# EMBEDDING SAFETY INTO RL: A NEW TAKE ON TRUST REGION METHODS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Reinforcement Learning (RL) agents are capable of solving a wide variety of tasks, but are prone to produce unsafe behaviour. Constrained Markov Decision Processes (CMDPs) are a popular framework for incorporating safety constraints. However, common solution methods often compromise reward maximization by being overly conservative or by allowing unsafe behaviour during training. We propose Constrained Trust Region Policy Optimization (C-TRPO), a novel approach that modifies the geometry of the policy space based on the safety constraints, yielding trust regions composed exclusively of safe policies and ensuring constraint satisfaction throughout training. We theoretically study the convergence and update properties of C-TRPO and highlight connections to TRPO, Natural Policy Gradient (NPG), and Constrained Policy Optimization (CPO). We demonstrate experimentally that C-TRPO significantly reduces constraint violations while achieving competitive return compared to state-of-the-art algorithms.

023 024 025

026

004

010 011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

Reinforcement Learning (RL) has emerged as a highly successful paradigm in machine learning for solving sequential decision and control problems, with policy gradient (PG) algorithms as a popular approach (Williams, 1992; Sutton et al., 1999; Konda & Tsitsiklis, 1999). Policy gradients are especially appealing for high-dimensional continuous control because they can be easily extended to function approximation. Due to their flexibility and generality, there has been significant progress in enhancing PGs to work robustly with deep neural network-based approaches. Variants of natural policy gradient methods such as Trust Region Policy Optimization (TRPO) and Proximal Policy Optimization (PPO) are among the most widely used general-purpose reinforcement learning algorithms (Schulman et al., 2017a;b).

While flexibility makes PGs popular among practitioners, it comes at a cost: the policy can explore
 any behavior during training, posing significant risks in real-world applications. Many methods have
 been proposed to improve the safety of policy gradients, often based on the Constrained Markov
 Decision Process (CMDP) framework. However, existing methods either struggle to minimize con straint violations during training or severely limit the agent's performance.

This work introduces a simple strategy to enhance constraint satisfaction in trust-region-based safe
 policy gradient approaches without compromising performance. We propose a novel family of pol icy divergences, inspired by barrier function methods in optimization and safe control, that modify
 the policy geometry to ensure that the trust regions consist only of safe policies.

045 This approach is motivated by the observation that TRPO and related methods base their trust region 046 on the state-average Kullback-Leibler (KL) divergence. It can be derived as the Bregman divergence 047 induced by the negative conditional entropy on the space of state-action occupancies, as shown by 048 Neu et al. (2017). The main insight of the present work is that safer trust regions can be derived by altering this function to incorporate the cost constraints. The resulting divergence is skewed away from the constraint surface, which is achieved by augmenting the negative conditional entropy by 051 another convex barrier-like function. Manipulating the policy divergence in this way allows us to obtain a provably safe trust region-based policy optimization algorithm that retains most of TRPO's 052 mechanisms and guarantees, simplifying existing methods, while achieving competitive returns with less constraint violations throughout training.

Related work Classic solution methods for CMDPs rely on linear programming techniques, see
Altman (1999). However, they struggle to scale, making them unsuitable for high-dimensional
or continuous control problems. While there are numerous approaches to CMDPs, we focus on
model-free, direct policy optimization methods. Model-based approaches, like those popularized by
Berkenkamp et al. (2017), are attractive due to their strict safety guarantees, but require the learning
of a model, which is not always feasible.

060 Lagrangian methods are a widely adopted approach, where the optimization problem is reformu-061 lated as a weighted objective that balances rewards and penalties for constraint violations. This is 062 often motivated by Lagrangian duality, where the penalty coefficient is interpreted as the dual vari-063 able. Learning the coefficient with stochastic gradient descent presents a popular baseline (Achiam 064 et al., 2017; Ray et al., 2019; Chow et al., 2019; Stooke et al., 2020). However, a naively tuned Lagrange multiplier may not work well in practice due to oscillations and overshoot. To address 065 this issue, Stooke et al. (2020) apply PID control to tune the dual variable during training, which 066 achieves less oscillations around the constraint and faster convergence to a feasible policy. While 067 Lagrangian approaches are becoming increasingly popular, it is not entirely clear how to update the 068 dual variables during training, see Sohrabi et al. (2024). 069

Penalty methods such as IPO (Liu et al., 2020) and P3O (Zhang et al., 2022) propose weighted 071 penalty-based policy optimization objectives based on practical considerations. The penalties are weighted against the reward objective where the penalty coefficient is a hyper-parameter. This sim-072 plifies the Lagrangian approach since the penalty coefficients don't have to be optimized during 073 training, which results in improved stability. More recently, the approach to use (smoothed) log-074 barriers (Usmanova et al., 2024; Zhang et al., 2024a; Dey et al., 2024) became more popular, which 075 keeps the algorithm simple due to the penalty approach, but can offer certain constraint satisfac-076 tion guarantees, see e.g. Ni & Kamgarpour (2024). However, working with an explicit penalty 077 produces suboptimal policies w.r.t the original constrained MDP and thus introduces an additional error, which has to be controlled; see for example Geist et al. (2019); Müller & Cayci (2024) for 079 treatments of the regularization error in the unconstrained case, and Liu et al. (2020) for an example of an optimization gap in safe policy optimization. In contrast, changing the trust regions and 081 therefore the problem geometry does not change the objective function and the set of optimizers and therefore does not introduce an additional error. 082

083 Trust region methods are closely related to our approach, particularly Constrained Policy Optimiza-084 tion (CPO; Achiam et al. (2017)), which extends TRPO by restricting updates to the intersection of 085 the trust region and the safe policy set, which ensures safety during training. While CPO provides 086 constraint satisfaction guarantees, it tends to oscillate around the constraint boundary with high over-087 shoot as it only prevents the policy updates of TRPO from leaving the safe policy set. To address 880 constraint satisfaction, Projection-based CPO (PCPO; Yang et al. (2020)) projects updates onto the safe policy space between updates, improving stability but further hindering reward maximization. 089 Building on PCPO, Zhang et al. (2020) replace second-order updates with a computationally ef-090 ficient first-order approach, and Yang et al. (2022) further refine these methods with a different 091 projection approach, which achieves improved performance bounds by incorporating Generalized 092 Advantage Estimation (GAE; Schulman et al. (2018)). 093

094

**Rethinking safe trust region methods** We adopt a trust region approach that constructs trust 095 regions exclusively within the safe policy set, eliminating the need for projections or constrained 096 optimization in the inner loop. Trust region methods retain TRPO's update guarantees for both 097 reward and constraints but often underperform compared to barrier penalty methods. To address this, 098 we replace the state-average KL-divergence with policy divergences that act as barrier functions, 099 see Figure 1. This modification encourages updates of the resulting trust region method to move 100 more parallel to the constraint surfaces rather than directly toward and thereby improves constraint 101 satisfaction, simplifies optimization, and achieves competitive returns by maintaining policies within 102 the safe set for longer, see also Figure 6 in the Appendix.

103 104

**Contributions** We summarize our contributions as follows:

105 106 107

• In Section 3, we introduce a modified policy divergence such that every trust region consists of only safe policies. We introduce an idealized TRPO update based on the modified



Figure 1: Illustration of policy divergences (dashed) close to the constraint (red). a) TRPO (dotted for reference) and CPO. b) C-TRPO's divergence depends on the hyper-parameter  $\beta$ , which modulates the strength of the barrier towards the constraint surface. For  $\beta \searrow 0$  we obtain an update equivalent to CPO, and more conservative updates for larger values ( $\beta = 2$ ). The plots were generated with the toy MDP in Figure 2. c) Shown are the quadratic approximations of the divergence in parameter space, which is obtained by mapping the policy onto its occupancy measure, where a safe geometry can be defined using standard tools from convex optimization (safe region in white).

divergence, a tractable optimization algorithm for deep function approximation (C-TRPO), and a corresponding natural gradient method (C-NPG).

- We provide an efficient implementation of the proposed approximate C-TRPO method, see Section 3.2, which comes with a minimal overhead compared to TRPO (up to the estimation of the expected cost) and no overhead compared to CPO. We demonstrate experimentally that C-TRPO yields competitive returns with smaller constraint violations compared to common safe policy optimization algorithms, see Section 5.
  - In Section 4, we introduce C-TRPO's improvement guarantees and contrast to TRPO and CPO. Further, we show that the C-NPG method is the continuous time limit of C-TRPO and provides global convergence guarantees towards the optimal safe policy; this is in contrast to penalization or barrier methods, which introduce a bias

### 2 BACKGROUND

We consider the infinite-horizon discounted constrained Markov decision process (CMDP) and refer the reader to Altman (1999) for a general treatment. The CMDP is given by the tuple  $(S, A, P, r, \mu, \gamma, C)$ , where S and A are the finite state-space and action-space respectively and we refer to Appendix B.3 for a discussion of continuous state and action spaces. Further,  $P: S \times A \to \Delta_S$  is the transition kernel,  $r: S \times A \to \mathbb{R}$  is the reward function,  $\mu \in \Delta_S$  is the initial state distribution at time t = 0, and  $\gamma \in [0, 1)$  is the discount factor. The space  $\Delta_S$  is the set of categorical distributions over S. Further, define the constraint set  $C = \{(c_i, b_i)\}_{i=1}^m$ , where  $c_i: S \times A \to \mathbb{R}$  are the cost functions and  $b_i \in \mathbb{R}$  are the cost thresholds.

An agent interacts with the CMDP by selecting a policy  $\pi \in \Pi$  from the set of all Markov policies, i.e. an element from the Cartesian product of |S| probability simplicies on A. Given such a policy  $\pi$ , the value functions  $V_r^{\pi}, V_{c_i}^{\pi} : S \to \mathbb{R}$ , action-value functions  $Q_r^{\pi}, Q_{c_i}^{\pi} : S \times A \to \mathbb{R}$ , and advantage functions  $A_r^{\pi}, A_c^{\pi} : S \times A \to \mathbb{R}$  associated with the reward r and the *i*-th cost  $c_i$  are defined as

130 131

132

133

134

135

138

139

140

141

142 143

144 145

$$V_f^{\pi}(s) \coloneqq (1-\gamma) \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \middle| s_0 = s \right],$$

where the function f is either r or  $c_i$ , and the expectations are taken over trajectories of the Markov process, meaning with respect to the initial distribution  $s_0 \sim \mu$ , the policy  $a_t \sim \pi(\cdot|s_t)$  and the state transition  $s_{t+1} \sim P(\cdot | s_t, a_t)$ . Analogously, we set

163 164

166 167

168 169

178 179

185

188 189

190 191

197

199

200 201

202

206

207

208 209 210

215

162

$$Q_{f}^{\pi}(s,a) \coloneqq (1-\gamma) \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} f(s_{t},a_{t}) \Big| s_{0} = s, a_{0} = a \right] \text{ and } A_{f}^{\pi}(s,a) \coloneqq Q_{f}^{\pi}(s,a) - V_{f}^{\pi}(s).$$

The goal is to solve the following constrained optimization problem

$$\text{maximize}_{\pi \in \Pi} V_r^{\pi}(\mu) \quad \text{subject to} \quad V_{c_i}^{\pi}(\mu) \le b_i \quad \text{for all } i = 1, \dots, m, \tag{1}$$

170 where  $V_f^{\pi}(\mu)$  are the expected values under the initial state distribution  $V_f^{\pi}(\mu) := \mathbb{E}_{s \sim \mu}[V_f^{\pi}(s)]$ . 171 We will also write  $V_f^{\pi} = V_f^{\pi}(\mu)$ , and omit the explicit dependence on  $\mu$  for convenience, and we 172 write  $V_f(\pi)$  when we want to emphasize its dependence on  $\pi$ . We denote the set of safe policies by 173  $\Pi_{\text{safe}} = \bigcap_{i=1}^{m} \{\pi : V_{c_i}(\pi) \leq b_i\}$  and always assume that it is nontrivial.

175The Dual Linear Program for CMDPsAny stationary policy  $\pi$  induces a discounted state-action176(occupancy) measure  $d_{\pi} \in \Delta_{S \times A}$ , indicating the relative frequencies of visiting a state-action pair,177discounted by how far the event lies in the future. This probability measure is defined as

$$d_{\pi}(s,a) \coloneqq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\pi}(s_t = s) \pi(a|s),$$
(2)

where  $\mathbb{P}_{\pi}(s_t = s)$  is the probability of observing the environment in state s at time t given the agent follows policy  $\pi$ . For finite MDPs, it is well-known that maximizing the expected discounted return can be expressed as the linear program

maximize<sub>d</sub>  $r^{\top}d$  subject to  $d \in \mathscr{D}$ ,

where  $\mathscr{D}$  is the set of feasible state-action measures, which form a polytope (Kallenberg, 1994). Analogously to an MDP, the discounted cost CMDP can be expressed as the linear program

$$maximize_d \ r^{\top}d \quad subject \text{ to } d \in \mathscr{D}_{safe}, \tag{3}$$

where  $\mathscr{D}_{\text{safe}} = \bigcap_{i=1}^{m} \{ d : c_i^{\top} d \leq b_i \} \cap \mathscr{D} \text{ is the safe occupancy set, see Figure 4 in Appendix A.}$ 

Information Geometry of Policy Optimization Among the most successful policy optimization
 schemes are natural policy gradient (NPG) methods or variants thereof, such as trust-region and
 proximal policy optimization (TRPO and PPO, respectively). These methods assume a convex ge ometry and corresponding Bregman divergences in the state-action polytope, see Neu et al. (2017);
 Müller & Montúfar (2023) for more detailed discussions. A general trust region update is defined as

$$\pi_{k+1} \in \operatorname*{arg\,max}_{\pi \in \Pi} \mathbb{A}_r^{\pi_k}(\pi) \quad \text{sbj. to } D_{\Phi}(d_{\pi_k} || d_{\pi}) \le \delta, \tag{4}$$

where  $D_{\Phi} \colon \mathscr{D} \times \mathscr{D} \to \mathbb{R}$  is the Bregman divergence induced by a convex  $\Phi \colon \operatorname{int}(\mathscr{D}) \to \mathbb{R}$ , and

$$\mathbb{A}_r^{\pi_k}(\pi) = \mathbb{E}_{s,a \sim d_{\pi_k}} \left[ \frac{\pi(a|s)}{\pi_k(a|s)} A_r^{\pi_k}(s,a) \right],\tag{5}$$

is called the *policy advantage* or *surrogate advantage*. We can interpret A as a surrogate optimization objective for the expected return. In particular, for a parameterized policy  $\pi_{\theta}$ , it holds that  $\nabla_{\theta} \mathbb{A}_{r,\pi_{\theta_k}}(\pi_{\theta})|_{\theta=\theta_k} = \nabla_{\theta} V_r(\theta_k)$ , see Kakade & Langford (2002); Schulman et al. (2017a).

