $s$ . Our method exploits the exponential number of trajectories to  
 135 reduce the variance while improving performance. Indeed, in the next section we prove that the  
 136 gradient variance of SoftTreeMax decays exponentially fast as a function of the behavior policy  $\pi_b$   
 137 and trajectory length  $d$ . In the experiments in Section 6, we also show how the practical version  
 138 of SoftTreeMax achieves a significant reduction in the noise of the PG process and leads to faster  
 139 convergence and higher reward.

## 140 4 Theoretical Analysis

141 In this section, we first bound the variance of PG when using the SoftTreeMax policy. Later, we  
 142 discuss how the gradient bias resulting due to approximate forward models diminishes as a function  
 143 of the approximation error, while retaining the same variance decay.

144 We show that the variance decreases exponentially with the tree depth, and the rate is determined  
 145 by the second eigenvalue of the transition kernel induced by  $\pi_b$ . Specifically, we bound the same  
 146 expression for variance as appears in [Greensmith et al., 2004, Sec. 3.5] and [Wu et al., 2018, Sec. A,  
 147 Eq. (21)]. Other types of analysis could instead have focused on the estimation aspect in the context  
 148 of sampling [Zhang et al., 2021, Shen et al., 2019, Pham et al., 2020]. Indeed, in our implementation  
 149 in Section 5, we manage to avoid sampling and directly compute the expectations in Eqs. (3) and  
 150 (4). As we show later, we do so by leveraging efficient parallel simulation on the GPU in feasible  
 151 run-time. In our application, due to the nature of the finite action space and quasi-deterministic Atari  
 152 dynamics [Bellemare et al., 2013], our expectation estimator is noiseless. We encourage future work  
 153 to account for the finite-sample variance component. We defer all the proofs to Appendix A.

154 We begin with a general variance bound that holds for any parametric policy.

155 **Lemma 4.1** (Bound on the policy gradient variance). *Let  $\nabla_{\theta} \log \pi_{\theta}(\cdot|s) \in \mathbb{R}^{A \times \dim(\theta)}$  be a matrix  
 156 whose  $a$ -th row is  $\nabla_{\theta} \log \pi_{\theta}(a|s)^{\top}$ . For any parametric policy  $\pi_{\theta}$  and function  $Q^{\pi_{\theta}} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ ,*

$$\text{Var}(\nabla_{\theta} \log \pi_{\theta}(a|s)Q^{\pi_{\theta}}(s, a)) \leq \max_{s,a} [Q^{\pi_{\theta}}(s, a)]^2 \max_s \|\nabla_{\theta} \log \pi_{\theta}(\cdot|s)\|_F^2.$$

157 Hence, to bound (1), it is sufficient to bound the Frobenius norm  $\|\nabla_{\theta} \log \pi_{\theta}(\cdot|s)\|_F$  for any  $s$ .

158 Note that SoftTreeMax does not reduce the gradient uniformly, which would have been equivalent  
 159 to a trivial change in the learning rate. While the gradient norm shrinks, the gradient itself scales  
 160 differently along the different coordinates. This scaling occurs along different eigenvectors, as a  
 161 function of problem parameters  $(P, \theta)$  and our choice of behavior policy  $(\pi_b)$ , as can be seen in  
 162 the proof of the upcoming Theorem 4.4. This allows SoftTreeMax to learn a good “shrinkage” that,  
 163 while reducing the overall gradient, still updates the policy quickly enough. This reduction in norm  
 164 and variance resembles the idea of gradient clipping Zhang et al. [2019], where the gradient is scaled  
 165 to reduce its variance, thus increasing stability and improving overall performance.

166 A common assumption in the RL literature [Szepesvári, 2010] that we adopt for the remainder of  
 167 the section is that the transition matrix  $P^{\pi_b}$ , induced by the behavior policy  $\pi_b$ , is irreducible and  
 168 aperiodic. Consequently, its second highest eigenvalue satisfies  $|\lambda_2(P^{\pi_b})| < 1$ .

169 From now on, we divide the variance results for the two variants of SoftTreeMax into two subsec-  
 170 tions. For C-SoftTreeMax, the analysis is simpler and we provide an exact bound. The case of  
 171 E-SoftTreeMax is more involved and we provide for it a more general result. In both cases, we show  
 172 that the variance decays exponentially with the planning horizon.

173 **4.1 Variance of C-SoftTreeMax**

174 We express C-SoftTreeMax in vector form as follows.

175 **Lemma 4.2** (Vector form of C-SoftTreeMax). *For  $d \geq 1$ , (3) is given by*

$$\pi_{d,\theta}^C(\cdot|s) = \frac{\exp\left[\beta\left(C_{s,d} + P_s(P^{\pi_b})^{d-1}\Theta\right)\right]}{\mathbf{I}_A^\top \exp\left[\beta\left(C_{s,d} + P_s(P^{\pi_b})^{d-1}\Theta\right)\right]}, \quad (5)$$

176 *where*

$$C_{s,d} = \gamma^{-d}R_s + P_s \left[ \sum_{h=1}^{d-1} \gamma^{h-d} (P^{\pi_b})^{h-1} \right] R_{\pi_b}.$$

177 The vector  $C_{s,d} \in \mathbb{R}^A$  represents the cumulative discounted reward in expectation along the trajectory  
 178 of horizon  $d$ . This trajectory starts at state  $s$ , involves an initial reward dictated by  $R_s$  and an  
 179 initial transition as per  $P_s$ . Thereafter, it involves rewards and transitions specified by  $R_{\pi_b}$  and  $P^{\pi_b}$ ,  
 180 respectively. Once the trajectory reaches depth  $d$ , the score function  $\theta(s_d)$  is applied.

181 **Lemma 4.3** (Gradient of C-SoftTreeMax). *The C-SoftTreeMax gradient is given by*

$$\nabla_\theta \log \pi_{d,\theta}^C = \beta \left[ \mathbf{I}_A - \mathbf{I}_A (\pi_{d,\theta}^C)^\top \right] P_s (P^{\pi_b})^{d-1},$$

182 *in  $\mathbb{R}^{A \times S}$ , where for brevity, we drop the  $s$  index in the policy above, i.e.,  $\pi_{d,\theta}^C \equiv \pi_{d,\theta}^C(\cdot|s)$ .*

183 We are now ready to present our first main result:

184 **Theorem 4.4** (Variance decay of C-SoftTreeMax). *For every  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , the C-SoftTreeMax*  
 185 *policy gradient variance is bounded by*

$$\text{Var} \left( \nabla_\theta \log \pi_{d,\theta}^C(a|s) Q(s, a) \right) \leq 2 \frac{A^2 S^2 \beta^2}{(1-\gamma)^2} |\lambda_2(P^{\pi_b})|^{2(d-1)}.$$

186 We provide the full proof in Appendix A.4, and briefly outline its essence here.

187 *Proof outline.* Lemma 4.1 allows us to bound the variance using a direct bound on the gradient  
 188 norm. The gradient is given in Lemma 4.3 as a product of three matrices, which we now study from  
 189 right to left. The matrix  $P^{\pi_b}$  is a row-stochastic matrix. Because the associated Markov chain is  
 190 irreducible and aperiodic, it has a unique stationary distribution. This implies that  $P^{\pi_b}$  has one and  
 191 only one eigenvalue equal to 1; all others have magnitude strictly less than 1. Let us suppose that  
 192 all these other eigenvalues have multiplicity 1 (the general case with repeated eigenvalues can be  
 193 handled via Jordan decompositions as in [Pelletier, 1998, Lemmal]). Then,  $P^{\pi_b}$  has the spectral  
 194 decomposition  $P^{\pi_b} = \mathbf{1}_S \mu_{\pi_b}^\top + \sum_{i=2}^S \lambda_i v_i u_i^\top$ , where  $\lambda_i$  is the  $i$ -th eigenvalue of  $P^{\pi_b}$  (ordered in  
 195 descending order according to their magnitude) and  $u_i$  and  $v_i$  are the corresponding left and right  
 196 eigenvectors, respectively, and therefore  $(P^{\pi_b})^{d-1} = \mathbf{1}_S \mu_{\pi_b}^\top + \sum_{i=2}^S \lambda_i^{d-1} v_i u_i^\top$ .

197 The second matrix in the gradient relation in Lemma 4.3,  $P_s$ , is a rectangular transition ma-  
 198 trix that translates the vector of all ones from dimension  $S$  to  $A$  :  $P_s \mathbf{1}_S = \mathbf{1}_A$ . Lastly, the  
 199 first matrix  $\left[ \mathbf{I}_A - \mathbf{1}_A (\pi_{d,\theta}^C)^\top \right]$  is a projection whose null-space includes the vector  $\mathbf{1}_A$ , i.e.,  
 200  $\left[ \mathbf{I}_A - \mathbf{1}_A (\pi_{d,\theta}^C)^\top \right] \mathbf{1}_A = 0$ . Combining the three properties above when multiplying the three matri-  
 201 ces of the gradient, it is easy to see that the first term in the expression for  $(P^{\pi_b})^{d-1}$  gets canceled,  
 202 and we are left with bounded summands scaled by  $\lambda_i (P^{\pi_b})^{d-1}$ . Recalling that  $|\lambda_i(P^{\pi_b})| < 1$  and  
 203 that  $|\lambda_2| \geq |\lambda_3| \geq \dots$  for  $i = 2, \dots, S$ , we obtain the desired result.  $\square$

204 Theorem 4.4 guarantees that the variance of the gradient decays exponentially with  $d$ . It also provides  
 205 a novel insight for choosing the behavior policy  $\pi_b$  as the policy that minimizes the absolute second  
 206 eigenvalue of the  $P^{\pi_b}$ . Indeed, the second eigenvalue of a Markov chain relates to its connectivity  
 207 and its rate of convergence to the stationary distribution [Levin and Peres, 2017].

208 **Optimal variance decay.** For the strongest reduction in variance, the behavior policy  $\pi_b$  should be  
 209 chosen to achieve an induced Markov chain whose transitions are state-independent. In that case,  $P^{\pi_b}$

210 is a rank one matrix of the form  $\mathbf{1}_S \mu_{\pi_b}^\top$ , and  $\lambda_2(P^{\pi_b}) = 0$ . Then,  $\text{Var}(\nabla_\theta \log \pi_\theta(a|s)Q(s, a)) = 0$ .  
 211 Naturally, this can only be done for pathological MDPs; see Appendix C.1 for a more detailed  
 212 discussion. Nevertheless, as we show in Section 5, we choose our tree expansion policy to reduce the  
 213 variance as best as possible.

214 **Worst-case variance decay.** In contrast, and somewhat surprisingly, when  $\pi_b$  is chosen so that the  
 215 dynamics is deterministic, there is no guarantee that it will decay exponentially fast. For example, if  
 216  $\hat{P}^{\pi_b}$  is a permutation matrix, then  $\lambda_2(P^{\pi_b}) = 1$ , and advancing the tree amounts to only updating the  
 217 gradient of one state for every action, as in the basic softmax.

## 218 4.2 Variance of E-SoftTreeMax

219 The proof of the variance bound for E-SoftTreeMax is similar to that of C-SoftTreeMax, but more  
 220 involved. It also requires the assumption that the reward depends only on the state, i.e.  $r(s, a) \equiv r(s)$ .  
 221 This is indeed the case in most standard RL environments such as Atari and Mujoco.

222 **Lemma 4.5** (Vector form of E-SoftTreeMax). *For  $d \geq 1$ , (4) is given by*

$$\pi_{d,\theta}^E(\cdot|s) = \frac{E_{s,d} \exp(\beta\Theta)}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)}, \quad (6)$$

223 where

$$E_{s,d} = P_s \prod_{h=1}^{d-1} (D(\exp(\beta\gamma^{h-d}R)) P^{\pi_b}).$$

224 The vector  $R$  above is the  $S$ -dimensional vector whose  $s$ -th coordinate is  $r(s)$ .

225 The matrix  $E_{s,d} \in \mathbb{R}^{A \times S}$  has a similar role to  $C_{s,d}$  from (5), but it represents the exponentiated  
 226 cumulative discounted reward. Accordingly, it is a product of  $d$  matrices as opposed to a sum. It  
 227 captures the expected reward sequence starting from  $s$  and then iteratively following  $P^{\pi_b}$ . After  $d$   
 228 steps, we apply the score function on the last state as in (6).

229 **Lemma 4.6** (Gradient of E-SoftTreeMax). *The E-SoftTreeMax gradient is given by*

$$\nabla_\theta \log \pi_{d,\theta}^E = \beta [I_A - \mathbf{I}_A(\pi_{d,\theta}^E)^\top] \times \frac{D(\pi_{d,\theta}^E)^{-1} E_{s,d} D(\exp(\beta\Theta))}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \in \mathbb{R}^{A \times S},$$

230 where for brevity, we drop the  $s$  index in the policy above, i.e.,  $\pi_{d,\theta}^E \equiv \pi_{d,\theta}^E(\cdot|s)$ .

231 This gradient structure is harder to handle than that of C-SoftTreeMax in Lemma 4.3, but here we  
 232 also can bound the decay of the variance nonetheless.

233 **Theorem 4.7** (Variance decay of E-SoftTreeMax). *There exists  $\alpha \in (0, 1)$  such that,*

$$\text{Var}(\nabla_\theta \log \pi_{d,\theta}^E(a|s)Q(s, a)) \in \mathcal{O}(\beta^2 \alpha^{2d}),$$

234 for every  $Q$ . Further, if  $P^{\pi_b}$  is reversible or if the reward is constant, then  $\alpha = |\lambda_2(P^{\pi_b})|$ .

