# Learning to Score Behaviors for Guided Policy Optimization

**Aldo Pacchiano**[*]
UC Berkeley

**Jack Parker-Holder**[*]
University of Oxford

**Yunhao Tang**[*]
Columbia University

**Anna Choromanska**
NYU

**Krzysztof Choromanski**
Google Brain Robotics

**Michael I. Jordan**
UC Berkeley

## Abstract

We introduce a new approach for comparing reinforcement learning policies, using Wasserstein distances (WDs) in a newly defined latent behavioral space. We show that by utilizing the dual formulation of the WD, we can learn score functions over policy behaviors that can in turn be used to lead policy optimization towards (or away from) (un)desired behaviors. Combined with smoothed WDs, the dual formulation allows us to devise efficient algorithms that take stochastic gradient descent steps through WD regularizers. We incorporate these regularizers into two novel on-policy algorithms, Behavior-Guided Policy Gradient and Behavior-Guided Evolution Strategies, which we demonstrate can outperform existing methods in a variety of challenging environments. We also provide an open source demo[*].

## 1 Introduction

One of the key challenges in reinforcement learning (RL) is to efficiently incorporate the behaviors of learned policies into optimization algorithms [17, 21, 7]. The fundamental question we aim to shed light on in this paper is:

*What is the right measure of similarity between two policies acting on the same underlying MDP and how can we devise algorithms to leverage this information for RL?*

In simple terms, the main thesis motivating the methods we propose is that:

*Two policies may perform similar actions at a local level but result in very different global behaviors.*

We propose to define *behaviors* via so-called Behavioral Policy Embeddings (henceforth referred to as Policy Embeddings), which can be both on policy and off policy.

On policy embeddings are achieved via what we call Behavioral Embeddings Maps (BEMs) - functions mapping trajectories of a policy into a latent behavioral space representing trajectories in a compact way. We define the policy embedding as the pushforward distributions over trajectory embeddings as a result of applying a BEM to the policy's trajectories. Importantly, two policies with distinct distributions over trajectories may result in the same probabilistic embedding. Off policy

---

[*]Equal contribution.

[*]Available at `https://github.com/behaviorguidedRL/BGRL`. We emphasize this is the exact code from our experiments, but a demo to build intuition and clarify our methods.

embeddings in contrast correspond to state and policy evaluation pairs resulting of evaluating the policy on states sampled from a probing state distribution that can be chosen independently from the policy.

Both embedding mechanisms result in probabilistic Policy Embeddings, which allow us to identify a policy with a distribution with support on an embedding space. Policy Embeddings provide us a way to rigorously define dissimilarity between policies. We do this by equipping them with metrics defined on the manifold of probability measures, namely a class of Wasserstein distances (WDs, [35]). There are several reasons for choosing WDs:

- **Flexibility**. We can use any cost function between embeddings of trajectories, allowing the distance between policy embeddings to arise organically from an interpretable distance between embedding points.

- **Non-injective BEMs**. Different trajectories may be mapped to the same embedding point (for example in the case of the last-state embedding). This precludes the use of likelihood-based distances such as the KL divergence [16], which we discuss in Section 6.

- **Behavioral Test Functions**. Solving the dual formulation of the WD objective yields a pair of test functions over the space of embeddings, used to score trajectories or state policy pairs (see: Sec. 5.2).

The Behavioral Test Functions, underpin all our algorithms, directing optimization towards desired behaviors. To learn them, it suffices to define the embedding type and BEM (if required) and the cost function between points in the resulting behavioral manifold. To mitigate the computational burden of computing WDs, we rely on their entropy-regularized formulations. This allows us to update the learned test functions in a computationally efficient manner via stochastic gradient descent (SGD) on a Reproducing Kernel Hilbert Space (RKHS). We develop a novel method for stochastic optimal transport based on random feature maps [26] to produce compact and memory-efficient representations of learned behavioral test functions. Finally, having laid the groundwork for comparing policies via behavior-driven trajectory or state-policy pairs scores, we address our core question by introducing two new on-policy RL algorithms:

- **Behavior Guided Policy Gradients (BGPG)**: We propose to replace the KL-based trust region from [31] with a WD-based in the behavior space.

- **Behavior Guided Evolution Strategies (BGES)**: BGES improves on Novelty Search [7] by jointly optimizing for reward and *novelty* using the WD in the behavior space.

We also demonstrate a way to harness our methodology for imitation and repulsion learning (Section 5.2), showing the universality of the proposed techniques.

## 2 Motivating Behavior-Guided Reinforcement Learning

Throughout this paper we prompt the reader to think of a policy as a distribution over its behaviors, induced by the policy's (possibly stochastic) map from state to actions and the unknown environment dynamics. We care about summarizing (or embedding) behaviors into succinct representations that can be compared with each other (via a cost/metric). These comparisons arise naturally when answering questions such as: Has a given trajectory achieved a certain level of reward? Has it visited a certain part of the state space? We think of these summaries or embeddings as characterizing the behavior of the trajectory or relevant state policy-pairs. We formalize these notions in Section 3.