TRPO and the original NPG assume the same policy geometry (Kakade, 2001; Schulman et al., 2017a), since they employ an identical Bregman divergence

$$D_{\mathcal{K}}(d_{\pi_1}||d_{\pi_2}) \coloneqq \sum_{s,a} d_{\pi_1}(s,a) \log \frac{\pi_1(a|s)}{\pi_2(a|s)} = \sum_s d_{\pi_1}(s) D_{\mathcal{KL}}(\pi_1(\cdot|s)||\pi_2(\cdot|s)).$$

We refer to Appendix A for details on Bregman divergences. We call  $D_{\rm K}$  the *Kakade divergence* and informally write  $D_{\rm K}(\pi_1, \pi_2) := D_{\rm K}(d_{\pi_1}, d_{\pi_2})$ . This divergence can be shown to be the Bregman divergence induced by the negative conditional entropy

$$\Phi_{\mathcal{K}}(d_{\pi}) \coloneqq \sum_{s,a} d_{\pi}(s,a) \log \pi(a|s), \tag{6}$$

see Neu et al. (2017). It is well known that with a parameterized policy  $\pi_{\theta}$ , a linear approximation of A and a quadratic approximation of the Bregman divergence  $D_{\rm K}$  at  $\theta_k$ , one obtains the *natural policy gradient* step given by

$$\theta_{k+1} = \theta_k + \epsilon_k G_{\rm K}(\theta_k)^+ \nabla_\theta V_r(\pi_{\theta_k}),\tag{7}$$

where  $G_{\rm K}(\theta)^+$  denotes a pseudo-inverse of the generalized Fisher-information matrix of the policy with entries given by  $G_{\rm K}(\theta)_{ij} = \partial_{\theta_i} d_{\theta} \nabla^2 \Phi_{\rm K}(d_{\theta}) \partial_{\theta_j} d_{\theta}$ , see Schulman et al. (2017a); Müller & Montúfar (2023) and Appendix A for more detailed discussions.

### 3 A SAFE GEOMETRY FOR CONSTRAINED MDPS

To prevent the policy iterates from violating the constraints during optimization, we construct policy divergences for which the trust regions are contained in the safe policy set.

### 3.1 SAFE TRUST REGIONS

232 A Bregman divergence is induced by a mirror function that dictates the behavior of the divergence, 233 see A. Take for example the mirror function for TRPO and NPG in Equation (6). The divergence is 234 defined when both policies are in the interior of  $\mathcal{D}$ , and as either one of the policies approaches the 235 boundary of the state-action polytope, the divergence approaches infinity. Hence, TRPO and NPG 236 don't allow their policy iterates to become entirely deterministic during optimization. Since the 237 behavior of a Bregman divergences is dictated by the shape of its mirror function, we first construct 238 a family of *safe mirror functions*, that induce policy divergences that are finite only in the safe occupancy set  $\mathscr{D}_{safe}$  instead of the entire state-action polytope  $\mathscr{D}$ . Safe policy divergences, in turn, 239 let us derive safe trust region and natural policy gradient methods. 240

To this end, we consider mirror functions of the form

$$\Phi_{\mathcal{C}}(d) \coloneqq \Phi_{\mathcal{K}}(d) + \sum_{i=1}^{m} \beta_i \phi(b_i - c_i^{\top} d),$$
(8)

where  $\Phi_{\rm K}$  is the conditional entropy defined in Equation (6), and  $\phi: \mathbb{R}_{>0} \to \mathbb{R}$  is a convex function with  $\phi'(x) \to +\infty$  for  $x \searrow 0$ . This ensures that  $\Phi_{\rm C}: \operatorname{int}(\mathscr{D}_{\operatorname{safe}}) \to \mathbb{R}$  is strictly convex and has infinite curvature at the cost surface  $b_i - c_i^{\top} d = 0$ , which means  $\|\nabla \Phi_{\rm C}(d_k)\| \to +\infty$ , when  $b_i - c_i^{\top} d_k \searrow 0$ . Possible candidates for  $\phi$  are  $\phi(x) = -\log(x)$  and  $\phi(x) = x \log(x)$  corresponding to a logarithmic barrier and entropy, respectively.

251 The Bregman divergence induced by  $\Phi_{\rm C}$  is given by

$$D_{\rm C}(d_1||d_2) = D_{\rm K}(d_1||d_2) + \sum_{i=1}^m \beta_i D_{\phi_i}(d_1||d_2), \tag{9}$$

where

$$D_{\phi_i}(d_1||d_2) = \phi(b_i - V_{c_i}(\pi_1)) - \phi(b_i - V_{c_i}(\pi_2)) + \phi'(b_i - V_{c_i}(\pi_2))(V_{c_i}(\pi_1) - V_{c_i}(\pi_2)).$$
(10)

The corresponding trust-region scheme is given by

$$\pi_{k+1} \in \operatorname*{arg\,max}_{\pi \in \Pi} \mathbb{A}_r^{\pi_k}(\pi) \quad \text{sbj. to } D_{\mathcal{C}}(d_{\pi_k} || d_{\pi}) \le \delta, \tag{11}$$

where  $\mathbb{A}_r$  is defined in Equation (5). Note the constraint is only satisfied if  $d_1, d_2 \in \operatorname{int}(\mathscr{D}_{safe})$  and the divergence approaches  $+\infty$  as  $d_2$  approaches the boundary of the safe set. Thus, the trust region  $\{d \in \mathscr{D} : D_{\mathbb{C}}(d_k || d) \leq \delta\}$  is contained in the set of safe occupancy measures for any finite  $\delta$ . Analogously to the case of unconstrained TRPO the corresponding natural policy gradient scheme is given by

$$\theta_{k+1} = \theta_k + \epsilon_k G_{\rm C}(\theta_k)^+ \nabla V_r(\theta_k), \tag{12}$$

where  $G_{\rm C}(\theta)^+$  denotes an arbitrary pseudo-inverse of  $G_{\rm C}(\theta)_{ij} = \partial_{\theta_i} d_{\theta}^\top \nabla^2 \Phi_{\rm C}(d_{\theta}) \partial_{\theta_j} d_{\theta}$ .

5

252 253 254

255

260 261 262

268

269

243 244 245

220 221

222

223

224 225

226 227

228

229 230

#### 270 3.2 CONSTRAINED TRUST REGION POLICY OPTIMIZATION 271

272 If we could solve the optimization problem in Equation (11) exactly, we would obtain a provably 273 safe trust region policy optimization method with zero constraint violations, as long as we start with a safe policy. However, the exact trust region update Equation (11) cannot be computed. Firstly, the 274 divergence depends on expected cost values, which we can only estimate. The resulting estimation 275 errors of the divergence might cause the policy iterates to leave the safe set, in which case the 276 divergence becomes ill-defined. Finally, the divergence also depends on the expected cost value of 277 the proposal policy, which is not available during the updates. To address these issues, we propose 278 an update based on a surrogate divergence, similar to how surrogate objectives are used in policy 279 optimization. We propose the following update, which we call Constrained TRPO (C-TRPO). 280

283

284

285

286

287 288

289

290

291

292

293

295

296

297

298

299 300

301

307 308

314 315 316

317 318 319

$$\pi_{k+1} = \underset{\pi \in \Pi}{\arg\max} \mathbb{A}_r^{\pi_k}(\pi) \quad \text{sbj. to } \bar{D}_{\mathcal{C}}(\pi || \pi_k) \le \delta.$$
(13)

Here,  $\bar{D}_{\rm C}$  is the surrogate divergence, defined below. Algorithm 1 shows the implementation of C-TRPO, which performs a constrained trust region update if the current policy is safe or a recovery step that minimizes the cost if the policy is unsafe. For the trust region update, we follow a similar implementation to the original TRPO, estimating the divergence, using a linear approximation of the surrogate objective, and a quadratic approximation of the trust region.

#### Algorithm 1 Constrained TRPO (C-TRPO); differences from TRPO in blue

1: Input: Initial policy  $\pi_0 \in \Pi_{\theta}$ , safety parameter  $\beta > 0$ , recovery parameter  $0 < b_{\rm H} \leq b$ 2: for  $k = 0, 1, 2, \dots$  do Sample a set of trajectories following  $\pi_k = \pi_{\theta_k}$ 3: if  $\pi_k \in \Pi_{\text{safe}}^{\text{H}}$  then 4:  $A \leftarrow A_r; D \leftarrow \bar{D}_C = \bar{D}_{KL} + \beta \bar{D}_{\Phi}$  {Constrained trust region update} 5: 6: else  $A \leftarrow -A_c; D \leftarrow \bar{D}_{\mathrm{KL}} \{ \text{Recovery} \}$ 7: end if 8: Compute  $\pi_{k+1}$  using TRPO with A as advantage estimate and with D as policy divergence. 9: 10: end for

**Surrogate Divergence** To aid in clarity, we focus on the case with a single constraint, but the results are easily extended to multiple constraints by summation of the individual constraint terms. 302 In practice, the exact constrained KL-Divergence  $D_{\rm C}$  cannot be evaluated, because it depends on 303 the cost-return of the optimized policy  $V_c(\pi)$ . However, we can approximate it locally around the 304 policy of the k-th iteration,  $\pi_k$ , using a surrogate divergence. This surrogate can be expressed as a 305 function of the policy cost advantage 306

$$\mathbb{A}_{c}^{\pi_{k}}(\pi) = \mathbb{E}_{d_{\pi_{k}}}\left[\frac{\pi(a|s)}{\pi_{k}(a|s)}A_{c}^{\pi_{k}}(s,a)\right],\tag{14}$$

309 which approximates  $V_c(\pi) - V_c^{\pi_k}$  up to first order in the policy parameters (Kakade & Langford, 310 2002; Schulman et al., 2017a; Achiam et al., 2017). Assume  $\pi_k \in \Pi_{\text{safe}}$  and define the *constraint* 311 margin  $\delta_b = b - V_c^{\pi_k}$ , which is positive if  $\pi_k \in \Pi_{\text{SAFE}}$ . Further, define the surrogate divergence 312  $\bar{D}_{\mathrm{C}}(\pi || \pi_k) = \bar{D}_{\mathrm{KL}}(\pi || \pi_k) + \beta \bar{D}_{\phi}(\pi || \pi_k)$ , where 313

$$\bar{D}_{\mathrm{KL}}(\pi || \pi_k) = \sum_{s \in \mathcal{S}} d_{\pi_k}(s) D_{\mathrm{KL}}(\pi || \pi_k)$$
(15)

and

$$\bar{D}_{\phi}(\pi_{\theta}||\pi_{\theta_{k}}) = \begin{cases} \phi(\delta_{b} - \mathbb{A}_{c}^{\pi_{k}}(\pi)) - \phi(\delta_{b}) + \phi'(\delta_{b})\mathbb{A}_{c}^{\pi_{k}}(\pi), & \text{if } \delta_{b} - \mathbb{A}_{c}^{\pi_{k}} \in \operatorname{dom}(\phi) \\ \infty & \text{otherwise} . \end{cases}$$
(16)

320 The surrogate  $D_{\phi}$  is closely related to the Bregman divergence  $D_{\phi}$ . They are equivalent up to the 321 substitution  $V_c(\pi) - V_c(\pi_k) \to \mathbb{A}_c^{\pi_k}(\pi)$ , see Appendix B.1. The surrogate can be estimated from samples of the MDP. In the practical implementation, we estimate  $\delta_b$ , and the policy cost advantage 322 from trajectory samples using GAE- $\lambda$  estimates Schulman et al. (2018). The consequences of the 323 substitution in the surrogate will be discussed in Section 4.

**Recovery with Hysteresis** The iterate may still leave the safe policy set  $\Pi_{safe}$ , either due to approximation errors of the divergence, or because we started outside the safe set. In this case, we perform a recovery step, where we only minimize the cost with TRPO as by Achiam et al. (2017). In tasks where the policy starts in the unsafe set, C-TRPO can get stuck at the cost surface. This is easily mitigated by including a hysteresis condition for returning to the safe set. If  $\pi_{k-1}$  is the previous policy, then  $\pi_k \in \Pi_{safe}^{\rm H}$  with  $\Pi_{safe}^{\rm H} = \{\pi_{\theta} \in \Pi_{\theta} \text{ and } V_c(\pi_{\theta}) \leq b_{\rm H}\}$  where  $b_{\rm H} = b$  if  $\pi_{k-1} \in \Pi_{safe}^{\rm H}$  and a user-specified fraction of b otherwise.

**Computational Complexity** The C-TRPO implementation adds no computational overhead compared to CPO, since  $D_{\phi}$  is just a function of the cost advantage estimate, and is simply added to the divergence of TRPO. Compared to TRPO, the cost value function must be approximated.

### 3.3 CONSTRAINED NATURAL POLICY GRADIENT

Practically, the C-TRPO optimization problem in Equation (13) is solved like traditional TRPO: the objective is approximated linearly, and the constraint is approximated quadratically in the policy parameters using automatic differentiation and the conjugate gradients algorithm. This leads to the policy parameter update

$$\theta_{k+1} = \theta_k + \alpha^i \sqrt{\frac{2\delta}{g_k^\top H_k^{-1} g_k}} \cdot H_k^{-1} g_k, \tag{17}$$

where

$$g_k = \nabla_{\theta} \mathbb{A}_c^{\theta_k}(\pi_{\theta})|_{\theta = \theta_k} \quad \text{and } H_k = \bar{H}_{\mathcal{C}}(\theta_k) = \nabla_{\theta}^2 \bar{D}_{\mathcal{C}}(\pi_{\theta} || \pi_{\theta_k})|_{\theta = \theta_k}$$
(18)

are finite sample estimates, and  $H^{-1}g$  is approximated using conjugate gradients. The  $\alpha^i \in [0, 1]$  are the coefficients for backtracking line search, which ensures  $\bar{D}_{\rm C}(\pi_{\theta}||\pi_{\theta_k}) \leq \delta$ .