235 **Theory versus Practice.** We demonstrate the above result in simulation. We draw a random finite  
 236 MDP, parameter vector  $\Theta \in \mathbb{R}_+^S$ , and behavior policy  $\pi_b$ . We then empirically compute the PG  
 237 variance of E-SoftTreeMax as given in (1) and compare it to  $|\lambda_2(P^{\pi_b})|^d$ . We repeat this experiment  
 238 three times for different  $P^{\pi_b}$ : (i) close to uniform, (ii) drawn randomly, and (iii) close to a permutation  
 239 matrix. As seen in Figure 1, the empirical variance and our bound match almost identically. This  
 240 also suggests that  $\alpha = |\lambda_2(P^{\pi_b})|$  in the general case and not only when  $P^{\pi_b}$  is reversible or when  
 241 the reward is constant.

## 242 4.3 Bias with an Approximate Forward Model

243 The definition of the two SoftTreeMax variants involves the knowledge of the underlying environment,  
 244 in particular the value of  $P$  and  $r$ . However, in practice, we often can only learn approximations of  
 245 the dynamics from interactions, e.g., using NNs [Ha and Schmidhuber, 2018, Schrittwieser et al.,  
 246 2020]. Let  $\hat{P}$  and  $\hat{r}$  denote the approximate kernel and reward functions, respectively. In this section,  
 247 we study the consequences of the approximation error on the C-SoftTreeMax gradient.

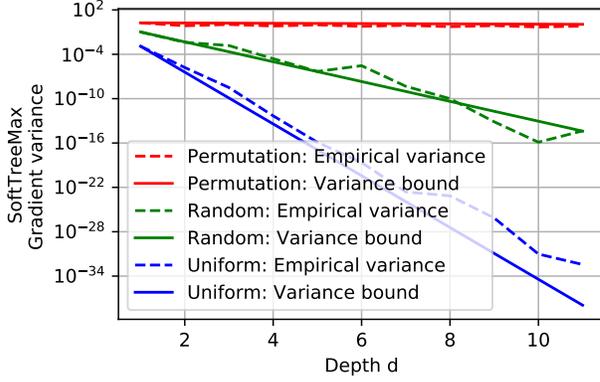


Figure 1: A comparison of the empirical PG variance and our bound for E-SoftTreeMax on randomly drawn MDPs. We present three cases for  $P^{\pi_b}$ : (i) close to uniform, (ii) drawn randomly, and (iii) close to a permutation matrix. This experiment verifies the optimal and worse-case rate decay cases. The variance bounds here are taken from Theorem 4.7 where we substitute  $\alpha = |\lambda_2(P^{\pi_b})|$ . To account for the constants, we match the values for the first point in  $d = 1$ .

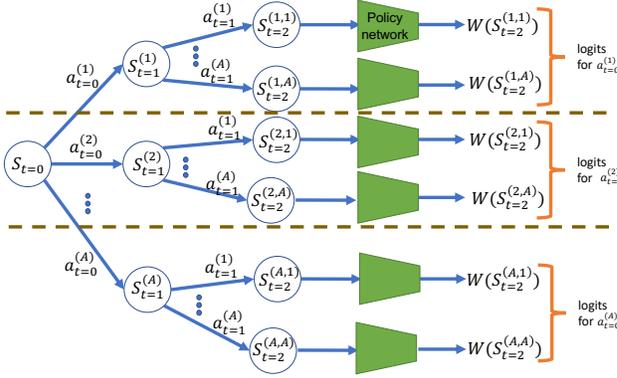


Figure 2: **SoftTreeMax policy.** Our exhaustive parallel tree expansion iterates on all actions at each state up to depth  $d (= 2$  here). The leaf state of every trajectory is used as input to the policy network. The output is then added to the trajectory’s cumulative reward as described in (2). I.e., instead of the standard softmax logits, we add the cumulative discounted reward to the policy network output. This policy is differentiable and can be easily integrated into any PG algorithm. In this work, we build on PPO and use its loss function to train the policy network.

248 Let  $\hat{\pi}_{d,\theta}^C$  be the C-SoftTreeMax policy defined given the approximate forward model introduced  
 249 above. That is, let  $\hat{\pi}_{d,\theta}^C$  be defined exactly as in (5), but using  $\hat{R}_s, \hat{P}_s, \hat{R}_{\pi_b}$  and  $\hat{P}^{\pi_b}$ , instead of their  
 250 unperturbed counterparts from Section 2. Then, the variance of the corresponding gradient again  
 251 decays exponentially with a decay rate of  $\lambda_2(\hat{P}^{\pi_b})$ . However, a gradient bias is introduced. In the  
 252 following, we bound this bias in terms of the approximation error and other problem parameters. The  
 253 proof is provided in Appendix A.9.

254 **Theorem 4.8.** *Let  $\epsilon$  be the maximal model mis-specification, i.e., let  $\max\{\|P - \hat{P}\|, \|r - \hat{r}\|\} = \epsilon$ .*  
 255 *Then the policy gradient bias due to  $\hat{\pi}_{d,\theta}^C$  satisfies*

$$\left\| \frac{\partial}{\partial \theta} \left( \nu^\top V^{\pi_{d,\theta}^c} \right) - \frac{\partial}{\partial \theta} \left( \nu^\top V^{\hat{\pi}_{d,\theta}^c} \right) \right\| = \mathcal{O} \left( \frac{1}{(1-\gamma)^2} S \beta^2 d \epsilon \right). \quad (7)$$

256 To the best of our knowledge, Theorem 4.8 is the first result that bounds the bias of the gradient  
 257 of a parametric policy due to an approximate model. It states that if the learned model is accurate  
 258 enough, we expect similar convergence properties for C-SoftTreeMax as we would have obtained  
 259 with the true dynamics. It also suggests that higher temperature (lower  $\beta$ ) reduces the bias. In this  
 260 case, the logits get less weight, with the extreme of  $\beta = 0$  corresponding to a uniform policy that has  
 261 no bias. Lastly, the error scales linearly with  $d$ : the policy suffers from cumulative error as it relies  
 262 on further-looking states in the approximate model.

## 263 5 SoftTreeMax: Deep Parallel Implementation

264 Following impressive successes of deep RL [Mnih et al., 2015, Silver et al., 2016], using deep NNs  
 265 in RL is standard practice. Depending on the RL algorithm, a loss function is defined and gradients  
 266 on the network weights can be calculated. In PG methods, the scoring function used in the softmax is  
 267 commonly replaced by a neural network  $W_\theta: \pi_\theta(a|s) \propto \exp(W_\theta(s, a))$ . Similarly, we implement  
 268 SoftTreeMax by replacing  $\theta(s)$  in (2) with a neural network  $W_\theta(s)$ . Although both variants of

269 SoftTreeMax from Section 3 involve computing an expectation, this can be hard in general. One  
270 approach to handle it is with sampling, though these introduce estimation variance into the process.  
271 We leave the question of sample-based theory and algorithmic implementations for future work.

272 Instead, in finite action space environments such as Atari, we compute the exact expectation in  
273 SoftTreeMax with an exhaustive TS of depth  $d$ . Despite the exponential computational cost of  
274 spanning the entire tree, recent advancements in parallel GPU-based simulation allow efficient  
275 expansion of all nodes at the same depth simultaneously [Dalal et al., 2021, Rosenberg et al., 2022].  
276 This is possible when a simulator is implemented on GPU [Dalton et al., 2020, Makoviychuk  
277 et al., 2021, Freeman et al., 2021], or when a forward model is learned [Kim et al., 2020, Ha and  
278 Schmidhuber, 2018]. To reduce the complexity to be linear in depth, we apply tree pruning to a  
279 limited width in all levels. We do so by sub-sampling only the most promising branches at each level.  
280 Limiting the width drastically improves runtime, and enables respecting GPU memory limits, with  
281 only a small sacrifice in performance.

282 To summarize, in the practical SoftTreeMax algorithm we perform an exhaustive tree expansion with  
283 pruning to obtain trajectories up to depth  $d$ . We expand the tree with equal weight to all actions, which  
284 corresponds to a uniform tree expansion policy  $\pi_b$ . We apply a neural network on the leaf states, and  
285 accumulate the result with the rewards along each trajectory to obtain the logits in (2). Finally, we  
286 aggregate the results using C-SoftTreeMax. We leave experiments E-SoftTreeMax for future work  
287 on risk-averse RL. During training, the gradient propagates to the NN weights of  $W_\theta$ . When the  
288 gradient  $\nabla_\theta \log \pi_{d,\theta}$  is calculated at each time step, it updates  $W_\theta$  for all leaf states, similarly to  
289 Siamese networks [Bertinetto et al., 2016]. An illustration of the policy is given in Figure 2.

## 290 6 Experiments

291 We conduct our experiments on multiple games from the Atari simulation suite [Bellemare et al.,  
292 2013]. As a baseline, we train a PPO [Schulman et al., 2017] agent with 256 GPU workers in parallel  
293 [Dalton et al., 2020]. For the tree expansion, we employ a GPU breadth-first as in [Dalal et al., 2021].  
294 We then train C-SoftTreeMax <sup>1</sup> for depths  $d = 1 \dots 8$ , with a single worker. For depths  $d \geq 3$ ,  
295 we limited the tree to a maximum width of 1024 nodes and pruned trajectories with low estimated  
296 weights. Since the distributed PPO baseline advances significantly faster in terms of environment  
297 steps, for a fair comparison, we ran all experiments for one week on the same machine. For more  
298 details see Appendix B.

299 In Figure 3, we plot the reward and variance of SoftTreeMax for each game, as a function of depth.  
300 The dashed lines are the results for PPO. Each value is taken after convergence, i.e., the average  
301 over the last 20% of the run. The numbers represent the average over five seeds per game. The plot  
302 conveys three intriguing conclusions. First, in all games, SoftTreeMax achieves significantly higher  
303 reward than PPO. Its gradient variance is also orders of magnitude lower than that of PPO. Second,  
304 the reward and variance are negatively correlated and mirror each other in almost all games. This  
305 phenomenon demonstrates the necessity of reducing the variance of PG for improving performance.  
306 Lastly, each game has a different sweet spot in terms of optimal tree depth. Recall that we limit the  
307 run-time in all experiments to one week. The deeper the tree, the slower each step and the run consists  
308 of less steps. This explains the non-monotone behavior as a function of depth. For a more thorough  
309 discussion on the sweet spot of different games, see Appendix B.3.

## 310 7 Related Work

311 **Softmax Operator.** The softmax policy became a canonical part of PG to the point where theoretical  
312 results of PG focus specifically on it [Zhang et al., 2021, Mei et al., 2020b, Li et al., 2021, Ding et al.,  
313 2022]. Even though we focus on a tree extension to the softmax policy, our methodology is general  
314 and can be easily applied to other discrete or continuous parameterized policies as in [Mei et al.,  
315 2020a, Miahhi et al., 2021, Silva et al., 2019]. **Tree Search.** One famous TS algorithm is Monte-Carlo  
316 TS (MCTS; [Browne et al., 2012]) used in AlphaGo [Silver et al., 2016] and MuZero [Schrittwieser  
317 et al., 2020]. Other algorithms such as Value Iteration, Policy Iteration and DQN were also shown to

<sup>1</sup>We also experimented with E-SoftTreeMax and the results were almost identical. This is due to the quasi-deterministic nature of Atari, which causes the trajectory logits (2) to have almost no variability. We encourage future work on E-SoftTreeMax using probabilistic environments that are risk-sensitive.

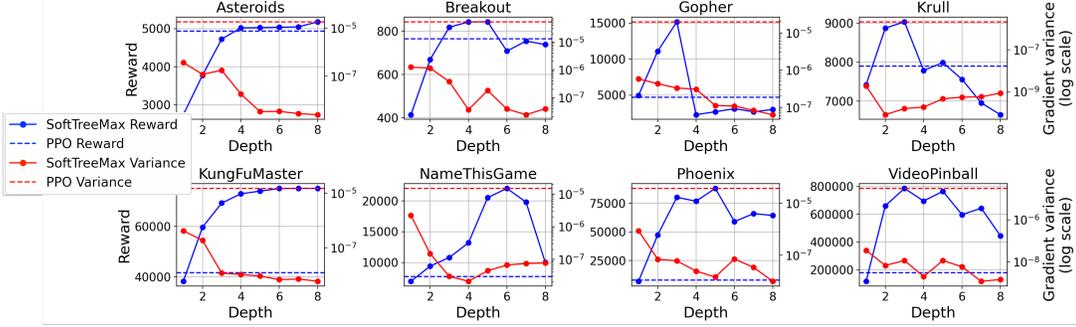


Figure 3: **Reward and Gradient variance: GPU SoftTreeMax (single worker) vs PPO (256 GPU workers).** The blue reward plots show the average of 50 evaluation episodes. The red variance plots show the average gradient variance of the corresponding training runs, averaged over five seeds. The dashed lines represent the same for PPO. Note that the variance y-axis is in log-scale.

318 give an improved performance with a tree search extensions [Efroni et al., 2019, Dalal et al., 2021].  
 319 **Parallel Environments.** In this work we used accurate parallel models that are becoming more  
 320 common with the increasing popularity of GPU-based simulation [Makoviychuk et al., 2021, Dalton  
 321 et al., 2020, Freeman et al., 2021]. Alternatively, in relation to Theorem 4.8, one can rely on recent  
 322 works that learn the underlying model [Ha and Schmidhuber, 2018, Schrittwieser et al., 2020] and  
 323 use an approximation of the true dynamics. **Risk Aversion.** Previous work considered exponential  
 324 utility functions for risk aversion [Chen et al., 2007, Garcia and Fernández, 2015, Fei et al., 2021].  
 325 This utility function is the same as E-SoftTreeMax formulation from (4), but we have it directly  
 326 in the policy instead of the objective. **Reward-free RL.** We showed that the gradient variance is  
 327 minimized when the transitions induced by the behavior policy  $\pi_b$  are uniform. This is expressed by  
 328 the second eigenvalue of the transition matrix  $P^{\pi_b}$ . This notion of uniform exploration is common to  
 329 the reward-free RL setup [Jin et al., 2020]. Several such works also considered the second eigenvalue  
 330 in their analysis [Liu and Brunskill, 2018, Tarbouriech and Lazaric, 2019].