We show that by identifying policies with the embedding distributions that result of applying the embedding function (summary) to their trajectories, and combining this with the provided cost metric, we can induce a topology over the space of policies given by the Wasserstein distance over their embedding distributions. The methods we propose can be thought of as ways to leverage this "behavior" geometry for a variety of downstream applications such as policy optimization and imitation learning.

This topology emerges naturally from the sole definition of an embedding map (behavioral summary) and a cost function. Crucially these choices occur in the semantic space of behaviors as opposed to parameters or visitation frequencies[*]. One of the advantages of choosing a Wasserstein geometry is

---

[*]If we choose an appropriate embedding map our framework handles visitation frequencies as well.

that non-surjective trajectory embedding maps are allowed. This is not possible with a KL induced one (in non-surjective cases, computing the likelihood ratios in the KL definition is in general intractable). In Sections 4 and 5 we show that in order to get a handle on this geometry, we can use the dual formulation of the Wasserstein distance to learn functions (Behavioral Test Functions) that can provide scores on trajectories which then can be added to the reward signal (in policy optimization) or used as a reward (in Imitation Learning).

In summary, by defining an embedding map of trajectories into a behavior embedding space equipped with a metric[*], our framework allows us to learn "reward" signals (Behavioral Test Functions) that can serve to steer policy search algorithms through the "behavior geometry" either in conjunction with a task specific reward (policy optimization) or on their own (e.g. Imitation Learning). We develop versions of on policy RL algorithms which we call Behavior Guided Policy Gradient (BGPG) and Behavior Guided Evolution Strategies (BGES) that enhance their baseline versions by the use of learned Behavioral Test Functions. Our experiments in Section 7 show this modification is useful. We also provide a simple example for repulsion learning and Imitation Learning, where we only need access to an expert's embedding. Our framework also has obvious applications to safety, learning policies that avoid undesirable behaviors.

A final important note is that in this work we only consider simple heuristics for the embeddings, as used in the existing literature. For BGES, these embeddings are those typically used in Quality Diversity algorithms [24], while for BGPG we reinterpret the action distribution currently used in KL-based trust regions [32, 31]. We emphasize the focus of this paper is on introducing the framework to score these behaviors to guide policy optimization.

## 3  Defining Behavior in Reinforcement Learning

A Markov Decision Process (MDP) is a tuple $(\mathcal{S}, \mathcal{A}, \mathrm{P}, \mathrm{R})$. Here $\mathcal{S}$ and $\mathcal{A}$ stand for the sets of states and actions respectively, such that for $s, s' \in \mathcal{S}$ and $a \in \mathcal{A}$: $\mathrm{P}(s'|a, s)$ is the probability that the system/agent transitions from $s$ to $s'$ given action $a$ and $\mathrm{R}(s', a, s)$ is a reward obtained by an agent transitioning from $s$ to $s'$ via $a$. A policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ is a (possibly randomized) mapping (parameterized by $\theta \in \mathbb{R}^d$) from $\mathcal{S}$ to $\mathcal{A}$. Let $\Gamma = \{\tau = s_0, a_0, r_0, \cdots s_H, a_H, r_H \text{ s.t. } s_i \in \mathcal{S}, a_i \in \mathcal{A}, r_i \in \mathbb{R}\}$ be the set of possible trajectories enriched by sequences of partial rewards under some policy $\pi$. The undiscounted reward function $\mathcal{R} : \Gamma \rightarrow \mathbb{R}$ (which expectation is to be maximized by optimizing $\theta$) satisfies $\mathcal{R}(\tau) = \sum_{i=0}^{H} r_i$, where $r_i = R(s_{i+1}, a_i, s_i)$.

### 3.1  Behavioral Embeddings

In this work we identify a policy with what we call a Policy Embedding. We focus on two types of Policy Embeddings both of which are probabilistic in nature, on policy and off policy embeddings, the first being trajectory based and the second ones state-based.

#### 3.1.1  On Policy Embeddings

We start with a Behavioral Embedding Map (BEM), $\Phi : \Gamma \rightarrow \mathcal{E}$, mapping trajectories to embeddings (Fig. 1), where $\mathcal{E}$ can be seen as a behavioral manifold. On Policy Embeddings can be for example: a) **State-Based**, such as the final state $\Phi_1(\tau) = s_H$ b) **Action-based:** such as the concatenation of actions $\Phi_4(\tau) = [a_0, ..., a_H]$ or c) **Reward-based:** the total reward $\Phi_5(\tau) = \sum_{t=0}^{H} r_t$, reward-to-go vector $\Phi_6(\tau) = \sum_{t=0}^{H} r_t \left( \sum_{i=0}^{t} e_i \right)$ (where $e_i \in \mathbb{R}^{H+1}$ is a one-hot vector corresponding to $i$ with dimension index from $0$ to $H$). Importantly, the mapping does not need to be surjective, as we see on the example of the final state embedding.