We show in Appendix B.2.3 that the Hessian

$$\nabla_{\theta}^{2} \bar{D}_{\mathcal{C}}(\theta || \hat{\theta})|_{\theta = \hat{\theta}} = G_{\mathcal{K}}(\theta_{k}) + \beta \phi''(b - V_{c}^{\hat{\theta}}(\theta)) \nabla_{\theta} V_{c}^{\hat{\theta}}(\theta)^{\top} \nabla_{\theta} V_{c}^{\hat{\theta}}(\theta),$$

is equivalent to the Gramian  $G_{\rm C}(\theta_k)$  of the natural gradient update in Equation (19). We call the resulting policy gradient

$$\theta_{k+1} = \theta_k + \epsilon_k \bar{H}_{\rm C}(\theta_k)^+ \nabla V_r(\theta_k), \tag{19}$$

the *Constrained NPG* (C-NPG). In particular, this shows that the C-TRPO update can be interpreted as a natural policy gradient step with an adaptive step size, see Appendix A. We emphasize that the idealized safe trust region update in Equation (11) and the C-TRPO update of Equation (13) agree up to second order in the policy parameters. This justifies the surrogate divergence in C-TRPO and motivates the discussion of the C-NPG flow in Section 4.2. We show in Theorem 5 that  $int(\mathcal{D}_{safe})$ is invariant under the dynamics of the C-NPG. This implies that if the trust region radius  $\delta$  is small, and the advantage estimation is accurate enough, the iterates under C-TRPO never leave the safe set.

# 4 ANALYSIS

Here, we provide a theoretical analysis of the updates of C-TRPO and study the convergence properties of the time-continuous version of C-NPG. All proofs are deferred to the Appendix C.

### 4.1 PROPERTIES OF THE C-TRPO UPDATE

The practical C-TRPO algorithm is implemented using the surrogate divergence introduced in Equation (13), which is identical to the theoretical divergence  $D_{\rm C}$  introduced in Equation (11) up to a mismatch between the policy advantage and the performance difference. The motivation for substituting the policy cost advantage for the performance difference is their equivalence up to first order and that we can estimate the advantage from samples of  $d_{\pi}$ . Similar to CPO, we can guarantee an almost improvement of the return (Achiam et al., 2017), despite the new divergence. 381 382

387

393

394

395

396 397 398

399

404

405

410

**Proposition 1** (C-TRPO reward update). Set  $\epsilon_r = \max_s |\mathbb{E}_{a \sim \pi_{k+1}} A_r^{\pi_k}(s, a)|$ . The expected reward of a policy updated with C-TRPO is bounded from below by

$$V_r(\pi_{k+1}) \ge V_r(\pi_k) - \frac{\sqrt{2\delta\gamma\epsilon_r}}{1-\gamma}.$$
(20)

Constraint violation, however, behaves slightly differently for the two algorithms. To see this, we establish a more concrete relation between C-TRPO and CPO. As  $\beta \searrow 0$ , the solution to Equation (13) approaches the constraint surface in the worst case, and we recover CPO, see Figure 1.

**Proposition 2.** The approximate C-TRPO update approaches the CPO update in the limit as  $\beta \searrow 0$ .

Intuitively, solving the C-TRPO problem with successively smaller values of  $\beta$ , would be similar to solving CPO with the interior point method using  $\bar{D}_{\phi}(\cdot || \pi_k)$  as the barrier function.

Further, C-TRPO is more conservative than CPO for any  $\beta > 0$  and as  $\beta \to +\infty$  the updated is maximally constrained in the cost-increasing direction. This is formalized as follows.

**Proposition 3** (C-TRPO worst-case constraint violation). Consider  $\Psi: [0, \delta_b) \to [0, \infty)$  defined by  $\Psi(x) = \phi(\delta_b - x) - \phi(\delta_b) - \phi'(\delta_b) \cdot x$  such that  $D_{\phi}(\pi || \pi_k) = \Psi(\mathbb{A}_c^{\pi_k}(\pi))$ . Further, set  $\epsilon_c = \max_s |\mathbb{E}_{a \sim \pi_{k+1}} A_c^{\pi_k}(s, a)|$ , and choose a strictly convex  $\phi$ . The worst-case constraint violation for C-TRPO is

$$V_c(\pi_{k+1}) \le V_c(\pi_k) + \Psi^{-1}(\delta/\beta) + \frac{\sqrt{2\delta\gamma\epsilon_c}}{1-\gamma}.$$
(21)

Further, it holds that  $\lim_{\beta \to +\infty} \Psi^{-1}(\delta/\beta) = 0$  and  $\Psi^{-1}(\delta/\beta) < b - V_c(\pi_k)$  for all  $\beta \in (0, \infty)$ .

This result is analogous to the worst-case constraint violation for CPO (Achiam et al., 2017, Proposition 2), except that it depends on the choice of  $\beta$  and is tighter than the corresponding guarantee for CPO, because  $\Psi^{-1}(\delta/\beta) < b - V_c(\pi_k)$  for all  $\beta \in (0, \infty)$ .

### 4.2 INVARIANCE AND CONVERGENCE OF CONSTRAINED NATURAL POLICY GRADIENTS

It is well known that TRPO is equivalent to a natural policy gradient method with an adaptive step
size, see also Appendix A. We study the time-continuous limit of C-TRPO and guarantee safety
during training and global convergence. In the context of constrained TRPO in Equation (11), we
study the natural policy gradient flow

$$\partial_t \theta_t = G_{\rm C}(\theta_t)^+ \nabla V_r(\theta_t), \tag{22}$$

where  $G_{\rm C}(\theta)^+$  denotes a pseudo-inverse of  $G_{\rm C}(\theta)_{ij} = \partial_{\theta_i} d_{\theta}^\top \nabla^2 \Phi_{\rm C}(d_{\theta}) \partial_{\theta_j} d_{\theta}$  and  $\theta \mapsto \pi_{\theta}$  is a differentiable policy parametrization. Moreover, we assume that  $\theta \mapsto \pi_{\theta}$  is regular, that it is surjective and the Jacobian is of maximal rank everywhere. This assumption implies overparametrization but is satisfied for common models like tabular softmax, tabular escort, or expressive log-linear policy parameterizations (Agarwal et al., 2021a; Mei et al., 2020a; Müller & Montúfar, 2023).

**Theorem 4** (Safety during training). Assume that  $\phi : \mathbb{R}_{>0} \to \mathbb{R}$  satisfies  $\phi'(x) \to +\infty$  for  $x \searrow 0$ and consider a regular policy parameterization. Then the set  $\Theta_{\mathcal{C}}$  is invariant under Equation (22).

A visualization of policies obtained by C-NPG for different safe initializations and varying choices of  $\beta$  is shown in Figure 2 for a toy MDP. We see that for even small choices of  $\beta$  the trajectories don't cross the constraint surface and the updates become more conservative for larger choices of  $\beta$ . Theorem 5. Assume that  $\frac{1}{2}(\alpha) = 1$ ,  $\alpha = 1$ ,  $\beta = 0$ , act  $V_{1,2}^{*}$ ,  $\gamma = 0$ ,  $\gamma = 0$ ,  $\gamma = 0$ .

**Theorem 5.** Assume that  $\phi'(x) \to +\infty$  for  $x \searrow 0$ , set  $V_{r,C}^* \coloneqq \max_{\pi \in \Pi_{\text{safe}}} V_r(\pi)$  and denote the set of optimal constrained policies by  $\Pi_{\text{safe}}^* = \{\pi \in \Pi_{\text{safe}} : V_r(\pi) = V_{r,C}^*\}$ , consider a regular policy parametrization and let  $(\theta_t)_{t\ge 0}$  solve Equation (22). It holds that  $V_r(\pi_{\theta_t}) \to V_{r,C}^*$  and **431** 

$$\lim_{t \to +\infty} \pi_t = \pi_{\text{safe}}^* = \arg\min\{D_{\mathcal{C}}(\pi^*, \pi_0) : \pi^* \in \Pi_{\text{safe}}^*\}.$$
(23)



Figure 2: Shown is the policy set  $\Pi \cong [0,1]^2$  for an MDP with two states and two actions with a heatmap of the reward  $V_r$ ; the constraint surface is shown in black with the safe policies below; optimization trajectories are shown for 10 safe initialization and for  $\beta = 0, 10^{-4}, 10^{-2}, 1$ .

In case of multiple optimal policies, Equation (23) identifies the optimal policy of the CMDP that the natural policy gradient method converges to as the projection of the initial policy  $\pi_0$  to the set of optimal safe policies  $\Pi_{safe}^*$  with respect to the constrained divergence  $D_C$ . In particular, this implies that the limiting policy  $\pi_{safe}^*$  satisfies as few constraints with equality as required to be optimal. To see this, note that  $\Pi_{safe}^*$  forms a face of  $\mathscr{D}_{safe}$  and that Bregman projections lie at the interior of faces (Müller et al., 2024, Lemma A.2) and hence satisfy as few linear constraints as required.

# 5 COMPUTATIONAL EXPERIMENTS

456 **Setup and main results** We benchmark C-TRPO against 9 common safe policy optimization al-457 gorithms (CPO Achiam et al. (2017), PCPO Yang et al. (2020), CPPO-PID Stooke et al. (2020), PPO-Lag and TRPO-Lag Achiam et al. (2017); Ray et al. (2019), FOCOPS Zhang et al. (2020), 458 CUP Yang et al. (2022), IPO Liu et al. (2020) and P3O Zhang et al. (2022)) on 8 tasks (4 Naviga-459 tion and 4 Locomotion) from the Safety Gymnasium (Ji et al., 2023) benchmark. The locomotion 460 tasks reward distance traveled, while penalizing high velocities, and the navigation tasks reward 461 goal reaching and penalize certain unsafe states. For the C-TRPO implementation we fix the con-462 vex generator  $\phi(x) = x \log(x)$ , motivated by its superior performance in our experiments, see 463 Appendix B.2.1, and  $b_{\rm H} = 0.8b$  and  $\beta = 1$  across all experiments.<sup>1</sup> We train each algorithm for 464 10 million environment steps and evaluate on 10 runs after training, see Table 1 in Appendix D. 465 Furthermore, each algorithm is trained with 5 seeds, and the cost regret is monitored throughout 466 training for every run. To get a better sense of the safety of the algorithms during training, we take 467 an online learning perspective and include as a metric the cumulative cost violation (Efroni et al., 2020; Müller et al., 2024) 468

472

442 443

444 445 446

447

448

449

450 451

452 453

454 455

$$CUMCOST_{+}(K,c) \coloneqq \sum_{k=0}^{K-1} [V_{c}^{\pi_{k}} - b]_{+}, \qquad (24)$$

473 where  $[x]_{+} = \max\{0, x\}$ , and K is the number of the training iterations.

We observe that C-TRPO is competitive with the leading algorithms of the benchmark in terms of
expected return (CPO, TRPO-Lagrangian), see Figure 3. Furthermore, it achieves notably lower cost
regret throughout training than the high-return algorithms, even comparable to the more conservative
PCPO algorithm. In Figure 3, we visualize the interquartile mean (IQM) of normalized scores across
training for expected returns of reward and cost and for the cost regret, including their stratified
bootstrap confidence intervals (Agarwal et al., 2021b).

481 **Discussion** For completeness, we also report environment-wise sample efficiency curves and eval-482 uation performances in Appendix D.3. Our experiments reveal that the algorithm's performance is 483 closely tied to the accuracy of divergence estimation, which hinges on the precise estimation of the 484 cost advantage and value functions. The safety parameter  $\beta$  modulates the stringency with which

<sup>485</sup> 

<sup>&</sup>lt;sup>1</sup>Code available at: (will be released after double-blind review)



Figure 3: Comparison of safe policy optimization algorithms based on the Inter Quartile Mean across 5 seeds and 8 tasks. From left to right: episode return of the reward (PPO normalized), episode return of the cost (threshold normalized), and cumulative cost violation (CPO normalized).

C-TRPO satisfies the constraint, and can do so without limiting the expected return on most environments at least for  $\beta \le 1$ , see Figure 7. For higher values, the expected return starts to degrade, partly due to  $\bar{D}_{\phi}$  being relatively noisy compared to  $\bar{D}_{\text{KL}}$  and thus we recommend the choice  $\beta = 1$ .

Further, we observe that constraint satisfaction is stable across different choices of cost threshold b,
 see Figure 8, and that in most environments, constraint violations seem to reduce as the algorithm
 converges, meaning that the regret flattens over time. This behavior suggests that the divergence estimation becomes increasingly accurate over time, potentially allowing C-TRPO to achieve sublinear
 regret. However, we leave regret analysis of the finite sample regime for future research.

We attribute the improved constraint satisfaction compared to CPO to a slowdown and reduction 510 in the frequency of oscillations around the cost threshold, which mitigates overshoot behaviors that 511 could otherwise violate constraints. The modified gradient preconditioner appears to deflect the pa-512 rameter trajectory away from the constraint, see Figure 2. This effect may also be partially attributed 513 to the hysteresis-based recovery mechanism, which helps smooth updates by leading the iterate away 514 from the boundary of the safe set. Employing a hysteresis fraction  $0 < b_{\rm H} < b$  might also be ben-515 eficial because C-TRPO's divergence estimates tend to be more reliable for strictly safe policies. 516 The effect of the choice of  $b_{\rm H}$  is shown in Figure 10 in the appendix. Finally, we present ablations 517 in Appendix D.2, which support our claims that both components—the modified trust region and 518 hysteresis—are effective in reducing safety violations.