## 331 8 Discussion

332 In this work, we introduced for the first time a differentiable parametric policy that combines TS with  
 333 PG. We proved that SoftTreeMax is essentially a variance reduction technique and explained how to  
 334 choose the expansion policy to minimize the gradient variance. It is an open question whether optimal  
 335 variance reduction corresponds to the appealing regret properties that were put forward by UCT  
 336 [Kocsis and Szepesvári, 2006]. We believe that this can be answered by analyzing the convergence  
 337 rate of SoftTreeMax, relying on the bias and variance results we obtained here.

338 As the learning process continues, the norm of the gradient and the variance *both* become smaller.  
 339 On the face of it, one can ask if the gradient becomes small as fast as the variance or even faster can  
 340 there be any meaningful learning? As we showed in the experiments, learning happens because the  
 341 variance reduces fast enough (a variance of 0 represents deterministic learning, which is fastest).

342 Finally, our work can be extended to infinite action spaces. The analysis can be extended to infinite-  
 343 dimension kernels that retain the same key properties used in our proofs. In the implementation, the  
 344 tree of continuous actions can be expanded by maintaining a parametric distribution over actions that  
 345 depend on  $\theta$ . This approach can be seen as a tree adaptation of MPPI [Williams et al., 2017].

## 346 9 Reproducibility and Limitations

347 In this submission, we include the code as part of the supplementary material. We also include a  
 348 docker file for setting up the environment and a README file with instructions on how to run both  
 349 training and evaluation. The environment engine is an extension of Atari-CuLE [Dalton et al., 2020],  
 350 a CUDA-based Atari emulator that runs on GPU. Our usage of a GPU environment is both a novelty  
 351 and a current limitation of our work.

352 **References**

- 353 A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. On the theory of policy gradient methods:  
354 Optimality, approximation, and distribution shift. *J. Mach. Learn. Res.*, 22(98):1–76, 2021.
- 355 M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An  
356 evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279,  
357 2013.
- 358 L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr. Fully-convolutional siamese  
359 networks for object tracking. In *European conference on computer vision*, pages 850–865. Springer,  
360 2016.
- 361 S. Bhatnagar, R. S. Sutton, M. Ghavamzadeh, and M. Lee. Natural actor–critic algorithms. *Automatica*,  
362 45(11):2471–2482, 2009.
- 363 C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener,  
364 D. Perez, S. Samothrakis, and S. Colton. A survey of monte carlo tree search methods. *IEEE*  
365 *Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.
- 366 S. Chatterjee and E. Seneta. Towards consensus: Some convergence theorems on repeated averaging.  
367 *Journal of Applied Probability*, 14(1):89–97, 1977.
- 368 X. Chen, M. Sim, D. Simchi-Levi, and P. Sun. Risk aversion in inventory management. *Operations*  
369 *Research*, 55(5):828–842, 2007.
- 370 G. Dalal, A. Hallak, S. Dalton, S. Mannor, G. Chechik, et al. Improve agents without retraining:  
371 Parallel tree search with off-policy correction. *Advances in Neural Information Processing Systems*,  
372 34:5518–5530, 2021.
- 373 S. Dalton et al. Accelerating reinforcement learning through gpu atari emulation. *Advances in Neural*  
374 *Information Processing Systems*, 33:19773–19782, 2020.
- 375 Y. Ding, J. Zhang, and J. Lavaei. On the global optimum convergence of momentum-based policy  
376 gradient. In *International Conference on Artificial Intelligence and Statistics*, pages 1910–1934.  
377 PMLR, 2022.
- 378 Y. Efroni, G. Dalal, B. Scherrer, and S. Mannor. How to combine tree-search methods in reinforcement  
379 learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages  
380 3494–3501, 2019.
- 381 Y. Fei, Z. Yang, Y. Chen, and Z. Wang. Exponential bellman equation and improved regret bounds  
382 for risk-sensitive reinforcement learning. *Advances in Neural Information Processing Systems*, 34:  
383 20436–20446, 2021.
- 384 C. D. Freeman, E. Frey, A. Raichuk, S. Girgin, I. Mordatch, and O. Bachem. Brax-a differenti-  
385 able physics engine for large scale rigid body simulation. In *Thirty-fifth Conference on Neural*  
386 *Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.
- 387 J. Garcia and F. Fernández. A comprehensive survey on safe reinforcement learning. *Journal of*  
388 *Machine Learning Research*, 16(1):1437–1480, 2015.
- 389 E. Greensmith, P. L. Bartlett, and J. Baxter. Variance reduction techniques for gradient estimates in  
390 reinforcement learning. *Journal of Machine Learning Research*, 5(9), 2004.
- 391 D. Ha and J. Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 392 R. A. Howard and J. E. Matheson. Risk-sensitive markov decision processes. *Management science*,  
393 18(7):356–369, 1972.
- 394 C. Jin, A. Krishnamurthy, M. Simchowitz, and T. Yu. Reward-free exploration for reinforcement  
395 learning. In *International Conference on Machine Learning*, pages 4870–4879. PMLR, 2020.
- 396 S. W. Kim, Y. Zhou, J. Phillion, A. Torralba, and S. Fidler. Learning to simulate dynamic environments  
397 with gamegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
398 *Recognition*, pages 1231–1240, 2020.

- 399 L. Kocsis and C. Szepesvári. Bandit based monte-carlo planning. In *European conference on machine*  
400 *learning*, pages 282–293. Springer, 2006.
- 401 D. A. Levin and Y. Peres. *Markov chains and mixing times*, volume 107. American Mathematical  
402 Soc., 2017.
- 403 G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen. Softmax policy gradient methods can take exponential  
404 time to converge. In *Conference on Learning Theory*, pages 3107–3110. PMLR, 2021.
- 405 Y. Liu and E. Brunskill. When simple exploration is sample efficient: Identifying sufficient conditions  
406 for random exploration to yield pac rl algorithms. *arXiv preprint arXiv:1805.09045*, 2018.
- 407 Y. Liu, K. Zhang, T. Basar, and W. Yin. An improved analysis of (variance-reduced) policy gradient  
408 and natural policy gradient methods. *Advances in Neural Information Processing Systems*, 33:  
409 7624–7636, 2020.
- 410 V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,  
411 A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot  
412 learning. *arXiv preprint arXiv:2108.10470*, 2021.
- 413 A. S. Mathkar and V. S. Borkar. Nonlinear gossip. *SIAM Journal on Control and Optimization*, 54  
414 (3):1535–1557, 2016.
- 415 J. Mei, C. Xiao, B. Dai, L. Li, C. Szepesvári, and D. Schuurmans. Escaping the gravitational pull of  
416 softmax. *Advances in Neural Information Processing Systems*, 33:21130–21140, 2020a.
- 417 J. Mei, C. Xiao, C. Szepesvari, and D. Schuurmans. On the global convergence rates of softmax  
418 policy gradient methods. In *International Conference on Machine Learning*, pages 6820–6829.  
419 PMLR, 2020b.
- 420 E. Miah, R. MacQueen, A. Ayoub, A. Masoumzadeh, and M. White. Resmax: An alternative  
421 soft-greedy operator for reinforcement learning. 2021.
- 422 V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Ried-  
423 miller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement  
424 learning. *nature*, 518(7540):529–533, 2015.
- 425 E. Noorani and J. S. Baras. Risk-sensitive reinforce: A monte carlo policy gradient algorithm for  
426 exponential performance criteria. In *2021 60th IEEE Conference on Decision and Control (CDC)*,  
427 pages 1522–1527. IEEE, 2021.
- 428 M. Papini, D. Binaghi, G. Canonaco, M. Pirotta, and M. Restelli. Stochastic variance-reduced policy  
429 gradient. In *International conference on machine learning*, pages 4026–4035. PMLR, 2018.
- 430 M. Pelletier. On the almost sure asymptotic behaviour of stochastic algorithms. *Stochastic processes*  
431 *and their applications*, 78(2):217–244, 1998.
- 432 N. Pham, L. Nguyen, D. Phan, P. H. Nguyen, M. Dijk, and Q. Tran-Dinh. A hybrid stochastic  
433 policy gradient algorithm for reinforcement learning. In *International Conference on Artificial*  
434 *Intelligence and Statistics*, pages 374–385. PMLR, 2020.
- 435 A. Raffin, A. Hill, M. Ernestus, A. Gleave, A. Kanervisto, and N. Dormann. Stable baselines3, 2019.
- 436 A. Rosenberg, A. Hallak, S. Mannor, G. Chechik, and G. Dalal. Planning and learning with adaptive  
437 lookahead. *arXiv preprint arXiv:2201.12403*, 2022.
- 438 J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart,  
439 D. Hassabis, T. Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned  
440 model. *Nature*, 588(7839):604–609, 2020.
- 441 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
442 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 443 Z. Shen, A. Ribeiro, H. Hassani, H. Qian, and C. Mi. Hessian aided policy gradient. In *International*  
444 *conference on machine learning*, pages 5729–5738. PMLR, 2019.

- 445 A. Silva, T. Killian, I. D. J. Rodriguez, S.-H. Son, and M. Gombolay. Optimization methods for inter-  
446 pretable differentiable decision trees in reinforcement learning. *arXiv preprint arXiv:1903.09338*,  
447 2019.
- 448 D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser,  
449 I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural  
450 networks and tree search. *nature*, 529(7587):484–489, 2016.
- 451 R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement  
452 learning with function approximation. *Advances in neural information processing systems*, 12,  
453 1999.
- 454 C. Szepesvári. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence  
455 and machine learning*, 4(1):1–103, 2010.
- 456 J. Tarbouriech and A. Lazaric. Active exploration in markov decision processes. In *The 22nd  
457 International Conference on Artificial Intelligence and Statistics*, pages 974–982. PMLR, 2019.
- 458 P. S. Thomas and E. Brunskill. Policy gradient methods for reinforcement learning with function  
459 approximation and action-dependent baselines. *arXiv preprint arXiv:1706.06643*, 2017.
- 460 G. Williams, N. Wagener, B. Goldfain, P. Drews, J. M. Rehg, B. Boots, and E. A. Theodorou.  
461 Information theoretic mpc for model-based reinforcement learning. In *2017 IEEE International  
462 Conference on Robotics and Automation (ICRA)*, pages 1714–1721. IEEE, 2017.
- 463 C. Wu, A. Rajeswaran, Y. Duan, V. Kumar, A. M. Bayen, S. Kakade, I. Mordatch, and P. Abbeel.  
464 Variance reduction for policy gradient with action-dependent factorized baselines. In *International  
465 Conference on Learning Representations*, 2018.
- 466 P. Xu, F. Gao, and Q. Gu. An improved convergence analysis of stochastic variance-reduced policy  
467 gradient. In *Uncertainty in Artificial Intelligence*, pages 541–551. PMLR, 2020.
- 468 J. Zhang, T. He, S. Sra, and A. Jadbabaie. Why gradient clipping accelerates training: A theoretical  
469 justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019.
- 470 J. Zhang, C. Ni, C. Szepesvari, M. Wang, et al. On the convergence and sample efficiency of  
471 variance-reduced policy gradient method. *Advances in Neural Information Processing Systems*, 34:  
472 2228–2240, 2021.

## 473 Appendix

### 474 A Proofs

#### 475 A.1 Proof of Lemma 4.1 – Bound on the policy gradient variance

476 For any parametric policy  $\pi_\theta$  and function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ ,

$$\text{Var}(\nabla_\theta \log \pi_\theta(a|s)Q(s, a)) \leq \max_{s,a} [Q(s, a)]^2 \max_s \|\nabla_\theta \log \pi_\theta(\cdot|s)\|_F^2,$$

477 where  $\nabla_\theta \log \pi_\theta(\cdot|s) \in \mathbb{R}^{A \times \dim(\theta)}$  is a matrix whose  $a$ -th row is  $\nabla_\theta \log \pi_\theta(a|s)^\top$ .