519 520 521

522

497

498

499 500 501

### 6 CONCLUSION AND OUTLOOK

In this paper, we introduced C-TRPO and C-NPG, two novel methods for solving Constrained 523 Markov Decision Processes (CMDPs). C-TRPO can be viewed as an extension or relaxation of 524 Constrained Policy Optimization (CPO), from which a natural policy gradient method, C-NPG, is 525 derived. C-TRPO represents a significant step toward safe, model-free reinforcement learning by 526 integrating constraint handling directly into the geometry of the policy space. Meanwhile, C-NPG 527 provides a provably safe natural policy gradient method for CMDPs, offering a foundational ap-528 proach to direct policy optimization in constrained settings-similar to how NPG is a cornerstone in 529 the theory of policy gradients for unconstrained MDPs. However, there are several limitations to ad-530 dress. First, the divergence estimation remains challenging, and we did not investigate the properties 531 of the finite sample estimates of the divergence. In addition, the CMDP framework may be somewhat limited in modeling safe exploration and control. Because CMDPs constrain the average cost 532 return, it can be difficult to model trajectory-wise or state-wise safety constraints. Several promising 533 directions for future research remain open. One avenue is to combine these methods with model-534 based policy optimization to improve cost return estimates, or with policy mirror descent to improve 535 computational efficiency, see e.g. Tomar et al. (2022). Additionally, integrating the proposed di-536 vergence with other safe policy optimization algorithms that utilize trust regions, e.g. PCPO, could 537 lead to stronger performance guarantees. 538

539 Overall, the proposed algorithms, C-TRPO and C-NPG, present a step forward in general-purpose CMDP algorithms and move us closer to deploying RL in high-stakes, real-world applications.

# 540 REFERENCES

555

571

572

573

574

578

579

580 581

582

- Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization, 2017.
- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *The Journal of Machine Learning Research*, 22(1):4431–4506, 2021a.
- 547 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare.
   548 Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Informa-*549 *tion Processing Systems*, 34, 2021b.
- Eitan Altman. Constrained Markov Decision Processes. CRC Press, Taylor & Francis Group, 1999.
   URL https://api.semanticscholar.org/CorpusID:14906227.
- Felipe Alvarez, Jérôme Bolte, and Olivier Brahic. Hessian riemannian gradient flows in convex
   programming. *SIAM journal on control and optimization*, 43(2):477–501, 2004.
- 556 Shun-ichi Amari. Information geometry and its applications, volume 194. Springer, 2016.
- Aaron D. Ames, Xiangru Xu, Jessy W. Grizzle, and Paulo Tabuada. Control barrier function based quadratic programs for safety critical systems. *IEEE Transactions on Automatic Control*, 62(8): 3861–3876, August 2017. ISSN 1558-2523. doi: 10.1109/tac.2016.2638961. URL http: //dx.doi.org/10.1109/TAC.2016.2638961.
- Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. Safe modelbased reinforcement learning with stability guarantees. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper\_files/paper/2017/ file/766ebcd59621e305170616ba3d3dac32-Paper.pdf.
- 568 Richard Cheng, Gabor Orosz, Richard M. Murray, and Joel W. Burdick. End-to-end safe reinforce 569 ment learning through barrier functions for safety-critical continuous control tasks, 2019. URL
   570 https://arxiv.org/abs/1903.08792.
  - Yinlam Chow, Ofir Nachum, Aleksandra Faust, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. Lyapunov-based safe policy optimization for continuous control, 2019. URL https://arxiv.org/abs/1901.10031.
- Sumanta Dey, Pallab Dasgupta, and Soumyajit Dey. P2bpo: Permeable penalty barrier-based policy optimization for safe rl. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
  volume 38, pp. 21029–21036, 2024.
  - Yonathan Efroni, Shie Mannor, and Matteo Pirotta. Exploration-exploitation in constrained mdps, 2020. URL https://arxiv.org/abs/2003.02189.
  - Eugene A Feinberg and Adam Shwartz. *Handbook of Markov decision processes: methods and applications*, volume 40. Springer Science & Business Media, 2012.
- Matthieu Geist, Bruno Scherrer, and Olivier Pietquin. A theory of regularized markov decision
   processes. In *International Conference on Machine Learning*, pp. 2160–2169. PMLR, 2019.
- Jiaming Ji, Borong Zhang, Jiayi Zhou, Xuehai Pan, Weidong Huang, Ruiyang Sun, Yiran Geng, Yifan Zhong, Josef Dai, and Yaodong Yang. Safety gymnasium: A unified safe reinforcement learning benchmark. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL https://openreview.net/forum?id= WZmlxIuIGR.
- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In
   *Proceedings of the Nineteenth International Conference on Machine Learning*, ICML '02, pp. 267–274, San Francisco, CA, USA, 2002. Morgan Kaufmann Publishers Inc. ISBN 1558608737.

594 Sham M Kakade. A natural policy gradient. Advances in neural information processing systems, 595 14, 2001. 596 Lodewijk CM Kallenberg. Survey of linear programming for standard and nonstandard markovian 597 control problems. part i: Theory. Zeitschrift für Operations Research, 40:1-42, 1994. 598 Vijay Konda and John Tsitsiklis. Actor-critic algorithms. In S. Solla, T. Leen, and 600 K. Müller (eds.), Advances in Neural Information Processing Systems, volume 12. MIT Press, 601 1999. URL https://proceedings.neurips.cc/paper\_files/paper/1999/ 602 file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf. 603 Yongshuai Liu, Jiaxin Ding, and Xin Liu. Ipo: Interior-point policy optimization under constraints. 604 Proceedings of the AAAI Conference on Artificial Intelligence, 34:4940–4947, 04 2020. doi: 605 10.1609/aaai.v34i04.5932. 606 607 Jincheng Mei, Chenjun Xiao, Bo Dai, Lihong Li, Csaba Szepesvári, and Dale Schuurmans. Escaping 608 the gravitational pull of softmax. Advances in Neural Information Processing Systems, 33:21130-21140, 2020a. 609 610 Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the global convergence 611 rates of softmax policy gradient methods. In International Conference on Machine Learning, pp. 612 6820-6829. PMLR, 2020b. 613 Johannes Müller and Semih Cayci. Essentially sharp estimates on the entropy regularization error 614 in discrete discounted markov decision processes. arXiv preprint arXiv:2406.04163, 2024. 615 616 Johannes Müller and Guido Montúfar. Geometry and convergence of natural policy gradient meth-617 ods. Information Geometry, pp. 1-39, 2023. 618 Johannes Müller, Semih Çaycı, and Guido Montúfar. Fisher-rao gradient flows of linear programs 619 and state-action natural policy gradients. arXiv preprint arXiv:2403.19448, 2024. 620 621 Adrian Müller, Pragnya Alatur, Volkan Cevher, Giorgia Ramponi, and Niao He. Truly no-regret 622 learning in constrained mdps, 2024. URL https://arxiv.org/abs/2402.15776. 623 Gergely Neu, Anders Jonsson, and Vicenç Gómez. A unified view of entropy-regularized markov 624 decision processes. arXiv preprint arXiv:1705.07798, 2017. 625 626 Tingting Ni and Maryam Kamgarpour. A safe exploration approach to constrained markov deci-627 sion processes. In ICML 2024 Workshop: Foundations of Reinforcement Learning and Control-628 Connections and Perspectives, 2024. 629 Matteo Pirotta, Marcello Restelli, Alessio Pecorino, and Daniele Calandriello. Safe policy itera-630 tion. In Sanjoy Dasgupta and David McAllester (eds.), Proceedings of the 30th International 631 Conference on Machine Learning, volume 28 of Proceedings of Machine Learning Research, pp. 632 307-315, Atlanta, Georgia, USA, 17-19 Jun 2013. PMLR. URL https://proceedings. 633 mlr.press/v28/pirotta13.html. 634 Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement 635 learning. arXiv preprint arXiv:1910.01708, 7(1):2, 2019. 636 637 John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region 638 policy optimization, 2017a. 639 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy 640 optimization algorithms, 2017b. URL https://arxiv.org/abs/1707.06347. 641 642 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-643 dimensional continuous control using generalized advantage estimation, 2018. URL https: 644 //arxiv.org/abs/1506.02438. 645 Motahareh Sohrabi, Juan Ramirez, Tianyue H. Zhang, Simon Lacoste-Julien, and Jose Gallego-646 Posada. On pi controllers for updating lagrange multipliers in constrained optimization, 2024. 647 URL https://arxiv.org/abs/2406.04558.

679

680

681 682

683

684

696 697

699 700

648	Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning by
649	pid lagrangian methods. 2020. URL https://arxiv.org/abs/2007.03964.
650	I 6 6 6 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1

- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller (eds.), Advances in Neural Information Processing Systems, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper\_files/paper/1999/ file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
- Manan Tomar, Lior Shani, Yonathan Efroni, and Mohammad Ghavamzadeh. Mirror descent policy
   optimization. In *International Conference on Learning Representations*, 2022. URL https:
   //openreview.net/forum?id=aB05SvgSt1.
- Ilnura Usmanova, Yarden As, Maryam Kamgarpour, and Andreas Krause. Log barriers for safe black-box optimization with application to safe reinforcement learning. *Journal of Machine Learning Research*, 25(171):1–54, 2024.
- Jesse van Oostrum, Johannes Müller, and Nihat Ay. Invariance properties of the natural gradient in overparametrised systems. *Information geometry*, 6(1):51–67, 2023.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256, 1992.
- Long Yang, Jiaming Ji, Juntao Dai, Linrui Zhang, Binbin Zhou, Pengfei Li, Yaodong Yang, and
   Gang Pan. Constrained update projection approach to safe policy optimization. Advances in
   *Neural Information Processing Systems*, 35:9111–9124, 2022.
- Tsung-Yen Yang, Justinian Rosca, Karthik Narasimhan, and Peter J. Ramadge. Projection-based constrained policy optimization, 2020. URL https://arxiv.org/abs/2010.03152.
- Baohe Zhang, Yuan Zhang, Lilli Frison, Thomas Brox, and Joschka Bödecker. Constrained reinforcement learning with smoothed log barrier function. *arXiv preprint arXiv:2403.14508*, 2024a.
- Linrui Zhang, Li Shen, Long Yang, Shixiang Chen, Bo Yuan, Xueqian Wang, and Dacheng Tao.
   Penalized proximal policy optimization for safe reinforcement learning, 2022. URL https: //arxiv.org/abs/2205.11814.
  - Yiming Zhang, Quan Vuong, and Keith Ross. First order constrained optimization in policy space. Advances in Neural Information Processing Systems, 33:15338–15349, 2020.
  - Yufeng Zhang, Jialu Pan, Li Ken Li, Wanwei Liu, Zhenbang Chen, Xinwang Liu, and Ji Wang. On the properties of kullback-leibler divergence between multivariate gaussian distributions. *Advances in Neural Information Processing Systems*, 36, 2024b.

# 702 A EXTENDED BACKGROUND

We consider the infinite-horizon discounted Markov decision process (MDP), given by the tuple ( $S, A, P, r, \mu, \gamma$ ). Here, S and A are the finite state-space and action-space respectively. Here, we make the restriction to finite MDPs as this simplifies the presentation. For a discussion of continuous state and action spaces, we refer to Appendix B.3. Further,  $P: S \times A \to \Delta_S$  is the transition kernel,  $r: S \times A \to \mathbb{R}$  is the reward function,  $\mu \in \Delta_S$  is the initial state distribution at time t = 0, and  $\gamma \in [0, 1)$  is the discount factor. The space  $\Delta_S$  is the set of categorical distributions over S.

The Reinforcement Learning (RL) protocol is usually described as follows: At time t = 0, an initial state  $s_0$  is drawn from  $\mu$ . At each integer time-step t, the agent chooses an action according to it's (stochastic) behavior policy  $a_t \sim \pi(\cdot|s_t)$ . A reward  $r_t = r(s_t, a_t)$  is given to the agent, and a new state  $s_{t+1} \sim P(\cdot|s_t, a_t)$  is sampled from the environment. Given a policy  $\pi$ , the value function  $V_r^{\pi}: S \to \mathbb{R}$ , action-value function  $Q_r^{\pi}: S \times \mathcal{A} \to \mathbb{R}$ , and advantage function  $A_r^{\pi}: S \times \mathcal{A} \to \mathbb{R}$ associated with the reward r are defined as

$$V_r^{\pi}(s) \coloneqq (1-\gamma) \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0 = s \right],$$

721

724 725

729

733 734 735

739

740 741

744 745

747 748 749

752

753 754

716

$$Q_{r}^{\pi}(s,a) \coloneqq (1-\gamma) \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} r(s_{t},a_{t}) \Big| s_{0} = s, a_{0} = a \right] \text{ and } A_{r}^{\pi}(s,a) \coloneqq Q_{r}^{\pi}(s,a) - V_{r}^{\pi}(s)$$

where and the expectations are taken over trajectories of the Markov process resulting from starting at s and following policy  $\pi$ . The goal is to

$$\operatorname{maximize}_{\pi \in \Pi} V_r^{\pi}(\mu) \tag{25}$$

where  $V_r^{\pi}(\mu)$  is the expected value under the initial state distribution  $V_r^{\pi}(\mu) \coloneqq \mathbb{E}_{s \sim \mu}[V_r^{\pi}(s)]$ . We will also write  $V_r^{\pi} = V_r^{\pi}(\mu)$ , and omit the explicit dependence on  $\mu$  for convenience, and we write  $V_r(\pi)$  when we want to emphasize its dependence on  $\pi$ .

**The Dual Linear Program for MDPs** Any stationary policy  $\pi$  induces a discounted state-action (occupancy) measure  $d_{\pi} \in \Delta_{S \times A}$ , indicating the relative frequencies of visiting a state-action pair, discounted by how far the visitation lies in the future. It is a probability measure defined as

$$d_{\pi}(s,a) \coloneqq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\pi}(s_t = s) \pi(a|s),$$

$$(26)$$

where  $\mathbb{P}_{\pi}(s_t = s)$  is the probability of observing the environment in state *s* at time *t* given the agent follows policy  $\pi$ . For finite MDPs, it is well-known that maximizing the expected discounted return can be expressed as the linear program

$$\max_{d} r^{\top} d \quad \text{subject to } d \in \mathscr{D}, \tag{27}$$

where  $\mathscr{D}$  is the set of feasible state-action measures Feinberg & Shwartz (2012). This set is also known as the *state-action polytope*, defined by

$$\mathscr{D} = \left\{ d \in \mathbb{R}_{\geq 0}^{\mathcal{S} \times \mathcal{A}} : \ell_s(d) = 0 \text{ for all } s \in \mathcal{S} \right\}$$

where the linear constraints  $\ell_s(d)$  are given by the *Bellman flow equations* 

$$\ell_s(d) = d(s) - \gamma \sum_{s',a'} d(s',a') P(s|s',a') - (1-\gamma)\mu(s),$$

where  $d(s) = \sum_{a} d(s, a)$  denotes the state-marginal of d. For any state-action measure d we obtain the associated policy via conditioning, meaning

$$\pi(a|s) \coloneqq \frac{d(s,a)}{\sum_{a'} d(s,a')} \tag{28}$$

<sup>755</sup> in case this is well-defined. This provides a one-to-one correspondence between policies and the state-action distributions under the following assumption.

Figure 4: The dual linear program for a CMDP of two states and two actions.

Assumption 6 (Exploration). For any policy  $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$  we have  $d_{\pi}(s) > 0$  for all  $s \in \mathcal{S}$ .

This assumption is standard in linear programming approaches and policy gradient methods where it is necessary for global convergence Kallenberg (1994); Mei et al. (2020b). Note that  $d \in \partial D$  if and only if d(s, a) = 0 for some s, a and hence the boundary of D is given by

$$\partial \mathscr{D} = \left\{ d_{\pi} : \pi(a|s) = 0 \text{ for some } s \in \mathcal{S}, a \in \mathcal{A} \right\}.$$

**Constrained Markov Decision Processes** Where MDPs aim to maximize the return, constrained MDPs (CMDPs) aim to maximize the return subject to a number of costs not exceeding certain thresholds. For a general treatment of CMDPs, we refer the reader to Altman (1999). An important application of CMDPs is in safety-critical reinforcement learning where the costs incorporate safety constraints. An infinite-horizon discounted CMDP is defined by the tuple  $(S, A, P, r, \mu, \gamma, C)$ , consisting of the standard elements of an MDP and an additional constraint set  $C = \{(c_i, b_i)\}_{i=1}^m$ , where  $c_i: S \times A \to \mathbb{R}$  are the cost functions and  $b_i \in \mathbb{R}$  are the cost thresholds.

In addition to the value functions and the advantage functions of the reward that are defined for the MDP, we define the same quantities  $V_{c_i}$ ,  $Q_{c_i}$ , and  $A_{c_i}$  w.r.t the *i*th cost  $c_i$ , simply by replacing rwith  $c_i$ . The objective is to maximize the discounted return, as before, but we restrict the space of policies to the safe policy set

791

792 793 794

796 797 798

803

756

768 769 770

771

775 776 777

 $\Pi_{\text{safe}} = \bigcap_{i=1}^{m} \Big\{ \pi : V_{c_i}(\pi) \le b_i \Big\},\tag{29}$ 

where

$$V_{c_i}^{\pi}(\mu) \coloneqq \mathbb{E}_{s \sim \mu}[V_{c_i}^{\pi}(s)].$$
(30)

is the expected discounted cumulative cost associated with the cost function  $c_i$ . Like the MDP, the discounted cost CMPD can be expressed as the linear program

$$\max_{d} r^{\top} d \quad \text{sbj. to } d \in \mathscr{D}_{\text{safe}}, \tag{31}$$

where

$$\mathscr{D}_{\text{safe}} = \bigcap_{i=1}^{m} \left\{ d \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}} : c_i^\top d \le b_i \right\} \cap \mathscr{D}$$
(32)

is the safe occupancy set, see Figure 4.

804 **Bregman divergences** Here, we give a short introduction to the concept of Bregman divergences, 805 which is required for the formulation of trust region methods. For this, we consider a convex subset 806 of Euclidean space  $C \subseteq \mathbb{R}^d$  with a non-empty interior int(C) and a strictly convex function  $\phi: C \rightarrow$ 807  $\mathbb{R}$  which we assume to be differentiable on the interior int(C). Then, the *Bregman divergence* 808 induced by  $\phi$  is given by

$$D_{\phi}(x||y) \coloneqq \phi(x) - \phi(y) - \nabla \phi(y)^{\top}(x-y), \tag{33}$$



which is well defined for  $x \in C, y \in int(C)$ . Intuitively, the Bregman divergence measures the difference between  $\phi$  and its linearization at y. The strict convexity of  $\phi$  ensures that  $D_{\phi}(x||y) \ge 0$ and  $D_{\phi}(x||y) = 0$  if and only if x = y. Therefore, Bregman divergences are commonly interpreted as a generalized measure for the distance between points, however, it is important to notice that it is not generally symmetric. An important example is the Euclidean distance  $D_{\phi}(x||y) = ||x - y||_2^2$ which arises from the choice  $\phi(x) := \|x\|_2^2$ . Another important Bregman divergence is the Kullback-Leibler (KL) divergence 

$$D_{\rm KL}(p||q) \coloneqq \sum_{i=1}^{d} p_i \log \frac{p_i}{q_i} - \sum_{i=1}^{d} p_i + \sum_{i=1}^{d} q_i,$$
(34)

where we use the common convention  $0 \log \frac{0}{0} \coloneqq 0$ . Then, the KL divergence is defined for  $p \in \mathbb{R}^d_{\geq 0}$ and  $q \in \mathbb{R}^d_{\geq 0}$  which is absolutely continuous with respect to p, meaning that  $p_i = 0$  implies  $q_i = 0$ . Note that if both p and q are probability vectors, meaning that  $\sum_i p_i = \sum_i q_i = 1$ , we obtain 

$$D_{\mathrm{KL}}(p||q) \coloneqq \sum_{i=1}^{d} p_i \log \frac{p_i}{q_i}.$$
(35)

> Information Geometry of Policy Optimization Among the most successful policy optimization schemes are natural policy gradient (NPG) methods or variants thereof like trust-region and proximal policy optimization (TRPO and PPO, respectively). These methods assume a convex geometry and corresponding Bregman divergences in the state-action polytope, where we refer to Neu et al. (2017); Müller & Montúfar (2023) for a more detailed discussion.

In general, a trust region update is defined as

$$\pi_{k+1} \in \operatorname*{arg\,max}_{\pi \in \Pi} \mathbb{A}_r^{\pi_k}(\pi) \quad \text{sbj. to } D_{\Phi}(d_{\pi_k} || d_{\pi}) \le \delta,$$
(36)

where  $D_{\Phi}: \mathscr{D} \times \mathscr{D} \to \mathbb{R}$  is a Bregman divergence induced by a suitably convex function  $\Phi: \operatorname{int}(\mathscr{D}) \to \mathbb{R}$ . The functional

$$\mathbb{A}_r^{\pi_k}(\pi) = \mathbb{E}_{s \sim d_{\pi_k}, a \sim \pi_\theta(\cdot|s)} \big[ A_r^{\pi_k}(s, a) \big], \tag{37}$$

as introduced in (Kakade & Langford, 2002), is called the policy advantage. As a loss function, it is also known as the surrogate advantage (Schulman et al., 2017a), since we can interpret  $\mathbb{A}$  as a surrogate optimization objective of the return. In particular, it holds for a parameterized policy  $\pi_{\theta}$ , that  $\nabla_{\theta} \mathbb{A}_{r}^{\pi_{\theta_{k}}}(\pi_{\theta})|_{\theta=\theta_{k}} = \nabla_{\theta} V_{r}(\theta_{k})$ , see Kakade & Langford (2002); Schulman et al. (2017a). TRPO and the original NPG assume the same geometry (Kakade, 2001; Schulman et al., 2017a), since they employ an identical Bregman divergence 

$$D_{\mathcal{K}}(d_{\pi_1}||d_{\pi_2}) \coloneqq \sum_{s,a} d_{\pi_1}(s,a) \log \frac{\pi_1(a|s)}{\pi_2(a|s)} = \sum_s d_{\pi_1}(s) D_{\mathcal{KL}}(\pi_1(\cdot|s)||\pi_2(\cdot|s)).$$

We refer to  $D_{\rm K}$  as the Kakade divergence and informally write  $D_{\rm K}(\pi_1, \pi_2) := D_{\rm K}(d_{\pi_1}, d_{\pi_2})$ . This divergence can be shown to be the Bregman divergence induced by the negative conditional entropy

$$\Phi_{\mathrm{K}}(d_{\pi}) \coloneqq \sum_{s,a} d_{\pi}(s,a) \log \pi(a|s), \tag{38}$$

see Neu et al. (2017). It is well known that with a parameterized policy  $\pi_{\theta}$ , a linear approximation of A and a quadratic approximation of the Bregman divergence  $D_{\rm K}$  at  $\theta$ , one obtains the *natural policy gradient* step given by 

$$\theta_{k+1} = \theta_k + \epsilon_k G_{\rm K}(\theta_k)^+ \nabla R(\theta_k), \tag{39}$$

where  $G_{\rm K}(\theta)^+$  denotes a pseudo-inverse of the Gramian matrix with entries equal to the stateaveraged Fisher-information matrix of the policy

$$G_{\mathrm{K}}(\theta)_{ij} \coloneqq \mathbb{E}_{s \sim d_{\pi_{\theta}}} \left[ \sum_{a} \frac{\partial_{\theta_{i}} \pi_{\theta}(a|s) \partial_{\theta_{j}} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} \right]$$
(40)

$$= \mathbb{E}_{d_{\pi_{\theta}}}[\partial_{\theta_i} \log \pi_{\theta}(a|s)\partial_{\theta_j} \log \pi_{\theta}(a|s)], \tag{41}$$

where we refer to Schulman et al. (2017a) for a more detailed discussion.

Consider a convex potential  $\Phi: \mathscr{D} \to \mathbb{R}$  or  $\Phi: \mathscr{D}_{safe} \to \mathbb{R}$  and the TRPO update

$$\theta_{k+1} \in \arg \max \mathbb{A}_r^{\pi_{\theta_k}}(\pi_{\theta}) \quad \text{sbj. to } D_{\Phi}(d_{\theta_k} || d_{\theta}) \le \epsilon.$$
 (42)

In practice, one uses a linear approximation of  $\mathbb{A}_r^{\pi_{\theta_k}}(\pi_{\theta})$  and a quadratic approximation of  $D_{\Phi}$  to compute the TRPO update. This gives the following approximation of TRPO

$$\theta_{k+1} \in \arg\max_{\theta} \nabla_{\theta} \mathbb{A}_r^{\theta_k}(\theta)|_{\theta=\theta_k} \cdot (\theta - \theta_k) \quad \text{sbj. to } \|\theta - \theta_k\|_{G(\theta_k)}^2 \le \epsilon,$$
(43)

where

$$G(\theta)_{ij} = \partial_{\theta_i} d_{\theta}^{\top} \nabla^2 \Phi(d_{\theta}) \partial_{\theta_j} d_{\theta}.$$
(44)

876 Note that by the policy gradient theorem, it holds that

$$\nabla_{\theta} \mathbb{A}_{r}^{\theta_{k}}(\theta)|_{\theta=\theta_{k}} = \nabla V_{r}(\theta_{k}).$$
(45)

Thus, the approximate TRPO update is equivalent to

$$\theta_{k+1} = \theta_k + \epsilon_k G(\theta_k)^+ \nabla V_r(\theta), \tag{46}$$

882 where

$$\epsilon_k = \frac{\sqrt{\epsilon}}{\|G(\theta_k)^+ \nabla V_r(\theta_k)\|_{G(\theta_k)}}.$$
(47)

Hence, the approximation TRPO update corresponds to a natural policy gradient update with an adaptively chosen step size.

### B DETAILS ON THE SAFE GEOMETRY FOR CMDPS

B.1 SAFE TRUST REGIONS

The safe mirror function for a single constraint is given by

$$\Phi_{\mathcal{C}}(d) \coloneqq \Phi_{\mathcal{K}}(d) + \sum_{i=1}^{m} \beta \, \phi(b - c^{\top} d), \tag{48}$$

and the resulting Bregman divergence

=

$$D_{\rm C}(d_1||d_2) = \Phi_{\rm C}(d_1) - \Phi_{\rm C}(d_2) - \langle \nabla \Phi_{\rm C}(d_2), d_1 - d_2 \rangle.$$
(49)

is a linear operator in  $\Phi$ , hence

$$D_{\Phi(d)+\beta\phi(b-c^{\top}d)}(d_1||d_2) = D_{\Phi_{\mathrm{K}}}(d_1||d_2) + \beta D_{\phi}(d_1||d_2),$$
(50)

where

$$D_{\phi}(d_1||d_2) = \phi(b - c^{\top}d_1) - \phi(b - c^{\top}d_2) - \langle \nabla \phi(b - c^{\top}d_2), d_1 - d_2 \rangle$$
(51)

$$= \phi(b - c^{\top}d_1) - \phi(b - c^{\top}d_2) - \phi'(b - c^{\top}d_2)(c^{\top}d_1 - c^{\top}d_2).$$
(52)

$$=\phi(b-V_c(\pi))-\phi(b-V_c(\pi_k))+\phi'(b-V_c(\pi_k))(V_c(\pi)-V_c(\pi_k)).$$
(53)

The last expression can be interpreted as the one-dimensional Bregman divergence  $D_{\phi}(b - V_c(\pi)||b - V_c(\pi_k))$ , which is a (strictly) convex function in  $V_c(\pi)$  for fixed  $\pi_k$  if  $\phi$  is (strictly) convex.

912 B.2 DETAILS ON C-TRPO 

#### 914 B.2.1 SURROGATE DIVERGENCE

In practice, the exact constrained KL-Divergence  $D_{\rm C}$  cannot be evaluated, because it depends on the cost-return of the optimized policy  $V_c(\pi)$ . Therefore, we use the surrogate divergence

$$\bar{D}_{\phi}(\pi_{\theta}||\pi_{\theta_{k}}) = \phi(b - V_{c}^{\pi_{k}} - \mathbb{A}_{c}^{\pi_{k}}(\pi)) - \phi(b - V_{c}^{\pi_{k}}) + \phi'(b - V_{c}^{\pi_{k}})\mathbb{A}_{c}^{\pi_{k}}(\pi)$$
(54)

is obtained by the substitution  $V_c(\pi) - V_c^{\pi_k} \to \mathbb{A}_c^{\pi_k}(\pi)$  in  $D_{\phi}$ .