478 *Proof.* The variance for a parametric policy  $\pi_\theta$  is given as follows:

$$\begin{aligned} \text{Var}(\nabla_\theta \log \pi_\theta(a|s)Q(a, s)) &= \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s)^\top \nabla_\theta \log \pi_\theta(a|s)Q(s, a)^2] - \\ &\quad \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s)Q(s, a)]^\top \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s)Q(s, a)], \end{aligned}$$

479 where  $Q(s, a)$  is the currently estimated Q-function and  $d_{\pi_\theta}$  is the discounted state visitation frequency  
480 induced by the policy  $\pi_\theta$ . Since the second term we subtract is always positive (it is of quadratic form  
481  $v^\top v$ ) we can bound the variance by the first term:

$$\begin{aligned}
\text{Var}(\nabla_\theta \log \pi_\theta(a|s)Q(a,s)) &\leq \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s)^\top \nabla_\theta \log \pi_\theta(a|s)Q(s,a)^2] \\
&= \sum_s d_{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s)^\top \nabla_\theta \log \pi_\theta(a|s)Q(s,a)^2 \\
&\leq \max_{s,a} [Q(s,a)]^2 \pi_\theta(a|s) \sum_s d_{\pi_\theta}(s) \sum_a \nabla_\theta \log \pi_\theta(a|s)^\top \nabla_\theta \log \pi_\theta(a|s) \\
&\leq \max_{s,a} [Q(s,a)]^2 \max_s \sum_a \nabla_\theta \log \pi_\theta(a|s)^\top \nabla_\theta \log \pi_\theta(a|s) \\
&= \max_{s,a} [Q(s,a)]^2 \max_s \|\nabla_\theta \log \pi_\theta(\cdot|s)\|_F^2.
\end{aligned}$$

482

□

### 483 A.2 Proof of Lemma 4.2 – Vector form of C-SoftTreeMax

484 In vector form, (3) is given by

$$\pi_{d,\theta}^C(\cdot|s) = \frac{\exp\left[\beta\left(C_{s,d} + P_s(P^{\pi_b})^{d-1}\Theta\right)\right]}{\mathbf{1}_A^\top \exp\left[\beta\left(C_{s,d} + P_s(P^{\pi_b})^{d-1}\Theta\right)\right]}, \quad (8)$$

485 where

$$C_{s,d} = \gamma^{-d}R_s + P_s \left[ \sum_{h=1}^{d-1} \gamma^{h-d} (P^{\pi_b})^{h-1} \right] R_{\pi_b}. \quad (9)$$

486 *Proof.* Consider the vector  $\ell_{s,\cdot} \in \mathbb{R}^{|\mathcal{A}|}$ . Its expectation satisfies

$$\begin{aligned}
\mathbb{E}^{\pi_b} \ell_{s,\cdot}(d;\theta) &= \mathbb{E}^{\pi_b} \left[ \sum_{t=0}^{d-1} \gamma^{t-d} r_t + \theta(s_d) \right] \\
&= \gamma^{-d}R_s + \sum_{t=1}^{d-1} \gamma^{t-d} P_s (P^{\pi_b})^{t-1} R_{\pi_b} + P_s (P^{\pi_b})^{d-1} \Theta.
\end{aligned}$$

487 As required.

□

### 488 A.3 Proof of Lemma 4.3 – Gradient of C-SoftTreeMax

489 The C-SoftTreeMax gradient of dimension  $A \times S$  is given by

$$\nabla_\theta \log \pi_{d,\theta}^C = \beta [I_A - \mathbf{1}_A (\pi_{d,\theta}^C)^\top] P_s (P^{\pi_b})^{d-1},$$

490 where for brevity, we drop the  $s$  index in the policy above, i.e.,  $\pi_{d,\theta}^C \equiv \pi_{d,\theta}^C(\cdot|s)$ .

491 *Proof.* The  $(j,k)$ -th entry of  $\nabla_\theta \log \pi_{d,\theta}^C$  satisfies

$$\begin{aligned}
[\nabla_\theta \log \pi_{d,\theta}^C]_{j,k} &= \frac{\partial \log(\pi_{d,\theta}^C(a^j|s))}{\partial \theta(s^k)} \\
&= \beta [P_s (P^{\pi_b})^{d-1}]_{j,k} - \frac{\sum_a \left[ \exp\left[\beta\left(C_{s,d} + P_s(P^{\pi_b})^{d-1}\Theta\right)\right] \right]_a \beta [P_s (P^{\pi_b})^{d-1}]_{a,k}}{\mathbf{1}_A^\top \exp\left[\beta\left(C_{s,d} + P_s(P^{\pi_b})^{d-1}\Theta\right)\right]} \\
&= \beta [P_s (P^{\pi_b})^{d-1}]_{j,k} - \beta \sum_a \pi_{d,\theta}^C(a|s) [P_s (P^{\pi_b})^{d-1}]_{a,k} \\
&= \beta [P_s (P^{\pi_b})^{d-1}]_{j,k} - \beta [(\pi_{d,\theta}^C)^\top P_s (P^{\pi_b})^{d-1}]_k \\
&= \beta [P_s (P^{\pi_b})^{d-1}]_{j,k} - \beta [\mathbf{1}_A (\pi_{d,\theta}^C)^\top P_s (P^{\pi_b})^{d-1}]_{j,k}.
\end{aligned}$$

492 Moving back to matrix form, we obtain the stated result.

□

493 **A.4 Proof of Theorem 4.4 – Exponential variance decay of C-SoftTreeMax**

494 The C-SoftTreeMax policy gradient is bounded by

$$\text{Var} \left( \nabla_{\theta} \log \pi_{d,\theta}^{\text{C}}(a|s) Q(s, a) \right) \leq 2 \frac{A^2 S^2 \beta^2}{(1-\gamma)^2} |\lambda_2(P^{\pi_b})|^{2(d-1)}.$$

495 *Proof.* We use Lemma 4.1 directly. First of all, it is known that when the reward is bounded in  $[0, 1]$ ,  
 496 the maximal value of the Q-function is  $\frac{1}{1-\gamma}$  as the sum of infinite discounted rewards. Next, we  
 497 bound the Frobenius norm of the term achieved in Lemma 4.3, by applying the eigen-decomposition  
 498 on  $P^{\pi_b}$ :

$$P^{\pi_b} = \mathbf{1}_S \mu^{\top} + \sum_{i=2}^S \lambda_i u_i v_i^{\top}, \quad (10)$$

499 where  $\mu$  is the stationary distribution of  $P^{\pi_b}$ , and  $u_i$  and  $v_i$  are left and right eigenvectors correspond-  
 500 ingly.

$$\begin{aligned} \|\beta (I_{A,A} - \mathbf{1}_A \pi^{\top}) P_s (P^{\pi_b})^{d-1}\|_F &= \beta \|(I_{A,A} - \mathbf{1}_A \pi^{\top}) P_s \left( \mathbf{1}_S \mu^{\top} + \sum_{i=2}^S \lambda_i^{d-1} u_i v_i^{\top} \right)\|_F \\ \text{($P_s$ is stochastic)} &= \beta \|(I_{A,A} - \mathbf{1}_A \pi^{\top}) \left( \mathbf{1}_A \mu^{\top} + \sum_{i=2}^S \lambda_i^{d-1} P_s u_i v_i^{\top} \right)\|_F \\ \text{($projection$ nullifies $\mathbf{1}_A \mu^{\top}$)} &= \beta \|(I_{A,A} - \mathbf{1}_A \pi^{\top}) \left( \sum_{i=2}^S \lambda_i^{d-1} P_s u_i v_i^{\top} \right)\|_F \\ \text{($triangle$ inequality)} &\leq \beta \sum_{i=2}^S \|(I_{A,A} - \mathbf{1}_A \pi^{\top}) (\lambda_i^{d-1} P_s u_i v_i^{\top})\|_F \\ \text{($matrix$ norm sub-multiplicativity)} &\leq \beta |\lambda_2^{d-1}| \sum_{i=2}^S \|I_{A,A} - \mathbf{1}_A \pi^{\top}\|_F \|P_s\|_F \|u_i v_i^{\top}\|_F \\ &= \beta |\lambda_2^{d-1}| (S-1) \|I_{A,A} - \mathbf{1}_A \pi^{\top}\|_F \|P_s\|_F. \end{aligned}$$

501 Now, we can bound the norm  $\|I_{A,A} - \mathbf{1}_A \pi^{\top}\|_F$  by direct calculation:

$$\|I_{A,A} - \mathbf{1}_A \pi^{\top}\|_F^2 = \text{Tr} \left[ (I_{A,A} - \mathbf{1}_A \pi^{\top}) (I_{A,A} - \mathbf{1}_A \pi^{\top})^{\top} \right] \quad (11)$$

$$= \text{Tr} \left[ I_{A,A} - \mathbf{1}_A \pi^{\top} - \pi \mathbf{1}_A^{\top} + \pi^{\top} \pi \mathbf{1}_A \mathbf{1}_A^{\top} \right] \quad (12)$$

$$= A - 1 - 1 + A \pi^{\top} \pi \quad (13)$$

$$\leq 2A. \quad (14)$$

502 From the Cauchy-Schwartz inequality,

$$\|P_s\|_F^2 = \sum_a \sum_s [[P_s]_{a,s}]^2 = \sum_a \|[P_s]_{a,\cdot}\|_2^2 \leq \sum_a \|[P_s]_{a,\cdot}\|_1 \|[P_s]_{a,\cdot}\|_{\infty} \leq A.$$

503 So,

$$\begin{aligned} \text{Var} \left( \nabla_{\theta} \log \pi_{d,\theta}^{\text{C}}(a|s) Q(s, a) \right) &\leq \max_{s,a} [Q(s, a)]^2 \max_s \|\nabla_{\theta} \log \pi_{d,\theta}^{\text{C}}(\cdot|s)\|_F^2 \\ &\leq \frac{1}{(1-\gamma)^2} \|\beta (I_{A,A} - \mathbf{1}_A \pi^{\top}) P_s (P^{\pi_b})^{d-1}\|_F^2 \\ &\leq \frac{1}{(1-\gamma)^2} \beta^2 |\lambda_2(P^{\pi_b})|^{2(d-1)} S^2 (2A^2), \end{aligned}$$

504 which obtains the desired bound.  $\square$

505 **A.5 A lower bound on C-SoftTreeMax gradient (result not in the paper)**

506 For completeness we also supply a lower bound on the Frobenius norm of the gradient. Note that  
 507 this result does not translate to the a lower bound on the variance since we have no lower bound  
 508 equivalence of Lemma 4.1.

509 **Lemma A.1.** *The Frobenius norm on the gradient of the policy is lower-bounded by:*

$$\|\nabla_{\theta} \log \pi_{d,\theta}^C(\cdot|s)\|_F \geq C \cdot \beta |\lambda_2(P^{\pi_b})|^{(d-1)}. \quad (15)$$

510 *Proof.* We begin by moving to the induced  $l_2$  norm by norm-equivalence:

$$\|\beta (I_{A,A} - \mathbf{1}_A \pi^\top) P_s (P^{\pi_b})^{d-1}\|_F \geq \|\beta (I_{A,A} - \mathbf{1}_A \pi^\top) P_s (P^{\pi_b})^{d-1}\|_2.$$

511 Now, taking the vector  $u$  to be the eigenvector of the second eigenvalue of  $P^{\pi_b}$ :

$$\begin{aligned} \|\beta (I_{A,A} - \mathbf{1}_A \pi^\top) P_s (P^{\pi_b})^{d-1}\|_2 &\geq \|\beta (I_{A,A} - \mathbf{1}_A \pi^\top) P_s (P^{\pi_b})^{d-1} u\|_2 \\ &= \beta \|(I_{A,A} - \mathbf{1}_A \pi^\top) P_s u\|_2 \\ &= \beta |\lambda_2(P^{\pi_b})|^{(d-1)} \|(I_{A,A} - \mathbf{1}_A \pi^\top) P_s u\|_2. \end{aligned}$$

512 Note that even though  $P_s u$  can be 0, that is not the common case since we can freely change  $\pi_b$  (and  
 513 therefore the eigenvectors of  $P^{\pi_b}$ ).  $\square$

514 **A.6 Proof of Lemma 4.5 – Vector form of E-SoftTreeMax**

515 For  $d \geq 1$ , (4) is given by

$$\pi_{d,\theta}^E(\cdot|s) = \frac{E_{s,d} \exp(\beta \Theta)}{1_A^\top E_{s,d} \exp(\beta \Theta)}, \quad (16)$$

516 where

$$E_{s,d} = P_s \prod_{h=1}^{d-1} (D(\exp[\beta \gamma^{h-d} R]) P^{\pi_b}) \quad (17)$$

517 with  $R$  being the  $|S|$ -dimensional vector whose  $s$ -th coordinate is  $r(s)$ .

518 *Proof.* Recall that

$$\ell_{s,a}(d; \theta) = \gamma^{-d} \left[ r(s) + \sum_{t=1}^{d-1} \gamma^t r(s_t) + \gamma^d \theta(s_d) \right]. \quad (18)$$

519 and, hence,

$$\exp[\beta \ell_{s,a}(d; \theta)] = \exp \left[ \beta \gamma^{-d} \left( r(s) + \sum_{t=1}^{d-1} \gamma^t r(s_t) + \gamma^d \theta(s_d) \right) \right]. \quad (19)$$

520 Therefore,

$$\mathbb{E}[\exp \beta \ell_{s,a}(d; \theta)] = \mathbb{E} \left[ \exp \left[ \beta \gamma^{-d} \left( r(s) + \sum_{t=1}^{d-1} \gamma^t r(s_t) \right) \right] \mathbb{E}[\exp[\beta(\theta(s_d))]|s_1, \dots, s_{d-1}] \right] \quad (20)$$

$$= \mathbb{E} \left[ \exp \left[ \beta \gamma^{-d} \left( r(s) + \sum_{t=1}^{d-1} \gamma^t r(s_t) \right) \right] P^{\pi_b}(\cdot|s_{d-1}) \right] \exp(\beta \Theta) \quad (21)$$

$$= \mathbb{E} \left[ \exp \left[ \beta \gamma^{-d} \left( r(s) + \sum_{t=1}^{d-2} \gamma^t r(s_t) \right) \right] \exp[\beta \gamma^{-1} r(s_{d-1})] P^{\pi_b}(\cdot|s_{d-1}) \right] \exp(\beta \Theta). \quad (22)$$

521 By repeatedly using iterative conditioning as above, the desired result follows. Note that  
 522  $\exp(\beta \gamma^{-d} r(s))$  does not depend on the action and is therefore cancelled out with the denomi-  
 523 nator.  $\square$

524 **A.7 Proof of Lemma 4.6 – Gradient of E-SoftTreeMax**

525 The E-SoftTreeMax gradient of dimension  $A \times S$  is given by

$$\nabla_{\theta} \log \pi_{d,\theta}^E = \beta [I_A - \mathbf{1}_A (\pi_{d,\theta}^E)^\top] \frac{D(\pi_{d,\theta}^E)^{-1} E_{s,d} D(\exp(\beta\Theta))}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)},$$

526 where for brevity, we drop the  $s$  index in the policy above, i.e.,  $\pi_{d,\theta}^E \equiv \pi_{d,\theta}^E(\cdot|s)$ .