When we center this divergence around policy  $\pi_k$  and keep this policy fixed, it becomes a function of the policy cost advantage.

$$\bar{D}_{\phi}(\pi_{\theta}||\pi_{\theta_{k}}) = \phi(b - V_{c}^{\pi_{k}} - \mathbb{A}_{c}^{\pi_{k}}(\pi)) - \phi(b - V_{c}^{\pi_{k}}) + \phi'(b - V_{c}^{\pi_{k}})\mathbb{A}_{c}^{\pi_{k}}(\pi)$$
$$= \phi(\delta_{b} - \mathbb{A}_{c}^{\pi_{k}}(\pi)) - \phi(\delta_{b}) + \phi'(\delta_{b})\mathbb{A}_{c}^{\pi_{k}}(\pi)$$
$$= \Psi(\mathbb{A}_{c}^{\pi_{k}}).$$

929 Note that  $\bar{D}_{\phi}(\pi_{\theta}||\pi_{\theta_{k}}) = \Psi(\mathbb{A}_{c}^{\pi_{k}}(\pi))$ , where  $\Psi(x) = \phi(\delta_{b} - x) - \phi(\delta_{b}) - \phi'(\delta_{b}) \cdot x$  is a (strictly) 930 convex function if  $\phi$  is (strictly) convex, since it is equivalent to the one-dimensional Bregman divergence  $D_{\phi}(\delta_{b} - x||\delta_{b})$  on the domain of  $\phi(b - x)$ , see Figure 5.



Figure 5: The surrogate Constrained KL-Divergence as a function of the policy cost advantage.

**Example 7.** The function  $\phi(x) = x \log(x)$  induces the divergence

$$\bar{D}_{\phi}(\pi_{\theta}||\pi_{\theta_{k}}) = \mathbb{A}_{c}^{\pi_{k}}(\pi_{\theta}) - (\delta_{b} - \mathbb{A}_{c}^{\pi_{k}}(\pi_{\theta})) \log\left(\frac{\delta_{b}}{\delta_{b} - \mathbb{A}_{c}^{\pi_{k}}(\pi_{\theta})}\right).$$
(55)

#### **B.2.2** ESTIMATION

In the practical implementation, the expected KL-divergence between the policy of the previous iteration,  $\pi_k$ , and the proposal policy  $\pi$  is estimated from state samples  $s_i$  by running  $\pi_k$  in the environment

$$\sum_{s} d_{\pi_{k}}(s) D_{\mathrm{KL}}(\pi(\cdot|s)||\pi_{k}(\cdot|s)) \approx 1/N \sum_{i=0}^{N-1} D_{\mathrm{KL}}(\pi(\cdot|s_{i})||\pi_{k}(\cdot|s_{i}))$$
(56)

where  $D_{\text{KL}}$  can be computed in closed form for Gaussian policies, where N is the batch size.

For the constraint term, we estimate  $\delta_b$  from trajectory samples, as well as the policy cost advantage

$$\mathbb{A}_{c}^{\pi_{k}}(\pi) \approx \hat{\mathbb{A}} = \frac{1}{N} \sum_{i=0}^{N-1} \frac{\pi(a_{i}|s_{i})}{\pi_{k}(a_{i}|s_{i})} \hat{A}_{i}^{\pi_{k}}$$
(57)

where  $\hat{A}_i^{\pi_k}$  is the GAE- $\lambda$  estimate of the advantage function (Schulman et al., 2018). For any suitable  $\phi$ , the resulting divergence estimate is

$$\hat{D}_{\phi} = \phi(\delta_b - \hat{\mathbb{A}}) - \phi(\delta_b) - \phi'(\delta_b)\hat{\mathbb{A}}$$
(58)

and for the specific choice  $\phi(x) = x \log(x)$ 

$$\hat{D}_{\phi} = \hat{\mathbb{A}} - (\delta_b - \hat{\mathbb{A}}) \log\left(\frac{\delta_b}{\delta_b - \hat{\mathbb{A}}}\right).$$
(59)

#### 972 B.2.3 DETAILS ON C-NPG 973

In showin that TRPO with quadratic approximation agrees with a natural gradient step, see Appendix A, we have used that  $\nabla_{\theta} \mathbb{A}_{r}^{\theta_{k}}(\theta)|_{\theta=\theta_{k}} = \nabla V_{r}(\theta_{k})$ , which holds although  $\mathbb{A}_{r}$  is only a proxy of  $V_{r}$ . We now provide a similar property for the quadratic approximation of the surrogate divergences  $\overline{D}_{C}$ .

**Proposition 8.** For any parameter  $\theta$  with  $\pi_{\theta} \in \prod_{\text{safe}}$  it holds that

$$\nabla^2_{\theta} \bar{D}_{\phi}(\theta || \hat{\theta}) |_{\theta = \hat{\theta}} = \nabla^2_{\theta} D_{\phi}(\theta || \hat{\theta}) |_{\theta = \hat{\theta}}$$

$$\tag{60}$$

and hence

979 980 981

982 983 984

989

$$\nabla^2_{\theta} \bar{D}_{\mathrm{KL}}(\theta || \hat{\theta})|_{\theta = \hat{\theta}} + \beta \nabla^2_{\theta} \bar{D}_{\phi}(\theta || \hat{\theta})|_{\theta = \hat{\theta}} = G_{\mathrm{C}}(\hat{\theta})$$
(61)

where  $G_{\rm C}(\theta)$  denotes the Gramian matrix of C-NPG with entries

$$G_{\rm C}(\theta)_{ij} = \partial_{\theta_i} d_{\theta}^{\dagger} \nabla^2 \Phi_{\rm C}(\theta) \partial_{\theta_j} d_{\theta}.$$
(62)

*Proof.* Let  $\bar{H}_{KL}(\theta) = \nabla^2_{\theta} \bar{D}_{KL}(\theta || \hat{\theta})|_{\theta=\hat{\theta}}$  and  $\bar{H}_{\phi}(\theta) = \nabla^2_{\theta} \bar{D}_{\phi}(\theta || \hat{\theta})|_{\theta=\hat{\theta}}$ . One can show that  $\bar{H}_{KL} = G_K(\theta)$  (Schulman et al., 2017a). Further, we have

$$\bar{H}_{\phi}(\theta) = \nabla_{\theta} \mathbb{A}_{c}^{\pi_{k}}(\theta)^{\top} \Psi''(\mathbb{A}_{c}^{\pi_{k}}(\theta)) \nabla_{\theta} \mathbb{A}_{c}^{\pi_{k}}(\theta) + \Psi'(\mathbb{A}_{c}^{\pi_{k}}(\theta))^{\top} \nabla_{\theta}^{2} \mathbb{A}_{c}^{\pi_{k}}(\theta)$$

$$\stackrel{a)}{=} \nabla_{\theta} \mathbb{A}_{c}^{\pi_{k}}(\theta)^{\top} \Psi''(\mathbb{A}_{c}^{\pi_{k}}(\theta)) \nabla_{\theta} \mathbb{A}_{c}^{\pi_{k}}(\theta)$$

$$\stackrel{b)}{=} \nabla_{\theta} \mathbb{A}_{c}^{\pi_{k}}(\theta)^{\top} \phi''(b - V_{c}^{\pi_{k}}(\theta)) \nabla_{\theta} \mathbb{A}_{c}^{\pi_{k}}(\theta)$$

$$= \nabla_{\theta} V_{c}^{\pi_{k}}(\theta)^{\top} \phi''(b - V_{c}^{\pi_{k}}(\theta)) \nabla_{\theta} V_{c}^{\pi_{k}}(\theta),$$

where a) follows from  $\Psi'(\mathbb{A}_{c}^{\pi_{k}}(\theta)) = 0$  since  $\Psi(0) = 0$ ,  $\Psi \ge 0$  and  $\mathbb{A}_{c}^{\hat{\theta}}(\theta)|_{\theta=\hat{\theta}} = 0$ . Further, b) follows because  $\Psi''(x)|_{x=0} = \phi''(\delta_{b})$ . Thus,  $\bar{H}_{\phi}$  is equivalent to the Gramian

$$G_{\mathcal{C}}(\theta)_{ij} \coloneqq \partial_{\theta_i} d_{\theta}^{\top} \nabla^2 \Phi_{\mathcal{C}}(\theta) \partial_{\theta_j} d_{\theta}$$
(63)

$$= G_{\rm K}(\theta)_{ij} + \beta \phi''(b - c_k^{\rm T} d_\theta) \partial_{\theta_i} d_\theta^{\rm T} c c^{\rm T} \partial_{\theta_i} d_\theta$$
(64)

$$=\bar{H}_{\mathrm{KL}} + \beta \nabla_{\theta} V_{c}(\theta)^{\top} \phi^{\prime\prime}(b - V_{c}(\theta)) \nabla_{\theta} V_{c}(\theta), \tag{65}$$

$$=\bar{H}_{\rm KL} + \beta \bar{H}_{\phi}.$$
(66)

1005 1006

1011

1013

1000

1007

In particular, this shows that the C-TRPO update can be interpreted as a natural policy gradient step with an adaptive step size and that the updates with  $D_{\rm C}$  and  $\bar{D}_{\rm C}$  are equivalent if we use a quadratic approximation for both, justifying  $\bar{D}_{\rm C}$  as a surrogate for  $D_{\rm C}$ .

# 1012 B.3 BEYOND FINITE MDPs

For the sake of simplicity and as this is required for our theoretical analysis, we have introduced
C-TRPO only for finite MDPs. However, C-TRPO can also be used for problems with continuous
state and action spaces as we discuss here. In this case, the state-action and state distributions are
defined as

$$d_{\pi}(S \times A) \coloneqq (1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} \mathbb{P}_{\pi}(s_{t} \in S, a_{t} \in A) \quad \text{and}$$

Again, for multiple constraints, the statement follows analogously.

1021  
1022  
1022  

$$d_{\pi}(S) \coloneqq (1-\gamma) \sum_{t=0}^{\infty} \gamma^{t} \mathbb{P}_{\pi}(s_{t} \in S)$$

1024 for every measurable subsets 
$$A \subseteq \mathcal{A}$$
 and  $S \subseteq \mathcal{S}$ . Further, the Kakade divergence is then given by  
1025  $D_{\mathrm{K}}(d^{\pi_1}||d^{\pi_2}) \coloneqq \mathbb{E}_{s \sim d^{\pi_1}} [D_{\mathrm{KL}}(\pi_1(\cdot|s)||\pi_2(\cdot|s))],$  (67)

which is well defined if  $\pi_1(\cdot|s)$  is absolutely continuous with respect to  $\pi_2(\cdot|s)$  for  $d^{\pi_1}$  almost all  $s \in S$ . The Bregman divergence that C-TRPO is builds on is – just as in the finite case – given by

1033 1034

1035 1036 1037

1041 1042

1044

1045

1050

1057 1058

$$D_{\rm C}(d_1||d_2) = D_{\rm K}(d_1||d_2) + \sum_{i=1}^m \beta_i D_{\phi_i}(d_1||d_2), \tag{68}$$

1031 1032 where

$$D_{\phi_i}(d_1||d_2) = \phi(b_i - V_{c_i}(\pi_1)) - \phi(b_i - V_{c_i}(\pi_2)) + \phi'(b_i - V_{c_i}(\pi_2))(V_{c_i}(\pi_1) - V_{c_i}(\pi_2)).$$
(69)

Like in the finite case, the policy advantage is defined as

$$\mathbb{A}_r^{\pi_k}(\pi) = \mathbb{E}_{s,a \sim d_{\pi_k}} \left[ \frac{\pi(a|s)}{\pi_k(a|s)} A_r^{\pi_k}(s,a) \right],\tag{70}$$

(72)

where  $A_r^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$  denotes the advantage function, which is defined analoguously to the finite case. Now, the plain trust region update is given b y

$$\theta_{k+1} \in \operatorname*{arg\,max}_{\theta} \mathbb{A}^{\pi_k}_r(\pi) \quad \text{sbj. to } D_{\mathbb{C}}(d_{\pi_k} || d_{\pi}) \le \delta.$$
(71)

Just like in the finite case, we use a surrogate divergence  $D_{\rm C}$  and obtain the formulation of C-TRPO

$$\pi_{k+1} = \operatorname*{arg\,max}_{\pi \in \Pi} \mathbb{A}_r^{\pi_k}(\pi) \quad \text{sbj. to } \bar{D}_{\mathrm{C}}(\pi || \pi_k) \le \delta.$$

1046 Here, the differences to  $D_{\rm C}$  are that we use use samples from the state distribution  $d^{\pi_k}$  and use a 1047 surrogate for the cost advantage to estimate the divergence  $D_{\pi_i}$  as described in Section 3.2. Fur-1048 ther, we use a parametric policy model  $\pi_{\theta}$  and a linear approximation of  $\mathbb{A}^{\pi_k}$  as well as quadratic 1049 approximation of  $\overline{D}_{\rm C}(\pi || \pi_k)$  for our practical implementation.

**Expression for Gaussian policies** We test C-TRPO in various control tasks and hence, where we use Gaussian policies. More precisely, the state and action space consist of Euclidean spaces  $\mathcal{S} = \mathbb{R}^{d_s}$  and  $\mathcal{A} = R^{d_a}$ . Then, we consider a policy network  $\mu_{\theta} \colon \mathcal{S} \to \mathcal{A}$ , which predicts the mean action and assume parameterized but state independent diagonal Gaussian noise, meaning that  $\pi_{\theta}(\cdot|s) = \mathcal{N}(\mu_{\theta}(s), \Sigma_{\theta})$ , where  $\Sigma_{\theta}$  is diagonal. Consequently, we can use a closed-form expression for the KL divergence as

$$D_{\mathrm{KL}}(\pi_{\theta_{1}}(\cdot|s)||\pi_{\theta_{2}}(\cdot|s)) = \frac{1}{2} \left( \operatorname{tr}\left(\Sigma_{\theta_{2}}^{-1}\Sigma_{\theta_{1}}\right) - d_{\mathsf{a}} + \|\mu_{\theta_{1}}(s) - \mu_{\theta_{2}}(s)\|_{\Sigma_{\theta_{2}}^{-1}}^{2} + \ln\left(\frac{\det \Sigma_{\theta_{2}}}{\det \Sigma_{\theta_{1}}}\right) \right)$$

see Zhang et al. (2024b).

# C PROOFS OF SECTION 4

1064 C.1 PROOFS OF SECTION 4.1

Our theoretical analysis of C-TRPO is built on the following bounds on the performance difference of two policies.

**Theorem 9** (Performance Difference, Achiam et al. (2017)). For any function f(s, a), the following bounds hold

$$V_{f}(\pi_{1}) - V_{f}(\pi_{2}) \leq \mathbb{A}_{f}^{\pi_{2}}(\pi_{1}) \pm \frac{2\gamma\epsilon_{f}}{(1-\gamma)} \sqrt{\frac{1}{2}\mathbb{E}_{s \sim d_{\pi_{2}}} D_{\mathrm{KL}}(\pi_{1}(\cdot|s)||\pi_{2}(\cdot|s))}$$
(73)

1072 where  $\epsilon_f = \max_s |\mathbb{E}_{a \sim \pi_1} A_f^{\pi_2}(s, a)|.$ 

1074 Theorem 9 can be interpreted as a bound on the error incurred by replacing the difference in returns 1075  $V_f(\pi_1) - V_f(\pi)$  of any state-action function by its policy advantage  $\mathbb{A}_f^{\pi_2}(\pi_1)$ .

**Proposition 1** (C-TRPO reward update). Set  $\epsilon_r = \max_s |\mathbb{E}_{a \sim \pi_{k+1}} A_r^{\pi_k}(s, a)|$ . The expected reward of a policy updated with C-TRPO is bounded from below by

$$V_r(\pi_{k+1}) \ge V_r(\pi_k) - \frac{\sqrt{2\delta\gamma\epsilon_r}}{1-\gamma}.$$
(20)

1062 1063

1070 1071

*Proof.* It follows from the lower bound in Theorem 9 that

$$V_r(\pi_{k+1}) - V_r(\pi_k) \ge \mathbb{A}_r^{\pi_k}(\pi_{k+1}) - \frac{\gamma \epsilon_r}{(1-\gamma)} \sqrt{2\bar{D}_{\mathcal{C}}(\pi_{k+1}||\pi_k)}$$
(74)

where we choose f = r. The bound holds because  $\bar{D}_{\phi} \ge 0$ , and thus  $\bar{D}_{C} \ge \mathbb{E}D_{KL}$ . Further,  $\delta \ge D_{C}$  and  $\mathbb{A}_{r}^{\pi_{k}}(\pi_{k+1}) \ge 0$  by the update equation, which concludes the proof. See Appendix C.3 for a more detailed discussion.

**Proposition 2.** The approximate C-TRPO update approaches the CPO update in the limit as  $\beta \searrow 0$ .

**Proof.** Let us fix a strictly safe policy  $\pi_0 \in int(\Pi_{safe})$ . In both cases, we approximate the expected cost of a policy using  $V_c(\pi) \approx V_c(\pi_0) + \mathbb{A}_c^{\pi_0}(\pi)$ , which is off by the advantage mismatch term in Theorem 1. Hence, we maximize the surrogate of the expected value  $\mathbb{A}_r^{\pi_0}(\pi)$  over the regions

$$P_{\text{CPO}} \coloneqq \{\pi \in \Pi : \bar{D}_{\mathrm{K}}(\pi, \pi_0) \le \delta, V_c(\pi_0) + \mathbb{A}_c^{\pi_0}(\pi) \le b\}$$

in the case of CPO, and

$$P_{\beta} := \{ \pi \in \Pi : \overline{D}_{\mathcal{C}}(\pi, \pi_0) < \delta \}.$$

1097 with C-TRPO for some  $\beta > 0$ . Note that 

$$\bar{D}_{\rm C}(\pi,\pi_0) = \bar{D}_{\rm K}(\pi,\pi_0) + \beta \Psi(\mathbb{A}_c^{\pi_0}(\pi)),\tag{75}$$

1100 and  $\Psi: (-\infty, \delta_b) \to (0, +\infty)$  and  $\Psi(t) \to +\infty$  for  $t \nearrow \delta_b$ , where  $\delta_b = b - V_c(\pi_0)$ . Denote the 1101 corresponding updates by  $\hat{\pi}_{CPO}$  and the C-TRPO update by  $\hat{\pi}_{\beta}$ . Note that we have  $P_{\beta} \subseteq P_{\beta'} \subseteq P_{CPO}$  for  $\beta \ge \beta'$ . Further, we have

$$\bigcup_{\beta > 0} P_{\beta} = \{ \pi \in P : D_{\mathcal{K}}(\pi, \pi_0) < \delta, V_c(\pi_0) + \mathbb{A}_c^{\pi_0}(\pi) < b \}.$$

Hence, the trust regions  $P_{\beta}$  grow for  $\beta \searrow 0$  and fill the interior of the trust region  $P_{CPO}$ .

Proposition 3 (C-TRPO worst-case constraint violation). Consider  $\Psi : [0, \delta_b) \to [0, \infty)$  defined by  $\Psi(x) = \phi(\delta_b - x) - \phi(\delta_b) - \phi'(\delta_b) \cdot x$  such that  $D_{\phi}(\pi || \pi_k) = \Psi(\mathbb{A}_c^{\pi_k}(\pi))$ . Further, set  $\epsilon_c = \max_s |\mathbb{E}_{a \sim \pi_{k+1}} A_c^{\pi_k}(s, a)|$ , and choose a strictly convex  $\phi$ . The worst-case constraint violation for C-TRPO is

$$V_c(\pi_{k+1}) \le V_c(\pi_k) + \Psi^{-1}(\delta/\beta) + \frac{\sqrt{2\delta\gamma\epsilon_c}}{1-\gamma}.$$
(21)

1117 Further, it holds that  $\lim_{\beta \to +\infty} \Psi^{-1}(\delta/\beta) = 0$  and  $\Psi^{-1}(\delta/\beta) < b - V_c(\pi_k)$  for all  $\beta \in (0, \infty)$ .

1119 *Proof.* Setting f = c in the upper bound from Theorem 9, and replacing  $\mathbb{E}D_{\mathrm{KL}}$  with  $\delta$  as in Propo-1120 sition 1 results in

$$V_c(\pi_{k+1}) \le V_c(\pi_k) + \mathbb{A}_c^{\pi_k}(\pi_{k+1}) + \frac{\sqrt{2\delta\gamma\epsilon_c}}{1-\gamma}.$$
(76)

1123 Recall that  $\overline{D}_{C} = \overline{D}_{K} + \beta \overline{D}_{\phi}$  and that  $\overline{D}_{\phi}(\pi_{k+1} || \pi_{k}) = \Psi(\mathbb{A}_{c}^{\pi_{k}}(\pi_{k+1}))$ , where  $\Psi(x) = \phi(\delta_{b} - x) - \phi(\delta_{b}) - \phi'(\delta_{b}) \cdot x$ . By the definition of the update it holds that

$$\Psi(\mathbb{A}_{c}^{\pi_{k}}(\pi_{k+1})) < \delta/\beta.$$
(77)

1127 Since we are only interested in upper bounding the worst case, we can focus on  $\mathbb{A}_c^{\pi_k}(\pi_{k+1}) > 0$ , so 1128 we restrict  $\Psi : [0, \delta_b) \to [0, \infty)$ . Further, for strictly convex  $\phi$ ,  $\Psi$  is strictly convex and increasing 1129 with increasing inverse. It follows that

$$\mathbb{A}_c^{\pi_k}(\pi_{k+1}) < \Psi^{-1}(\delta/\beta),\tag{78}$$

1132 with  $\Psi^{-1}: [0, \infty) \to [0, \delta_b)$ . Because  $\Psi^{-1}$  is an increasing function of  $\beta$  on  $[0, \infty)$  with maximum 1133 at  $\delta_b = b - V_c(\pi_k)$ , it holds that  $\Psi^{-1}(\beta/\delta) < b - V_c(\pi_k)$  for any  $\beta > 0$ , which concludes the proof.

# 1134 C.2 DETAILS ON THE RESULTS IN SECTION 4.2

1136 Recall that we study the natural policy gradient flow

1137 1138

1140 1141 1142

$$\partial_t \theta_t = G_{\rm C}(\theta_t)^+ \nabla V_r(\theta_t), \tag{79}$$

where  $G_{\rm C}(\theta)^+$  denotes a pseudo-inverse of  $G_{\rm C}(\theta)$  with entries

$$G_{\mathcal{C}}(\theta)_{ij} \coloneqq \partial_{\theta_i} d_{\theta}^{\top} \nabla^2 \Phi_{\mathcal{C}}(d_{\theta}) \partial_{\theta_j} d_{\theta} = G_{\mathcal{K}}(\theta)_{ij} + \sum_k \beta_k \phi''(b_k - c_k^{\top} d_{\theta}) \partial_{\theta_i} d_{\theta}^{\top} c_k c_k^{\top} \partial_{\theta_i} d_{\theta}.$$
(80)

and  $\theta \mapsto \pi_{\theta}$  is a differentiable policy parametrization.

1145 Moreover, we assume that  $\theta \mapsto \pi_{\theta}$  is regular, that it is surjective and the Jacobian is of maximal rank 1146 everywhere. This assumption implies overparametrization but is satisfied for common models like 1147 tabular softmax, tabular escort, or expressive log-linear policy parameterizations (Agarwal et al., 2021a; Mei et al., 2020a; Müller & Montúfar, 2023).

We denote the set of safe parameters by  $\Theta_{safe} := \{\theta \in \mathbb{R}^p : \pi_\theta \in \Pi_{safe}\}$ , which is non-convex in general and say that  $\Theta_{safe}$  is *invariant* under Equation (22) if  $\theta_0 \in \Theta_{safe}$  implies  $\theta_t \in \Theta_{safe}$  for all t. Invariance is associated with safe control during optimization and is typically achieved via control barrier function methods (Ames et al., 2017; Cheng et al., 2019). We study the evolution of the state-action distributions  $d_t = d^{\pi_{\theta_t}}$  as this allows us to employ the linear programming formulation of CMPDs and we obtain the following convergence guarantees.

**Theorem 4** (Safety during training). Assume that  $\phi \colon \mathbb{R}_{>0} \to \mathbb{R}$  satisfies  $\phi'(x) \to +\infty$  for  $x \searrow 0$ and consider a regular policy parameterization. Then the set  $\Theta_{\mathbb{C}}$  is invariant under Equation (22).

1157

*Proof.* Consider a solution  $(\theta_t)_{t>0}$  of Equation (79). As the mapping  $\pi \mapsto d^{\pi}$  is a diffeomor-1158 phism (Müller & Montúfar, 2023) the parameterization  $\Theta_{\text{safe}} \to \mathscr{D}_{\text{safe}}, \theta \mapsto d^{\pi_{\theta}}$  is surjective and 1159 has a Jacobian of maximal rank everywhere. As  $G_{\rm C}(\theta)_{ij} = \partial_{\theta_i} d_{\theta} \nabla \Phi_{\rm C} \partial_{\theta_i} d_{\theta}$  this implies that the 1160 state-action distributions  $d_t = d^{\pi_{\theta_t}}$  solve the Hessian gradient flow with Legendre-type function 1161  $\Phi_{\rm C}$  and the linear objective  $d \mapsto r^{\top} d$ , see Amari (2016); van Oostrum et al. (2023); Müller & 1162 Montúfar (2023) for a more detailed discussion. It suffices to study the gradient flow in the space of 1163 state-action distributions  $d_t$ . It is easily checked that  $\Phi_C$  is a Legendre-type function for the convex domain  $\mathscr{D}_{\mathcal{C}}$ , meaning that it satisfies  $\|\nabla \Phi(d_n)\| \to +\infty$  for  $d_n \to d \in \partial \mathscr{D}_{safe}$ . Since the objective 1164 is linear, it follows from the general theory of Hessian gradient flows of convex programs that the 1165 flow is well posed, see Alvarez et al. (2004); Müller & Montúfar (2023).  $\square$ 1166

**Theorem 5.** Assume that  $\phi'(x) \to +\infty$  for  $x \searrow 0$ , set  $V_{r,C}^* \coloneqq \max_{\pi \in \Pi_{\text{safe}}} V_r(\pi)$  and denote the set of optimal constrained policies by  $\Pi_{\text{safe}}^* = \{\pi \in \Pi_{\text{safe}} : V_r(\pi) = V_{r,C}^*\}$ , consider a regular policy parametrization and let  $(\theta_t)_{t\ge 0}$  solve Equation (22). It holds that  $V_r(\pi_{\theta_t}) \to V_{r,C}^*$  and

$$\lim_{t \to +\infty} \pi_t = \pi_{\text{safe}}^{\star} = \arg\min\{D_{\mathcal{C}}(\pi^{\star}, \pi_0) : \pi^{\star} \in \Pi_{\text{safe}}^{\star}\}.$$
(23)

1173 1174 Proof. Just like in the proof of Theorem 5 we see that  $d_t = d^{\pi_{\theta_t}}$  solves the Hessian gradient flow 1175 with respect to the Legendre type function  $\Phi_C$ . Now the claims regarding convergence and the 1176 identification of the limit  $\lim_{t\to+\infty} \pi_{\theta_t}$  follows from the general theory of Hessian gradient flows, 1176 see Alvarez et al. (2004); Müller et al. (2024).

1177

1171 1172

# 1178 C.3 PERFORMANCE IMPROVEMENT BOUNDS AND CHOICE OF DIVERGENCE

In a series of works (Kakade & Langford, 2002; Pirotta et al., 2013; Schulman et al., 2017a; Achiam
et al., 2017), the following bound on policy performance difference between two policies has been
established.

$$V_f(\pi') - V_f(\pi) \stackrel{\leq}{>} \mathbb{A}_f^{\pi'}(\pi) \pm \frac{2\gamma\epsilon_f}{(1-\gamma)} \mathbb{E}_{s \sim d_\pi} D_{\mathrm{TV}}(\pi'||\pi)(s)$$
(81)

<sup>1185</sup> where  $D_{\rm TV}$  is the Total Variation Distance. Furthermore, by Pinsker's inequality, we have that

1186 1187

1183

$$D_{\rm TV}(\pi'||\pi) \le \sqrt{\frac{1}{2}} D_{\rm KL}(\pi'||\pi),$$
 (82)





1214

and by Jensen's inequality

1217 1218

1220

1221

1223

1225

1226

1227 1228

$$\mathbb{E}_{s \sim d_{\pi}} D_{\mathrm{TV}}(\pi'||\pi)(s) \leq \sqrt{\frac{1}{2}} \mathbb{E}_{s \sim d_{\pi}} D_{\mathrm{KL}}(\pi'||\pi)(s)},\tag{83}$$

1219 It follows that we can not only substitute the KL-divergence into the bound but any divergence

$$D_{\Phi}(d'_{\pi}||d_{\pi}) \ge \mathbb{E}_{s \sim d_{\pi}} D_{\mathrm{KL}}(\pi'||\pi)(s) \tag{84}$$

can be substituted, and still retains TRPO's and CPO's update guarantees.