527 *Proof.* The  $(j, k)$ -th entry of  $\nabla_{\theta} \log \pi_{d,\theta}^E$  satisfies

$$\begin{aligned} [\nabla_{\theta} \log \pi_{d,\theta}^E]_{j,k} &= \frac{\partial \log(\pi_{d,\theta}^E(a^j|s))}{\partial \theta(s^k)} \\ &= \frac{\partial}{\partial \theta(s^k)} \left( \log[(E_{s,d})_j^\top \exp(\beta\Theta)] - \log[\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)] \right) \\ &= \frac{\beta (E_{s,d})_{j,k} \exp(\beta\theta(s^k))}{(E_{s,d})_j^\top \exp(\beta\Theta)} - \frac{\beta \mathbf{1}_A^\top E_{s,d} e_k \exp(\beta\theta(s^k))}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \\ &= \frac{\beta (E_{s,d} e_k \exp(\beta\theta(s^k)))_j}{(E_{s,d})_j^\top \exp(\beta\Theta)} - \frac{\beta \mathbf{1}_A^\top E_{s,d} e_k \exp(\beta\theta(s^k))}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \\ &= \beta \left[ \frac{e_j^\top}{e_j^\top E_{s,d} \exp(\beta\Theta)} - \frac{\mathbf{1}_A^\top}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \right] E_{s,d} e_k \exp(\beta\theta(s^k)). \end{aligned}$$

528 Hence,

$$[\nabla_{\theta} \log \pi_{d,\theta}^E]_{\cdot,k} = \beta \left[ D(E_{s,d} \exp(\beta\Theta))^{-1} - (\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta))^{-1} \mathbf{1}_A \mathbf{1}_A^\top \right] E_{s,d} e_k \exp(\beta\theta(s^k))$$

529 From this, it follows that

$$\nabla_{\theta} \log \pi_{d,\theta}^E = \beta \left[ D(\pi_{d,\theta}^E)^{-1} - \mathbf{1}_A \mathbf{1}_A^\top \right] \frac{E_{s,d} D(\exp(\beta\Theta))}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)}. \quad (23)$$

530 The desired result is now easy to see. □

531 **A.8 Proof of Theorem 4.7 — Exponential variance decay of E-SoftTreeMax**

532 There exists  $\alpha \in (0, 1)$  such that, for any function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ ,

$$\text{Var}(\nabla_{\theta} \log \pi_{d,\theta}^E(a|s) Q(s, a)) \in \mathcal{O}(\beta^2 \alpha^{2d}).$$

533 If all rewards are equal ( $r \equiv \text{const}$ ), then  $\alpha = |\lambda_2(P^{\pi_b})|$ .

534 *Proof outline.* Recall that thanks to Lemma 4.1, we can bound the PG variance using a direct bound  
535 on the gradient norm. The definition of the induced norm is

$$\|\nabla_{\theta} \log \pi_{d,\theta}^E\| = \max_{z: \|z\|=1} \|\nabla_{\theta} \log \pi_{d,\theta}^E z\|,$$

536 with  $\nabla_{\theta} \log \pi_{d,\theta}^E$  given in Lemma 4.6. Let  $z \in \mathbb{R}^S$  be an arbitrary vector such that  $\|z\| = 1$ . Then,  
537  $z = \sum_{i=1}^S c_i z_i$ , where  $c_i$  are scalar coefficients and  $z_i$  are vectors spanning the  $S$ -dimensional space.  
538 In the full proof, we show our specific choice of  $z_i$  and prove they are linearly independent given that  
539 choice. We do note that  $z_1 = \mathbf{1}_S$ .

The first part of the proof relies on the fact that  $(\nabla_{\theta} \log \pi_{d,\theta}^E) z_1 = 0$ . This is easy to verify using  
Lemma 4.6 together with (6), and because  $[I_A - \mathbf{1}_A (\pi_{d,\theta}^E)^\top]$  is a projection matrix whose null-space  
is spanned by  $\mathbf{1}_S$ . Thus,

$$\nabla_{\theta} \log \pi_{d,\theta}^E z = \nabla_{\theta} \log \pi_{d,\theta}^E \sum_{i=2}^S c_i z_i.$$

540 In the second part of the proof, we focus on  $E_{s,d}$  from (6), which appears within  $\nabla_\theta \log \pi_{d,\theta}^E$ . Notice  
 541 that  $E_{s,d}$  consists of the product  $\prod_{h=1}^{d-1} (D(\exp(\beta\gamma^{h-d}R)P^{\pi_b}))$ . Even though the elements in this  
 542 product are not stochastic matrices, in the full proof we show how to normalize each of them to a  
 543 stochastic matrix  $B_h$ . We thus obtain that

$$E_{s,d} = P_s D(M_1) \prod_{h=1}^{d-1} B_h,$$

544 where  $M_1 \in \mathbb{R}^S$  is some strictly positive vector. Then, we can apply a result by Mathkar and Borkar  
 545 [2016], which itself builds on [Chatterjee and Seneta, 1977]. The result states that the product of  
 546 stochastic matrices  $\prod_{h=1}^{d-1} B_h$  of our particular form converges exponentially fast to a matrix of the  
 547 form  $\mathbf{1}_S \mu^\top$  s.t.  $\|\mathbf{1}_S \mu^\top - \prod_{h=1}^{d-1} B_h\| \leq C\alpha^d$  for some constant  $C$ .

548 Lastly,  $\mathbf{1}_S \mu_{\pi_b}^\top$  gets canceled due to our choice of  $z_i$ ,  $i = 2, \dots, S$ . This observation along with the  
 549 above fact that the remainder decays then shows that  $\nabla_\theta \log \pi_{d,\theta}^E \sum_{i=2}^S z_i = \mathcal{O}(\alpha^d)$ , which gives the  
 550 desired result.  $\square$

551 *Full technical proof.* Let  $d \geq 2$ . Recall that

$$E_{s,d} = P_s \prod_{h=1}^{d-1} (D(\exp[\beta\gamma^{h-d}R])P^{\pi_b}), \quad (24)$$

552 and that  $R$  refers to the  $S$ -dimensional vector whose  $s$ -th coordinate is  $r(s)$ . Define

$$B_i = \begin{cases} P^{\pi_b} & \text{if } i = d-1, \\ D^{-1}(P^{\pi_b} M_{i+1})P^{\pi_b} D(M_{i+1}) & \text{if } i = 1, \dots, d-2, \end{cases} \quad (25)$$

553 and the vector

$$M_i = \begin{cases} \exp(\beta\gamma^{-1}R) & \text{if } i = d-1, \\ \exp(\beta\gamma^{i-d}R) \circ P^{\pi_b} M_{i+1} & \text{if } i = 1, \dots, d-2, \end{cases} \quad (26)$$

554 where  $\circ$  denotes the element-wise product. Then,

$$E_{s,d} = P_s D(M_1) \prod_{i=1}^{d-1} B_i. \quad (27)$$

555 It is easy to see that each  $B_i$  is a row-stochastic matrix, i.e., all entries are non-negative and  
 556  $B_i \mathbf{1}_S = \mathbf{1}_S$ .

557 Next, we prove that all non-zeros entries of  $B_i$  are bounded away from 0 by a constant. This is  
 558 necessary to apply the next result from Chatterjee and Seneta [1977]. The  $j$ -th coordinate of  $M_i$   
 559 satisfies

$$(M_i)_j = \exp[\beta\gamma^{i-d}R_j] \sum_k [P^{\pi_b}]_{j,k}(M_{i+1})_k \leq \|\exp[\beta\gamma^{i-d}R]\|_\infty \|M_{i+1}\|_\infty. \quad (28)$$

560 Separately, observe that  $\|M_{d-1}\|_\infty \leq \|\exp(\beta\gamma^{-1}R)\|_\infty$ . Plugging these relations in (26) gives

$$\|M_1\|_\infty \leq \prod_{h=1}^{d-1} \|\exp[\beta\gamma^{h-d}R]\|_\infty = \prod_{h=1}^{d-1} \|\exp[\beta\gamma^{-d}R]\|_\infty^{\gamma^h} = \|\exp[\beta\gamma^{-d}R]\|_\infty^{\sum_{h=1}^{d-1} \gamma^h} \leq \|\exp[\beta\gamma^{-d}R]\|_\infty^{\frac{1}{1-\gamma}}. \quad (29)$$

561 Similarly, for every  $1 \leq i \leq d-1$ , we have that

$$\|M_i\|_\infty \leq \prod_{h=i}^{d-1} \|\exp[\beta\gamma^{-d}R]\|_\infty^{\gamma^h} \leq \|\exp[\beta\gamma^{-d}R]\|_\infty^{\frac{1}{1-\gamma}}. \quad (30)$$

562 The  $jk$ -th entry of  $B_i = D^{-1}(P^{\pi_b} M_{i+1})P^{\pi_b} D(M_{i+1})$  is

$$(B_i)_{jk} = \frac{P_{jk}^{\pi_b} [M_{i+1}]_k}{\sum_{\ell=1}^{|S|} P_{j\ell}^{\pi_b} [M_{i+1}]_\ell} \geq \frac{P_{jk}^{\pi_b}}{\sum_{\ell=1}^{|S|} P_{j\ell}^{\pi_b} [M_{i+1}]_\ell} \geq \frac{P_{jk}^{\pi_b}}{\|\exp[\beta\gamma^{-d}R]\|_\infty^{\frac{1}{1-\gamma}}}. \quad (31)$$

563 Hence, for non-zero  $P_{jk}^{\pi_b}$ , the entries are bounded away from zero by the same. We can now proceed  
 564 with applying the following result.

565 Now, by [Chatterjee and Seneta, 1977, Theorem 5] (see also (14) in [Mathkar and Borkar, 2016]),  
 566  $\lim_{d \rightarrow \infty} \prod_{i=1}^{d-1} B_i$  exists and is of the form  $\mathbf{1}_S \mu^\top$  for some probability vector  $\mu$ . Furthermore, there  
 567 is some  $\alpha \in (0, 1)$  such that  $\varepsilon(d) := \left( \prod_{i=1}^{d-1} B_i \right) - \mathbf{1}_S \mu^\top$  satisfies

$$\|\varepsilon(d)\| = O(\alpha^d). \quad (32)$$

568 Pick linearly independent vectors  $w_2, \dots, w_S$  such that

$$\mu^\top w_i = 0 \text{ for } i = 2, \dots, d. \quad (33)$$

569 Since  $\sum_{i=2}^S \alpha_i w_i$  is perpendicular to  $\mu$  for any  $\alpha_2, \dots, \alpha_S$  and because  $\mu^\top \exp(\beta\Theta) > 0$ , there  
 570 exists no choice of  $\alpha_2, \dots, \alpha_S$  such that  $\sum_{i=2}^S \alpha_i w_i = \exp(\beta\Theta)$ . Hence, if we let  $z_1 = \mathbf{1}_S$  and  
 571  $z_i = D(\exp(\beta\Theta))^{-1} w_i$  for  $i = 2, \dots, S$ , then it follows that  $\{z_1, \dots, z_S\}$  is linearly independent.  
 572 In particular, it implies that  $\{z_1, \dots, z_S\}$  spans  $\mathbb{R}^S$ .

573 Now consider an arbitrary unit norm vector  $z := \sum_{i=1}^S c_i z_i \in \mathbb{R}^S$  s.t.  $\|z\|_2 = 1$ . Then,

$$\nabla_\theta \log \pi_{d,\theta}^E z = \nabla_\theta \log \pi_{d,\theta}^E \sum_{i=2}^S c_i z_i \quad (34)$$

$$= \beta [I_A - \mathbf{1}_A (\pi_{d,\theta}^E)^\top] \frac{D(\pi_{d,\theta}^E)^{-1} E_{s,d} D(\exp(\beta\Theta))}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \sum_{i=2}^S c_i z_i \quad (35)$$

$$= \beta [I_A - \mathbf{1}_A (\pi_{d,\theta}^E)^\top] \frac{D(\pi_{d,\theta}^E)^{-1} E_{s,d}}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \sum_{i=2}^S c_i w_i \quad (36)$$

$$= \beta [I_A - \mathbf{1}_A (\pi_{d,\theta}^E)^\top] \frac{D(\pi_{d,\theta}^E)^{-1} [\mathbf{1}_S \mu^\top + \varepsilon(d)]}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \sum_{i=2}^S c_i w_i \quad (37)$$

$$= \beta [I_A - \mathbf{1}_A (\pi_{d,\theta}^E)^\top] \frac{D(\pi_{d,\theta}^E)^{-1} \varepsilon(d)}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \sum_{i=2}^S c_i w_i \quad (38)$$

$$= \beta [I_A - \mathbf{1}_A (\pi_{d,\theta}^E)^\top] \frac{D(\pi_{d,\theta}^E)^{-1} \varepsilon(d) D(\exp(\beta\Theta))}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} (z - c_1 \mathbf{1}_S), \quad (39)$$

574 where (34) follows from the fact that  $\nabla_\theta \log \pi_{d,\theta}^E z_1 = \nabla_\theta \log \pi_{d,\theta}^E \mathbf{1}_S = 0$ , (35) follows from  
 575 Lemma 4.6, (36) holds since  $z_i = D(\exp(\beta\Theta))^{-1} w_i$ , (38) because  $\mu$  is perpendicular  $w_i$  for each  $i$ ,  
 576 while (39) follows by reusing  $z_i = D(\exp(\beta\Theta))^{-1} w_i$  relation along with the fact that  $z_1 = \mathbf{1}_S$ .