### 1224 C.4 COMPARISON WITH CPO

In the approximate case of C-TRPO and CPO, where the reward is approximated linearly, and the trust region quadratically, the constraints differ in that C-TRPO's constraint is

$$(\theta - \theta_k)(\bar{H}_{\mathrm{KL}}(\theta) + \beta \bar{H}_{\phi}(\theta))(\theta - \theta_k) < \delta$$

1229 whereas CPO's is

1234

1235

$$(\theta - \theta_k) \bar{H}_{\mathrm{KL}}(\theta) (\theta - \theta_k) < \delta \text{ and } V_c^{\theta_k} + (\nabla_{\theta} \mathbb{A}_c^{\theta_k}(\theta))^\top (\theta - \theta_k) \le b.$$

Figure 6 illustrates the differences between CPO and C-TRPO.

# D ADDITIONAL EXPERIMENTS

1236 1237 D.1 Effects of the hyper-parameters

To better understand the effects of the two hyperparameters  $\beta$  and  $b_{\rm H}$ , we observe how they change the training dynamics through the example of the *AntVelocity* environment.

1241 The safety parameter  $\beta$  modulates the stringency with which C-TRPO satisfies the constraint, without limiting the expected return for values up to  $\beta = 1$ , see Figure 7. For higher values, the expected return starts to degrade, partly due to  $\bar{D}_{\phi}$  being relatively noisy compared to  $\bar{D}_{\text{KL}}$  and thus we recommend the choice  $\beta = 1$ .

Further, we observe that constraint satisfaction is stable across different choices of cost threshold *b*, see Figure 8, and that in most environments, constraint violations seem to reduce as the algorithm converges, meaning that the regret flattens over time. This behavior suggests that the divergence estimation becomes increasingly accurate over time, potentially allowing C-TRPO to achieve sublinear regret. However, we leave regret analysis of the finite sample regime for future research.

Finally, employing a hysteresis fraction  $0 < b_{\rm H} < b$  seems beneficial, possible because it leads the iterate away from the boundary of the safe set, and because divergence estimates tend to be more reliable for strictly safe policies. The effect of the choice of  $b_{\rm H}$  is visualized in Figure 10.





Figure 9: In difficult environments, e.g. those that start off in the unsafe policy set, it seems to be beneficial to set a fraction of the cost limit for hysteresis.

1293

1282

1283

1284

1285

1286

1287

1288

1289

1290 1291

# 1296 D.2 ABLATION STUDY: CPO vs. C-TRPO

We conduct an ablation study to rule out that our improvements of C-TRPO over CPO are only due to hysteresis. For this, we run both CPO and C-TRPO with and without hysteresis with the same hysteresis parameter as in our other experiments. We see that the hysteresis improves safety for both algorithm. Further, we find that the hysteresis slightly reduces the return of C-TRPO. Overall, we clearly see that C-TRPO itself is much safer compared to CPO as even C-TRPO without hysteresis achieves lower cost regret compared to CPO with hysteresis.



Figure 10: Ablation study on the core components of C-TRPO: Safe trust region (C-TRPO no hyst.) and recovery with hysteresis (CPO hyst.). Evaluation is based on the Inter Quartile Mean (IQM) normalized scores across 5 seeds and 8 tasks. From left to right: episode return of the reward (PPO normalized), episode return of the cost (threshold normalized), and cumulative cost violation (CPO normalized).

# 1350 D.3 PERFORMANCE ON INDIVIDUAL ENVIRONMENTS









1458Table 1: Average evaluation performance per task across 10 evaluation runs and 5 seeds each. We1459highlight the best performance with respect to the average return  $V_r$  in bold, and underline the1460lowest average cost  $V_c$ . Note that the table only contains information about the final evaluation1461performance, not about cost violations during training.

		C-TRPO	C-TRPO-HYST		CPO		CPO-HYST	
	$V_r$	$V_c$	$V_r$	$V_c$	$V_r$	$V_c$	$V_r$	$V_c$
AntVelocity	$2810.1\pm45.2$	$7.9 \pm 7.6$	$2786.5 \pm 75.8$	$9.0 \pm 8.2$	$2569.7 \pm 61.4$	$5.7 \pm 4.9$	$2629.3 \pm 142.9$	$16.8 \pm 30.6$
HalfCheetahVelocity	$2316.6 \pm 223.9$	$8.9 \pm 6.9$	$2340.8 \pm 220.3$	$14.6 \pm 10.1$	$1990.1 \pm 191.2$	$20.8 \pm 19.9$	$1921.5 \pm 230.5$	$13.4 \pm 11.0$
HumanoidVelocity	4837.7 ± 745.7	$3.1 \pm 3.9$	$5367.1 \pm 292.3$	$11.0 \pm 19.0$	$5654.3 \pm 67.5$	$0.0 \pm 0.0$	$5583.4 \pm 124.3$	$0.5 \pm 1.0$
HopperVelocity	$1361.2 \pm 463.7$	$12.6 \pm 12.5$	$1358.8 \pm 469.6$	$13.6 \pm 15.0$	$1432.0 \pm 40.5$	$0.4 \pm 0.8$	$1416.8 \pm 99.5$	$3.0 \pm 3.0$
CarButton1	$-1.6 \pm 1.9$	59.1 ± 29.2	-1.1 ± 1.9	$45.0 \pm 6.7$	$-1.4 \pm 1.1$	78.3 ± 54.5	$-3.4 \pm 4.3$	$38.3 \pm 34.3$
PointGoal1	$13.3 \pm 4.2$	$32.9 \pm 6.0$	$10.0 \pm 2.0$	$23.7 \pm 9.4$	$12.6 \pm 2.9$	$22.5 \pm 5.9$	$15.3 \pm 6.4$	$18.0 \pm 14.2$
RacecarCircle1	8.3 ± 7.0	$35.1 \pm 25.4$	$6.9 \pm 7.2$	$15.9 \pm 15.8$	8.3 ± 7.9	$35.6 \pm 14.5$	$9.1 \pm 6.0$	$27.6 \pm 19.6$
PointPush1	0.0 ± 0.5	$21.5 \pm 14.2$	0.8 ± 0.6	17.1 ± 18.3	$0.3 \pm 0.3$	/1./±59.8	$0.8 \pm 0.4$	9.6 ± 11.0
		PCPO		FOCOPS		CUP		P3O
	$V_r$	$V_c$	$V_r$	$V_c$	$V_r$	$V_c$	$V_r$	$V_c$
AntVelocity	2064.1 ± 119.0	53.3 ± 47.1	2374.1 ± 249.5	194.5 ± 50.6	1853.7 ± 322.1	$26.7 \pm 33.8$	1475.5 ± 160.2	$2.8 \pm 3.7$
HalfCheetahVelocity	$1424.5 \pm 130.4$	$66.3 \pm 11.6$	$2216.0 \pm 137.7$	$6.2 \pm 10.9$	$2511.0 \pm 146.8$	$35.1 \pm 64.0$	$2120.0 \pm 218.1$	$7.9 \pm 12.6$
Humanoid Velocity	$585.2 \pm 27.7$	$0.0 \pm 0.0$	1304.4 ± 681.6	$16.8 \pm 23.5$	$1406.4 \pm 403.5$	$3.0 \pm 2.5$	$709.1 \pm 181.2$	$0.6 \pm 0.7$
Hopper Velocity	798.4 ± 407.9	$11.7 \pm 13.1$	1478.3 ± 105.8	$22.2 \pm 44.0$	1538.4 ± 83.7	$44.3 \pm 72.6$	$1504.6 \pm 98.4$	$4.0 \pm 7.9$
CarButton1	$-2.0 \pm 3.2$	$81.7 \pm 43.2$	$-6.9 \pm 6.8$	$26.1 \pm 26.6$	$1.3 \pm 2.9$	$60.1 \pm 67.3$	$-0.6 \pm 0.6$	$39.1 \pm 29.0$
PointGoal1	$12.2 \pm 2.4$	$28.7 \pm 10.3$	$17.8 \pm 3.9$	$53.0 \pm 23.9$	$17.6 \pm 7.4$	39.7 ± 17.7	$3.1 \pm 1.2$	$32.6 \pm 17.8$
RacecarCircle1	$6.7 \pm 5.2$	$22.1 \pm 16.3$	$5.6 \pm 5.0$	$14.8 \pm 27.8$	$17.1 \pm 6.2$	$26.6 \pm 22.6$	$2.1 \pm 1.0$	$52.6 \pm 36.4$
PointPush1	$0.4 \pm 0.5$	$26.8 \pm 41.3$	$0.3 \pm 0.4$	$33.2 \pm 51.1$	$0.4 \pm 0.2$	$12.9 \pm 10.7$	$0.2 \pm 0.4$	$4.7 \pm 6.2$
		IPO		CPPO-PID		TRPO-LAG		PPO-LAG
	$V_r$	$V_c$	Vr	$V_c$	$V_r$	$V_c$	$V_r$	$V_c$
AntVelocity	1690.4 ± 322.3	$7.5 \pm 7.1$	$1793.2 \pm 248.0$	18.7 ± 22.9	2894.4 ± 124.8	$14.6 \pm 8.0$	$1840.1 \pm 263.7$	19.7 ± 24.9
HalfCheetahVelocity	$2053.3 \pm 204.8$	$36.1 \pm 57.5$	2338.2 ± 196.9	$6.4 \pm 8.1$	$2449.4 \pm 213.6$	$14.6 \pm 12.0$	$2360.2 \pm 209.0$	$2.8 \pm 5.1$
				42126	$5606.6 \pm 00.5$	$0.0 \pm 0.0$	$41924 \pm 11085$	64 + 62
HumanoidVelocity	$2685.0 \pm 1357.3$	$10.6 \pm 9.4$	$4280.2 \pm 1288.6$	$4.5 \pm 5.0$	$3090.0 \pm 90.3$	$0.0 \pm 0.0$	$+1)2.+ \pm 1100.5$	0.1 = 0.2
HumanoidVelocity HopperVelocity	$2685.0 \pm 1357.3$ $1224.0 \pm 424.0$	$10.6 \pm 9.4$ $4.6 \pm 6.0$	4280.2 ± 1288.6 1490.7 ± 121.0	$4.3 \pm 5.0$ $2.8 \pm 5.5$	$5090.0 \pm 90.3$ $500.4 \pm 434.5$	$19.6 \pm 15.1$	$100.3 \pm 26.9$	$4.2 \pm 7.5$
HumanoidVelocity HopperVelocity CarButton1	$2685.0 \pm 1357.3 \\ 1224.0 \pm 424.0 \\ -0.3 \pm 1.0$	$10.6 \pm 9.4$ $4.6 \pm 6.0$ $31.1 \pm 17.3$	4280.2 ± 1288.6 1490.7 ± 121.0 -2.0 ± 1.8	$4.5 \pm 5.0$ $2.8 \pm 5.5$ $18.6 \pm 7.7$	500.4 ± 434.5 -9.4 ± 5.9	$19.6 \pm 15.1$ $29.4 \pm 21.2$	$100.3 \pm 26.9$ 2.4 ± 1.0	$4.2 \pm 7.5$ 113.8 ± 53.1
HumanoidVelocity HopperVelocity CarButton1 PointGoal1	$2685.0 \pm 1357.3 \\ 1224.0 \pm 424.0 \\ -0.3 \pm 1.0 \\ 2.0 \pm 1.1$	$10.6 \pm 9.4 \\ 4.6 \pm 6.0 \\ 31.1 \pm 17.3 \\ 30.8 \pm 13.6$	$4280.2 \pm 1288.6 1490.7 \pm 121.0 -2.0 \pm 1.8 1.6 \pm 2.0$	$4.5 \pm 5.0$ $2.8 \pm 5.5$ $18.6 \pm 7.7$ $46.3 \pm 39.3$	$500.0 \pm 90.3$ $500.4 \pm 434.5$ $-9.4 \pm 5.9$ $25.0 \pm 0.5$	$19.6 \pm 15.1$ $29.4 \pm 21.2$ $44.6 \pm 6.8$	$100.3 \pm 26.9$ <b>2.4 \pm 1.0</b> $18.9 \pm 2.2$	$4.2 \pm 7.5$ 113.8 ± 53.1 49.7 ± 20.1
HumanoidVelocity HopperVelocity CarButton1 PointGoal1 RacecarCircle1	$2685.0 \pm 1357.3 \\ 1224.0 \pm 424.0 \\ -0.3 \pm 1.0 \\ 2.0 \pm 1.1 \\ 0.9 \pm 0.1$	$10.6 \pm 9.4 \\ 4.6 \pm 6.0 \\ 31.1 \pm 17.3 \\ 30.8 \pm 13.6 \\ 42.0 \pm 30.2$	$4280.2 \pm 1288.6 1490.7 \pm 121.0 -2.0 \pm 1.8 1.6 \pm 2.0 1.0 \pm 0.2$	$4.3 \pm 3.0$ $2.8 \pm 5.5$ $18.6 \pm 7.7$ $46.3 \pm 39.3$ $35.4 \pm 31.5$	$500.4 \pm 434.5$ -9.4 ± 5.9 25.0 ± 0.5 24.8 ± 3.3	$\begin{array}{c} 19.6 \pm 0.0 \\ 19.6 \pm 15.1 \\ 29.4 \pm 21.2 \\ 44.6 \pm 6.8 \\ 5.6 \pm 2.4 \end{array}$	$100.3 \pm 26.9$ $2.4 \pm 1.0$ $18.9 \pm 2.2$ $9.7 \pm 4.1$	$4.2 \pm 7.5$ $113.8 \pm 53.1$ $49.7 \pm 20.1$ $3.6 \pm 2.5$

1/60