577 From (39), it follows that

$$\|\nabla_{\theta} \log \pi_{d,\theta}^E z\| \leq \beta \|\varepsilon(d)\| \left\| [I_A - \mathbf{1}_A(\pi_{d,\theta}^E)^\top] \frac{D(\pi_{d,\theta}^E)^{-1}}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \right\| \|D(\exp(\beta\Theta))\| \|z - c_1 \mathbf{1}_S\| \quad (40)$$

$$\leq \beta \alpha^d (\|I_A\| + \|\mathbf{1}_A(\pi_{d,\theta}^E)^\top\|) \left\| \frac{D(\pi_{d,\theta}^E)^{-1}}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \right\| \exp(\beta \max_s \theta(s)) \|z - c_1 \mathbf{1}_S\| \quad (41)$$

$$\leq \beta \alpha^d (1 + \sqrt{A}) \left\| \frac{D(\pi_{d,\theta}^E)^{-1}}{\mathbf{1}_A^\top E_{s,d} \exp(\beta\Theta)} \right\| \exp(\beta \max_s \theta(s)) \|z - c_1 \mathbf{1}_S\| \quad (42)$$

$$\leq \beta \alpha^d (1 + \sqrt{A}) \|D^{-1}(E_{s,d} \exp(\beta\Theta))\| \exp(\beta \max_s \theta(s)) \|z - c_1 \mathbf{1}_S\| \quad (43)$$

$$\leq \beta \alpha^d (1 + \sqrt{A}) \frac{1}{\min_s [E_{s,d} \exp(\beta\Theta)]_s} \exp(\beta \max_s \theta(s)) \|z - c_1 \mathbf{1}_S\| \quad (44)$$

$$\leq \beta \alpha^d (1 + \sqrt{A}) \frac{\exp(\beta \max_s \theta(s))}{\exp(\beta \min_s \theta(s)) \min_s |M_1|} \|z - c_1 \mathbf{1}_S\| \quad (45)$$

$$\leq \beta \alpha^d (1 + \sqrt{A}) \frac{\exp(\beta \max_s \theta(s))}{\exp(\beta \min_s \theta(s)) \exp(\beta \min_s r(s))} \|z - c_1 \mathbf{1}_S\| \quad (46)$$

$$\leq \beta \alpha^d (1 + \sqrt{A}) \exp(\beta [\max_s \theta(s) - \min_s \theta(s) - \min_s r(s)]) \|z - c_1 \mathbf{1}_S\|. \quad (47)$$

578 Lastly, we prove that  $\|z - c_1 \mathbf{1}_S\|$  is bounded independently of  $d$ . First, denote by  $c = (c_1, \dots, c_S)^\top$   
579 and  $\tilde{c} = (0, c_2, \dots, c_S)^\top$ . Also, denote by  $Z$  the matrix with  $z_i$  as its  $i$ -th column. Now,

$$\|z - c_1 \mathbf{1}_S\| = \left\| \sum_{i=2}^S c_i z_i \right\| \quad (48)$$

$$= \|Z\tilde{c}\| \quad (49)$$

$$\leq \|Z\| \|\tilde{c}\| \quad (50)$$

$$\leq \|Z\| \|c\| \quad (51)$$

$$= \|Z\| \|Z^{-1}z\| \quad (52)$$

$$\leq \|Z\| \|Z^{-1}\|, \quad (53)$$

580 where the last relation is due to  $z$  being a unit vector. All matrix norms here are  $l_2$ -induced norms.

581 Next, denote by  $W$  the matrix with  $w_i$  in its  $i$ -th column. Recall that in (33) we only defined  
582  $w_2, \dots, w_S$ . We now set  $w_1 = \exp(\beta\Theta)$ . Note that  $w_1$  is linearly independent of  $\{w_2, \dots, w_S\}$   
583 because of (33) together with the fact that  $\mu^\top w_1 > 0$ . We can now express the relation between  $Z$   
584 and  $W$  by  $Z = D^{-1}(\exp(\beta\Theta))W$ . Substituting this in (53), we have

$$\|z - c_1 \mathbf{1}_S\| \leq \|D^{-1}(\exp(\beta\Theta))W\| \|W^{-1}D(\exp(\beta\Theta))\| \quad (54)$$

$$\leq \|W\| \|W^{-1}\| \|D(\exp(\beta\Theta))\| \|D^{-1}(\exp(\beta\Theta))\|. \quad (55)$$

585 It further holds that

$$\|D(\exp(\beta\Theta))\| \leq \max_s \exp(\beta\theta(s)) \leq \max\{1, \exp[\beta \max_s \theta(s)]\}, \quad (56)$$

586 where the last relation equals 1 if  $\theta(s) < 0$  for all  $s$ . Similarly,

$$\|D^{-1}(\exp(\beta\Theta))\| \leq \frac{1}{\min_s \exp(\beta\theta(s))} \leq \frac{1}{\min\{1, \exp[\beta \min_s \theta(s)]\}}. \quad (57)$$

587 Furthermore, by the properties of the  $l_2$ -induced norm,

$$\|W\|_2 \leq \sqrt{S}\|W\|_1 \quad (58)$$

$$= \sqrt{S} \max_{1 \leq i \leq S} \|w_i\|_1 \quad (59)$$

$$= \sqrt{S} \max\{\exp(\beta\Theta), \max_{2 \leq i \leq S} \|w_i\|_1\} \quad (60)$$

$$\leq \sqrt{S} \max\{1, \exp[\beta \max_s \theta(s)], \max_{2 \leq i \leq S} \|w_i\|_1\}. \quad (61)$$

588 Lastly,

$$\|W^{-1}\| = \frac{1}{\sigma_{\min}(W)} \quad (62)$$

$$\leq \left( \prod_{i=1}^{S-1} \frac{\sigma_{\max}(W)}{\sigma_i(W)} \right) \frac{1}{\sigma_{\min}(W)} \quad (63)$$

$$= \frac{(\sigma_{\max}(W))^{S-1}}{\prod_{i=1}^S \sigma_i(W)} \quad (64)$$

$$= \frac{\|W\|^{S-1}}{|\det(W)|}. \quad (65)$$

589 The determinant of  $W$  is a sum of products involving its entries. To upper bound (65) independently  
 590 of  $d$ , we lower bound its denominator by upper and lower bounds on the entries  $[W]_{i,1}$  that are  
 591 independent of  $d$ , depending on their sign:

$$\min\{1, \exp[\beta \min_s \theta(s)]\} \leq [W]_{i,1} \leq \max\{1, \exp[\beta \max_s \theta(s)]\}. \quad (66)$$

592 Using this, together with (53), (55), (56), (57), and (61), we showed that  $\|z - c_1 \mathbf{1}_S\|$  is upper bounded  
 593 by a constant independent of  $d$ . This concludes the proof.  $\square$

## 594 A.9 Bias Estimates

**Lemma A.2.** For any matrix  $A$  and  $\hat{A}$ ,

$$\hat{A}^k - A^k = \sum_{h=1}^k \hat{A}^{h-1} (\hat{A} - A) A^{k-h}.$$

595 *Proof.* The proof follows from first principles:

$$\sum_{h=1}^k \hat{A}^{h-1} (\hat{A} - A) A^{k-h} = \sum_{h=1}^k \hat{A}^{h-1} \hat{A} A^{k-h} - \sum_{h=1}^k \hat{A}^{h-1} A A^{k-h} \quad (67)$$

$$= \sum_{h=1}^k \hat{A}^h A^{k-h} - \sum_{h=1}^k \hat{A}^{h-1} A^{k-h+1} \quad (68)$$

$$= \hat{A}^k - A^k + \sum_{h=1}^{k-1} \hat{A}^h A^{k-h} - \sum_{h=2}^k \hat{A}^{h-1} A^{k-h+1} \quad (69)$$

$$= \hat{A}^k - A^k. \quad (70)$$

596  $\square$

597 Henceforth,  $\|\cdot\|$  will refer to  $\|\cdot\|_\infty$ , i.e. the induced infinity norm. Also, for brevity, we denote  $\pi_{d,\theta}^C$   
 598 and  $\hat{\pi}_{d,\theta}^C$  by  $\pi_\theta$  and  $\hat{\pi}_\theta$ , respectively. Similarly, we use  $d_{\pi_\theta}$  and  $d_{\hat{\pi}_\theta}$  to denote  $d_{\pi_{d,\theta}^C}$  and  $d_{\hat{\pi}_{d,\theta}^C}$ . As for  
 599 the induced norm of the matrix  $P$  and its perturbed counterpart  $\hat{P}$ , which are of size  $S \times A \times S$ ,  
 600 we slightly abuse notation and denote  $\|P - \hat{P}\| = \max_s \{\|P_s - \hat{P}_s\|\}$ , where  $P_s$  is as defined in  
 601 Section 2.

602 **Definition A.3.** Let  $\epsilon$  be the maximal model mis-specification, i.e.,  $\max\{\|P - \hat{P}\|, \|r - \hat{r}\|\} = \epsilon$ .

603 **Lemma A.4.** Recall the definitions of  $R_s, P_s, R_{\pi_b}$  and  $P^{\pi_b}$  from Section 2, and respectively denote  
 604 their perturbed counterparts by  $\hat{R}_s, \hat{P}_s, \hat{R}_{\pi_b}$  and  $\hat{P}^{\pi_b}$ . Then, for  $\epsilon$  defined in Definition A.3,

$$\max\{\|R_s - \hat{R}_s\|, \|P_s - \hat{P}_s\|, \|R_{\pi_b} - \hat{R}_{\pi_b}\|, \|P^{\pi_b} - \hat{P}^{\pi_b}\|\} = O(\epsilon). \quad (71)$$

605 *Proof.* The proof follows easily from the fact that the differences above are convex combinations of  
 606  $P - \hat{P}$  and  $r - \hat{r}$ .  $\square$

607 **Lemma A.5.** Let  $\pi_\theta$  be as in (5), and let  $\hat{\pi}_\theta$  also be defined as in (5), but with  $R_s, P_s, P^{\pi_b}$  replaced  
 608 by their perturbed counterparts  $\hat{R}_s, \hat{P}_s, \hat{P}^{\pi_b}$  throughout. Then,

$$\|\pi_{d,\theta}^C - \hat{\pi}_{d,\theta}^C\| = O(\beta d \epsilon). \quad (72)$$

609 *Proof.* To prove the desired result, we work with (5) to bound the error between  $R_s, P_s, P^{\pi_b}, R_{\pi_b}$   
 610 and their perturbed versions.

First, we apply Lemma A.2 together with Lemma A.4 to obtain that  $\|(P^{\pi_b})^k - (\hat{P}^{\pi_b})^k\| = O(k\epsilon)$ .  
 Next, denote by  $M$  the argument in the exponent in (5), i.e.

$$M := \beta[C_{s,d} + P_s(P^{\pi_b})^{d-1}\Theta].$$

611 Similarly, let  $\hat{M}$  be the corresponding perturbed sum that relies on  $\hat{P}$  and  $\hat{r}$ . Combining the bounds  
 612 from Lemma A.4, and using the triangle inequality, we have that  $\|\hat{M} - M\| = O(\beta d \epsilon)$ .

Eq. (5) states that the C-SoftTreeMax policy in the true environment is  $\pi_\theta = \exp(M)/(1^\top \exp(M))$ .  
 Similarly define  $\hat{\pi}_\theta$  using  $\hat{M}$  for the approximate model. Then,

$$\hat{\pi}_\theta = (\pi_\theta \circ \exp(M - \hat{M}))1^\top \exp(M)/(1^\top \exp(\hat{M})),$$

613 where  $\circ$  denotes element-wise multiplication. Using the above relation, we have that  $\|\hat{\pi}_\theta - \pi_\theta\| =$   
 614  $\|\pi_\theta\| \left\| \frac{\exp(M - \hat{M})1^\top \exp(M)}{1^\top \exp(M)} - 1 \right\|$ . Using the relation  $|e^x - 1| = O(x)$  as  $x \rightarrow 0$ , the desired result  
 615 follows.  $\square$

616  $\square$

617 **Theorem A.6.** Let  $\epsilon$  be as in Definition A.3. Further let  $\hat{\pi}_{d,\theta}^C$  being the corresponding approximate  
 618 policy as given in Lemma 4.2. Then, the policy gradient bias is bounded by

$$\left\| \frac{\partial}{\partial \theta} (\nu^\top V^{\pi_\theta}) - \frac{\partial}{\partial \theta} (\nu^\top V^{\hat{\pi}_\theta}) \right\| = \mathcal{O} \left( \frac{1}{(1-\gamma)^2} S \beta^2 d \epsilon \right). \quad (73)$$

619 We first provide a proof outline for conciseness, and only after it the complete proof.

*Proof outline.* First, we prove that  $\max\{\|R_s - \hat{R}_s\|, \|P_s - \hat{P}_s\|, \|R_{\pi_b} - \hat{R}_{\pi_b}\|, \|P^{\pi_b} - \hat{P}^{\pi_b}\|\} = \mathcal{O}(\epsilon)$ .  
 This follows from the fact that the differences above are suitable convex combinations of either the  
 rows of  $P - \hat{P}$  or  $r - \hat{r}$ . We use the above observation along with the definitions of  $\pi_{d,\theta}^C$  and  $\hat{\pi}_{d,\theta}^C$   
 given in (5) to show that  $\|\pi_{d,\theta}^C - \hat{\pi}_{d,\theta}^C\| = O(\beta d \epsilon)$ . The proof for the latter builds upon two key facts:  
 (a)  $\|(P^{\pi_b})^k - (\hat{P}^{\pi_b})^k\| \leq \sum_{h=1}^k \|\hat{P}^{\pi_b}\|^{h-1} \|\hat{P}^{\pi_b} - P^{\pi_b}\| \|P^{\pi_b}\|^{k-h} = O(k\epsilon)$  for any  $k \geq 0$ , and (b)  
 $|e^x - 1| = O(x)$  as  $x \rightarrow 0$ . Next, we decompose the LHS of (7) to get

$$\sum_s \left( \prod_{i=1}^4 X_i(s) - \prod_{i=1}^4 \hat{X}_i(s) \right) = \sum_s \sum_{i=1}^4 \hat{X}_1(s) \cdots \hat{X}_{i-1}(s) (X_i(s) - \hat{X}_i(s)) \times X_{i+1}(s) \cdots X_4(s),$$

620 where  $X_1(s) = d_{\pi_{d,\theta}^C}(s) \in \mathbb{R}$ ,  $X_2(s) = (\nabla_\theta \log \pi_{d,\theta}^C(\cdot|s))^\top \in \mathbb{R}^{S \times A}$ ,  $X_3(s) = D(\pi_{d,\theta}^C(\cdot|s)) \in$   
 621  $\mathbb{R}^{A \times A}$ ,  $X_4(s) = Q^{\pi_{d,\theta}^C}(s, \cdot) \in \mathbb{R}^{A \times A}$ , and  $\hat{X}_1(s), \dots, \hat{X}_4(s)$  are similarly defined with  $\pi_{d,\theta}^C$  re-  
 622 placed by  $\hat{\pi}_{d,\theta}^C$ . Then, we show that, for  $i = 1, \dots, 4$ , (i)  $\|X_i(s) - \hat{X}_i(s)\| = O(\epsilon)$  and (ii)  
 623  $\max\{\|X_i\|, \|\hat{X}_i\|\}$  is bounded by problem parameters. From this, the desired result follows.  $\square$

624 *Proof.* We have

$$\frac{\partial}{\partial \theta} (\nu^\top V^{\pi_\theta}) - \frac{\partial}{\partial \theta} (\nu^\top V^{\hat{\pi}_\theta}) \quad (74)$$

$$= \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a)] - \mathbb{E}_{s \sim d_{\hat{\pi}_\theta}, a \sim \hat{\pi}_\theta(\cdot|s)} [\nabla_\theta \log \hat{\pi}_\theta(a|s) Q^{\hat{\pi}_\theta}(s, a)] \quad (75)$$

$$= \sum_{s, a} (d_{\pi_\theta}(s) \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a) - d_{\hat{\pi}_\theta}(s) \hat{\pi}_\theta(a|s) \nabla_\theta \log \hat{\pi}_\theta(a|s) Q^{\hat{\pi}_\theta}(s, a)) \quad (76)$$

$$= \sum_s \left( d_{\pi_\theta}(s) (\nabla_\theta \log \pi_\theta(\cdot|s))^\top D(\pi_\theta(\cdot|s)) Q^{\pi_\theta}(s, \cdot) \right. \quad (77)$$

$$\left. - d_{\hat{\pi}_\theta}(s) (\nabla_\theta \log \hat{\pi}_\theta(\cdot|s))^\top D(\hat{\pi}_\theta(\cdot|s)) Q^{\hat{\pi}_\theta}(s, \cdot) \right) \quad (78)$$

$$= \sum_s \left( \prod_{i=1}^4 X_i(s) - \prod_{i=1}^4 \hat{X}_i(s) \right) \quad (79)$$

$$= \sum_s \sum_{i=1}^4 \hat{X}_1(s) \cdots \hat{X}_{i-1}(s) (X_i(s) - \hat{X}_i(s)) X_{i+1}(s) \cdots X_4(s), \quad (80)$$

625 where  $X_1(s) = d_{\pi_\theta}(s) \in \mathbb{R}$ ,  $X_2(s) = (\nabla_\theta \log \pi_\theta(\cdot|s))^\top \in \mathbb{R}^{S \times A}$ ,  $X_3(s) = D(\pi_\theta(\cdot|s)) \in \mathbb{R}^{A \times A}$ ,  
 626  $X_4(s) = Q^{\pi_\theta}(s, \cdot) \in \mathbb{R}^{A \times A}$ , and  $\hat{X}_1(s), \dots, \hat{X}_4(s)$  are similarly defined with  $\pi_\theta$  replaced by  $\hat{\pi}_\theta$ .

627 Therefore,

$$\left\| \frac{\partial}{\partial \theta} (\nu^\top V^{\pi_\theta}) - \frac{\partial}{\partial \theta} (\nu^\top V^{\hat{\pi}_\theta}) \right\| \leq \left( \max_s \Gamma(s) \right) S, \quad (81)$$

628 where

$$\Gamma(s) = \left\| \sum_s \sum_{i=1}^4 \hat{X}_1(s) \cdots \hat{X}_{i-1}(s) (X_i(s) - \hat{X}_i(s)) X_{i+1}(s) \cdots X_4(s) \right\|. \quad (82)$$

629 Next, since  $d_{\pi_\theta}$ ,  $d_{\hat{\pi}_\theta}$ ,  $\pi_\theta$ , and  $\hat{\pi}_\theta$  are all distributions, we have

$$\max\{|X_1(s)|, |\hat{X}_1(s)|, |X_3(s, a)|, |\hat{X}_3(s, a)|\} \leq 1. \quad (83)$$

630 Separately, using Lemma 4.3, we have

$$\|X_2\| = \|\nabla_\theta \log \pi_\theta(a|s)\| \leq \beta (\|I_A\| + \|\mathbf{1}_A \pi_\theta^\top\|) \|P_s\| \|(P^{\pi_\theta})^{d-1}\|. \quad (84)$$

631 Since all rows of the above matrices have non-negative entries that add up to 1, we get

$$\|Y\| \leq 2\beta. \quad (85)$$

632 In the rest of the proof, we bound each of  $\|X_1 - \hat{X}_1\|, \dots, \|X_4 - \hat{X}_4\|$ .

633 Finally,

$$\|X_4\| \leq \frac{1}{1-\gamma}. \quad (86)$$

634 Similarly, the same bounds hold for  $\hat{X}_1, \hat{X}_2, \hat{X}_3$  and  $\hat{X}_4$ .

635 From, we have

$$\|X_1 - \hat{X}_1\| \leq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \|\nu^\top (P^{\pi_\theta})^t - \nu^\top (P^{\hat{\pi}_\theta})^t\| \quad (87)$$

$$\leq (1-\gamma) \|\nu\| \sum_{t=0}^{\infty} \gamma^t t d\epsilon \quad (88)$$

$$\leq (1-\gamma) d\epsilon \sum_{t=0}^{\infty} \gamma^t t \quad (89)$$

$$= \frac{\gamma d\epsilon}{1-\gamma}. \quad (90)$$

636 The last relation follows from the fact that  $(1 - \gamma)^{-1} = \sum_{t=0}^{\infty} \gamma^t$ , which in turn implies

$$\gamma \frac{\partial}{\partial \gamma} \left( \frac{1}{1 - \gamma} \right) = \sum_{t=0}^{\infty} t \gamma^t. \quad (91)$$

637 From Lemma A.5, it follows that

$$\|X_3 - \hat{X}_3\| = O(\beta d \epsilon). \quad (92)$$

638 Next, recall that from Lemma 4.3 that

$$X_2(s, \cdot) = \beta [I_A - \mathbf{1}_A(\pi_\theta)^\top] P_s (P^{\pi_b})^{d-1}.$$

639 Then,

$$\|X_2(s, \cdot) - \hat{X}_2(s, \cdot)\| \leq \|\beta [I_A - \mathbf{1}_A(\pi_\theta)^\top] P_s\| \| (P^{\pi_b})^{d-1} - (\hat{P}^{\pi_b})^{d-1} \| \quad (93)$$

$$+ \|\beta [I_A - \mathbf{1}_A(\pi_\theta)^\top]\| \|P_s - \hat{P}_s\| \| (\hat{P}^{\pi_b})^{d-1} \| \quad (94)$$

$$+ \beta \|\mathbf{1}_A(\pi_\theta)^\top - \mathbf{1}_A(\hat{\pi}_\theta)^\top\| \| \hat{P}_s (\hat{P}^{\pi_b})^{d-1} \|. \quad (95)$$

640 Following the same argument as in (85) and applying Lemma A.2, we have that (93) is  $O(\beta d \epsilon)$ .  
 641 Similarly, from the argument of (85), Eq. (94) is  $O(\beta \epsilon)$ . Lastly, (95) is  $O(\beta d \epsilon)$  due to Lemma A.5.  
 642 Putting the above three terms together, we have that

$$\|X_2(s, \cdot) - \hat{X}_2(s, \cdot)\| = O(\beta d \epsilon). \quad (96)$$

643 Since the state-action value function satisfies the Bellman equation, we have

$$Q^{\pi_\theta} = r + \gamma P Q^{\pi_\theta} \quad (97)$$

644 and

$$Q^{\hat{\pi}_\theta} = \hat{r} + \gamma \hat{P} Q^{\hat{\pi}_\theta}. \quad (98)$$

645 Consequently,

$$\|Q^{\pi_\theta} - Q^{\hat{\pi}_\theta}\| \leq \|r - \hat{r}\| + \gamma \|P Q^{\pi_\theta} - P Q^{\hat{\pi}_\theta}\| + \gamma \|P Q^{\hat{\pi}_\theta} - \hat{P} Q^{\hat{\pi}_\theta}\| \quad (99)$$

$$\leq \epsilon + \gamma \|P\| \|Q^{\pi_\theta} - Q^{\hat{\pi}_\theta}\| + \gamma \|P - \hat{P}\| \|Q^{\hat{\pi}_\theta}\| \quad (100)$$

$$\leq \epsilon + \gamma \|Q^{\pi_\theta} - Q^{\hat{\pi}_\theta}\| + \frac{\gamma}{1 - \gamma} \epsilon, \quad (101)$$

646 which finally shows that

$$\|X_4 - \hat{X}_4\| = \|Q^{\pi_\theta} - Q^{\hat{\pi}_\theta}\| \leq \frac{\epsilon}{(1 - \gamma)^2}. \quad (102)$$

647

□

## 648 B Experiments

### 649 B.1 Implementation Details

650 The environment engine is the highly efficient Atari-CuLE [Dalton et al., 2020], a CUDA-based  
 651 version of Atari that runs on GPU. Similarly, we use Atari-CuLE for the GPU-based breadth-first TS  
 652 as done in Dalal et al. [2021]: In every tree expansion, the state  $S_t$  is duplicated and concatenated  
 653 with all possible actions. The resulting tensor is fed into the GPU forward model to generate the  
 654 tensor of next states  $(S_{t+1}^0, \dots, S_{t+1}^{A-1})$ . The next-state tensor is then duplicated and concatenated  
 655 again with all possible actions, fed into the forward model, etc. This procedure is repeated until the  
 656 final depth is reached, for which  $W_\theta(s)$  is applied per state.

657 We train SoftTreeMax for depths  $d = 1 \dots 8$ , with a single worker. We use five seeds for each  
 658 experiment.

659 For the implementation, we extend Stable-Baselines3 [Raffin et al., 2019] with all parameters taken  
 660 as default from the original PPO paper [Schulman et al., 2017]. For depths  $d \geq 3$ , we limited the  
 661 tree to a maximum width of 1024 nodes and pruned non-promising trajectories in terms of estimated  
 662 weights. Since the distributed PPO baseline advances significantly faster in terms of environment  
 663 steps, for a fair comparison, we ran all experiments for one week on the same machine and use the  
 664 wall-clock time as the x-axis. We use Intel(R) Xeon(R) CPU E5-2698 v4 @ 2.20GHz equipped with  
 665 one NVIDIA Tesla V100 32GB.

## 666 B.2 Time-Based Training Curves

667 We provide the training curves in Figure 4. For brevity, we exclude a few of the depths from the plots.  
 668 As seen, there is a clear benefit for SoftTreeMax over distributed PPO with the standard softmax  
 669 policy. In most games, PPO with the SoftTreeMax policy shows very high sample efficiency: it  
 670 achieves higher episodic reward although it observes much less episodes, for the same running time.

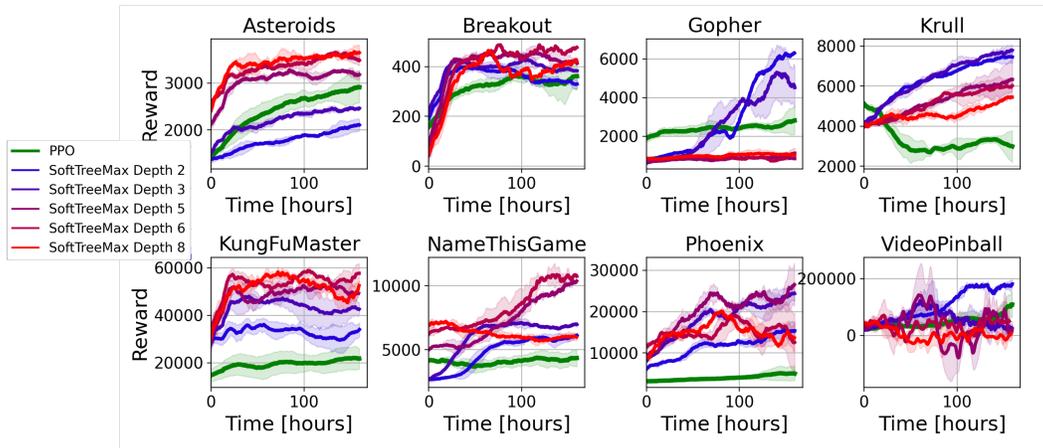


Figure 4: **Training curves: GPU SoftTreeMax (single worker) vs PPO (256 GPU workers).** The plots show average reward and standard deviation over 5 seeds. The x-axis is the wall-clock time. The runs ended after one week with varying number of time-steps. The training curves correspond to the evaluation runs in Figure 3.

## 671 B.3 Step-Based Training Curves

672 In Figure 5 we also provide the same convergence plots where the x-axis is now the number of online  
 673 interactions with the environment, thus excluding the tree expansion complexity. As seen, due to the  
 674 complexity of the tree expansion, less steps are conducted during training (limited to one week) as  
 675 the depth increases. In this plot, the monotone improvement of the reward with increasing tree depth  
 676 is noticeable in most games.

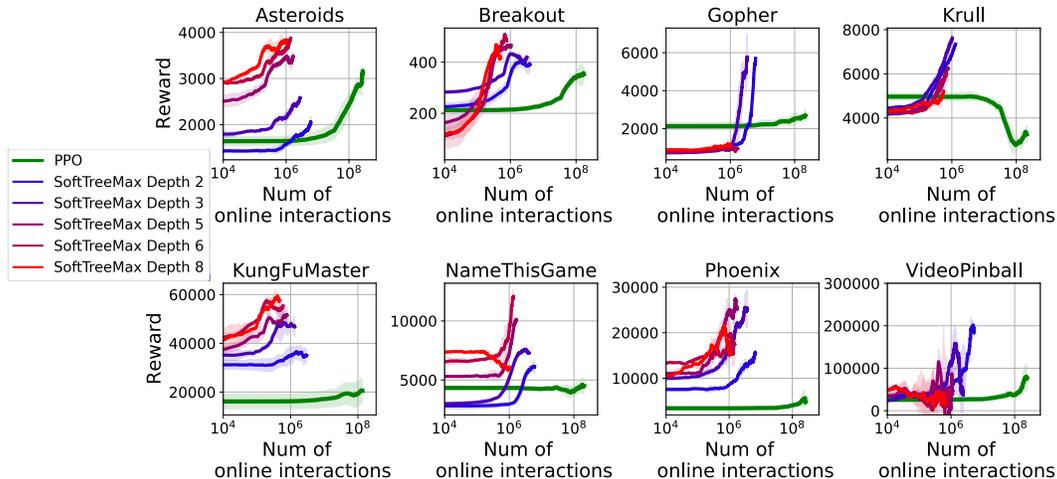


Figure 5: **Training curves: GPU SoftTreeMax (single worker) vs PPO (256 GPU workers)**. The plots show average reward and standard deviation over 5 seeds. The x-axis is the number of online interactions with the environment. The runs ended after one week with varying number of time-steps. The training curves correspond to the evaluation runs in Figure 3.

677 We note that not for all games we see monotonicity. Our explanation for this phenomenon relates to  
 678 how immediate reward contributes to performance compared to the value. Different games benefit  
 679 differently from long-term as opposed to short-term planning. Games that require longer-term  
 680 planning need a better value estimate. A good value estimate takes longer to obtain with larger depths,  
 681 in which we apply the network to states that are very different from the ones observed so far in the  
 682 buffer (recall that as in any deep RL algorithm, we train the model only on states in the buffer). If  
 683 the model hasn’t learned a good enough value function yet, and there is no guiding dense reward  
 684 along the trajectory, the policy becomes noisier, and can take more steps to converge – even more  
 685 than those we run in our week-long experiment.

686 For a concrete example, let us compare Breakout to Gopher. Inspecting Fig. 5, we observe that  
 687 Breakout quickly (and monotonically) gains from large depths since it relies on the short term goal  
 688 of simply keeping the paddle below the moving ball. In Gopher, however, for large depths ( $\geq 5$ ),  
 689 learning barely started even by the end of the training run. Presumably, this is because the task in  
 690 Gopher involves multiple considerations and steps: the agent needs to move to the right spot and  
 691 then hit the mallet the right amount of times, while balancing different locations. This task requires  
 692 long-term planning and thus depends more strongly on the accuracy of the value function estimate.  
 693 In that case, for depth 5 or more, we would require more train steps for the value to “kick in” and  
 694 become beneficial beyond the gain from the reward in the tree.

695 The figures above convey two key observations that occur for at least some non-zero depth: (1) The  
 696 final performance with the tree is better than PPO (Fig. 3); and (2) the intermediate step-based results  
 697 with the tree are better than PPO (Fig. 5). This leads to our main takeaway from this work — there  
 698 is no reason to believe that the vanilla policy gradient algorithm should be better than a multi-step  
 699 variant. Indeed, we show that this is not the case.

## 700 C Further discussion

### 701 C.1 The case of $\lambda_2(P^{\pi_b}) = 0$

702 When  $P^{\pi_b}$  is rank one, it is not only its variance that becomes 0, but also the norm of the gradient  
 703 itself (similarly to the case of  $d \rightarrow \infty$ ). Note that such a situation will happen rarely, in degenerate  
 704 MDPs. This is a local minimum for SoftTreeMax and it would cause the PG iteration to get stuck,  
 705 and to the optimum in the (desired but impractical) case where  $\pi_b$  is the optimal policy. However,  
 706 a similar phenomenon was also discovered in the standard softmax with deterministic policies:

707  $\theta(s, a) \rightarrow \infty$  for one  $a$  per  $s$ . PG with softmax would suffer very slow convergence near these  
708 local equilibria, as observed in Mei et al. [2020a]. To see this, note that the softmax gradient is  
709  $\nabla_{\theta} \log \pi_{\theta}(a|s) = e_a - \pi_{\theta}(\cdot|s)$ , where  $e_a \in [0, 1]^A$  is the vector with 0 everywhere except for the  
710  $a$ -th coordinate. I.e., it will be zero for a deterministic policy. SoftTreeMax avoids these local optima  
711 by integrating the reward into the policy itself (but may get stuck in another, as discussed above).

## 712 **NeurIPS Paper Checklist**

### 713 **1. Claims**

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

716 Answer: [Yes]

717 Justification: [NA]

718 Guidelines:

- 719 • The answer NA means that the abstract and introduction do not include the claims  
720 made in the paper.
- 721 • The abstract and/or introduction should clearly state the claims made, including the  
722 contributions made in the paper and important assumptions and limitations. A No or  
723 NA answer to this question will not be perceived well by the reviewers.
- 724 • The claims made should match theoretical and experimental results, and reflect how  
725 much the results can be expected to generalize to other settings.
- 726 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
727 are not attained by the paper.

### 728 **2. Limitations**

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

730 Answer: [Yes]

731 Justification: We included a relevant section at the end of the paper.

732 Guidelines:

- 733 • The answer NA means that the paper has no limitation while the answer No means that  
734 the paper has limitations, but those are not discussed in the paper.
- 735 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 736 • The paper should point out any strong assumptions and how robust the results are to  
737 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
738 model well-specification, asymptotic approximations only holding locally). The authors  
739 should reflect on how these assumptions might be violated in practice and what the  
740 implications would be.
- 741 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
742 only tested on a few datasets or with a few runs. In general, empirical results often  
743 depend on implicit assumptions, which should be articulated.
- 744 • The authors should reflect on the factors that influence the performance of the approach.  
745 For example, a facial recognition algorithm may perform poorly when image resolution  
746 is low or images are taken in low lighting. Or a speech-to-text system might not be  
747 used reliably to provide closed captions for online lectures because it fails to handle  
748 technical jargon.
- 749 • The authors should discuss the computational efficiency of the proposed algorithms  
750 and how they scale with dataset size.
- 751 • If applicable, the authors should discuss possible limitations of their approach to  
752 address problems of privacy and fairness.
- 753 • While the authors might fear that complete honesty about limitations might be used by  
754 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
755 limitations that aren't acknowledged in the paper. The authors should use their best  
756 judgment and recognize that individual actions in favor of transparency play an impor-  
757 tant role in developing norms that preserve the integrity of the community. Reviewers  
758 will be specifically instructed to not penalize honesty concerning limitations.

### 759 **3. Theory Assumptions and Proofs**

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

762 Answer: [Yes]

763 Justification: All proofs can be found in the appendix.

764 Guidelines:

- 765 • The answer NA means that the paper does not include theoretical results.
- 766 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 767 referenced.
- 768 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 769 • The proofs can either appear in the main paper or the supplemental material, but if
- 770 they appear in the supplemental material, the authors are encouraged to provide a short
- 771 proof sketch to provide intuition.
- 772 • Inversely, any informal proof provided in the core of the paper should be complemented
- 773 by formal proofs provided in appendix or supplemental material.
- 774 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 775 4. Experimental Result Reproducibility

776 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

777 perimental results of the paper to the extent that it affects the main claims and/or conclusions

778 of the paper (regardless of whether the code and data are provided or not)?

779 Answer: [Yes]

780 Justification: Yes. We also attached the repository, together with a docker environment, as

781 supplementary material.

782 Guidelines:

- 783 • The answer NA means that the paper does not include experiments.
- 784 • If the paper includes experiments, a No answer to this question will not be perceived
- 785 well by the reviewers: Making the paper reproducible is important, regardless of
- 786 whether the code and data are provided or not.
- 787 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 788 to make their results reproducible or verifiable.
- 789 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 790 For example, if the contribution is a novel architecture, describing the architecture fully
- 791 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 792 be necessary to either make it possible for others to replicate the model with the same
- 793 dataset, or provide access to the model. In general, releasing code and data is often
- 794 one good way to accomplish this, but reproducibility can also be provided via detailed
- 795 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 796 of a large language model), releasing of a model checkpoint, or other means that are
- 797 appropriate to the research performed.
- 798 • While NeurIPS does not require releasing code, the conference does require all submis-
- 799 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 800 nature of the contribution. For example
- 801 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 802 to reproduce that algorithm.
- 803 (b) If the contribution is primarily a new model architecture, the paper should describe
- 804 the architecture clearly and fully.
- 805 (c) If the contribution is a new model (e.g., a large language model), then there should
- 806 either be a way to access this model for reproducing the results or a way to reproduce
- 807 the model (e.g., with an open-source dataset or instructions for how to construct
- 808 the dataset).
- 809 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 810 authors are welcome to describe the particular way they provide for reproducibility.
- 811 In the case of closed-source models, it may be that access to the model is limited in
- 812 some way (e.g., to registered users), but it should be possible for other researchers
- 813 to have some path to reproducing or verifying the results.

#### 814 5. Open access to data and code

815 Question: Does the paper provide open access to the data and code, with sufficient instruc-

816 tions to faithfully reproduce the main experimental results, as described in supplemental

817 material?

818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868

Answer: [Yes]

Justification: We attached the repository, together with a docker environment, as supplementary material.

Guidelines:

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

## 6. Experimental Setting/Details

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

Answer: [Yes]

Justification: [NA]

Guidelines:

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

## 7. Experiment Statistical Significance

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

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 869 • It should be clear whether the error bar is the standard deviation or the standard error  
870 of the mean.
- 871 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
872 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
873 of Normality of errors is not verified.
- 874 • For asymmetric distributions, the authors should be careful not to show in tables or  
875 figures symmetric error bars that would yield results that are out of range (e.g. negative  
876 error rates).
- 877 • If error bars are reported in tables or plots, The authors should explain in the text how  
878 they were calculated and reference the corresponding figures or tables in the text.

## 879 8. Experiments Compute Resources

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

883 Answer: [Yes]

884 Justification: All relevant information can be found in the paper and the appendix.

885 Guidelines:

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

## 894 9. Code Of Ethics

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

897 Answer: [Yes]

898 Justification: [NA]

899 Guidelines:

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

## 905 10. Broader Impacts

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

908 Answer: [NA]

909 Justification: The paper has no societal impact.

910 Guidelines:

- 911 • The answer NA means that there is no societal impact of the work performed.
- 912 • If the authors answer NA or No, they should explain why their work has no societal  
913 impact or why the paper does not address societal impact.
- 914 • Examples of negative societal impacts include potential malicious or unintended uses  
915 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
916 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
917 groups), privacy considerations, and security considerations.

- 918 • The conference expects that many papers will be foundational research and not tied  
919 to particular applications, let alone deployments. However, if there is a direct path to  
920 any negative applications, the authors should point it out. For example, it is legitimate  
921 to point out that an improvement in the quality of generative models could be used to  
922 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
923 that a generic algorithm for optimizing neural networks could enable people to train  
924 models that generate Deepfakes faster.
- 925 • The authors should consider possible harms that could arise when the technology is  
926 being used as intended and functioning correctly, harms that could arise when the  
927 technology is being used as intended but gives incorrect results, and harms following  
928 from (intentional or unintentional) misuse of the technology.
- 929 • If there are negative societal impacts, the authors could also discuss possible mitigation  
930 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
931 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
932 feedback over time, improving the efficiency and accessibility of ML).

## 933 11. Safeguards

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

937 Answer: [NA]

938 Justification: The paper poses no such risks.

939 Guidelines:

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

## 950 12. Licenses for existing assets

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

954 Answer: [Yes]

955 Justification: [NA]

956 Guidelines:

- 957 • The answer NA means that the paper does not use existing assets.
- 958 • The authors should cite the original paper that produced the code package or dataset.
- 959 • The authors should state which version of the asset is used and, if possible, include a  
960 URL.
- 961 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 962 • For scraped data from a particular source (e.g., website), the copyright and terms of  
963 service of that source should be provided.
- 964 • If assets are released, the license, copyright information, and terms of use in the  
965 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
966 has curated licenses for some datasets. Their licensing guide can help determine the  
967 license of a dataset.
- 968 • For existing datasets that are re-packaged, both the original license and the license of  
969 the derived asset (if it has changed) should be provided.

970 • If this information is not available online, the authors are encouraged to reach out to  
971 the asset’s creators.

972 **13. New Assets**

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

975 Answer: [NA]

976 Justification: The paper does not release new assets.

977 Guidelines:

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

986 **14. Crowdsourcing and Research with Human Subjects**

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

990 Answer: [NA]

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

992 Guidelines:

- 993 • The answer NA means that the paper does not involve crowdsourcing nor research with  
994 human subjects.
- 995 • Including this information in the supplemental material is fine, but if the main contribu-  
996 tion of the paper involves human subjects, then as much detail as possible should be  
997 included in the main paper.
- 998 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
999 or other labor should be paid at least the minimum wage in the country of the data  
1000 collector.

1001 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**  
1002 **Subjects**

1003 Question: Does the paper describe potential risks incurred by study participants, whether  
1004 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1005 approvals (or an equivalent approval/review based on the requirements of your country or  
1006 institution) were obtained?

1007 Answer: [NA]

1008 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1009 Guidelines:

- 1010 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1011 human subjects.
- 1012 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1013 may be required for any human subjects research. If you obtained IRB approval, you  
1014 should clearly state this in the paper.
- 1015 • We recognize that the procedures for this may vary significantly between institutions  
1016 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1017 guidelines for their institution.
- 1018 • For initial submissions, do not include any information that would break anonymity (if  
1019 applicable), such as the institution conducting the review.