---

[*]The embedding space can be discrete or continuous and the metric need not be smooth, and can be for example a simple discrete $\{0, 1\}$ valued criterion
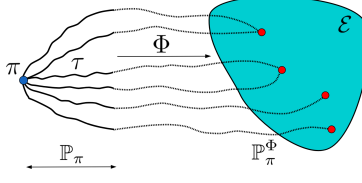
Figure 1: **Behavioral Embedding Maps (BEMs)** map trajectories to points in the behavior embedding space $\mathcal{E}$. Two trajectories may map to the same point in $\mathcal{E}$.

Given a policy $\pi$, we let $\mathbb{P}_\pi$ denote the distribution induced over the space of trajectories $\Gamma$ and by $\mathbb{P}_\pi^\Phi$ the corresponding pushforward distribution on $\mathcal{E}$ induced by $\Phi$. We call $P_\pi^\Phi$ the policy embeddings of a policy $\pi$. A policy $\pi$ can be fully characterized by the distribution $\mathbb{P}_\pi$ (see: Fig. 1).

Additionally, we require $\mathcal{E}$ to be equipped with a metric (or cost function) $C : \mathcal{E} \times \mathcal{E} \to \mathbb{R}$. Given two trajectories $\tau_1, \tau_2$ in $\Gamma$, $C(\Phi(\tau_1), \Phi(\tau_2))$ measures how different these trajectories are in the behavior space. We note that some embeddings are only for the tabular case ($|\mathcal{S}|, |\mathcal{A}| < \infty$) while others are universal.

### 3.1.2  Off Policy Embeddings

Let $\mathbb{P}_\mathcal{S}$ be some "probe" distribution over states $\mathcal{S}$ and $\pi$ be a policy. We define $\mathbb{P}_\pi^{\Phi_\mathcal{S}}$ to be the distribution of pairs $(s, \pi(s))$ for $s \sim \mathbb{P}_S$. We identify $\mathcal{E}$ with the product space $\mathcal{S} \times \Delta_\mathcal{A}$ (where $\Delta_\mathcal{A}$ denotes the set of distributions over $\mathcal{A}$) endowed with an appropriate metric $C : \mathcal{E} \times \mathcal{E} \to \mathbb{R}$. In our experiments we identify $\mathcal{C}$ with the $l_2$ norm over $\mathcal{E}$ and $\mathbb{P}_\mathcal{S}$ with a mechanism that samples states from a buffer of visited states. We only add an $\mathcal{S}$ to the notation for $\mathbb{P}_\pi^{\Phi_\mathcal{S}}$ when distinguishing from on-policy embeddings is needed.

This definition allows the "probing" distribution $\mathbb{P}_\mathcal{S}$ to be off policy, independent of the policy at hand. If $C$ is a norm and $\mathbb{P}_\mathcal{S}$ has mass only in user-relevant areas of the state space, a WD of zero between two policies (whose embeddings use the same probing distribution) implies they behave equally where the user cares. Our off Policy Embeddings are of the form $(s, \pi(s))$ but other choices are valid.

## 4  Wasserstein Distance & Optimal Transport Problem

Let $\mu, \nu$ be (Radon) probability measures over domains $\mathcal{X} \subseteq \mathbb{R}^m, \mathcal{Y} \subseteq \mathbb{R}^n$ and let $\mathcal{C} : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ be a cost function. For $\gamma > 0$, a *smoothed* Wasserstein Distance is defined as:

$$\mathrm{WD}_\gamma(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} C(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y}) + \Sigma, \tag{1}$$

where $\Sigma = \gamma \mathrm{KL}(\pi | \xi)$, $\Pi(\mu, \nu)$ is the space of couplings (joint distributions) over $\mathcal{X} \times \mathcal{Y}$ with marginal distributions $\mu$ and $\nu$, $\mathrm{KL}(\cdot | \cdot)$ denotes the KL divergence between distributions $\pi$ and $\rho$ with support $\mathcal{X} \times \mathcal{Y}$ defined as: $\mathrm{KL}(\pi | \rho) = \int_{\mathcal{X} \times \mathcal{Y}} \left( \log \left( \frac{d\pi}{d\xi}(\mathbf{x}, \mathbf{y}) \right) \right) d\pi(\mathbf{x}, \mathbf{y})$ and $\xi$ is a reference measure over $\mathcal{X} \times \mathcal{Y}$. When the cost is an $\ell_p$ distance and $\gamma = 0$, $\mathrm{WD}_\gamma$ is also known as the Earth mover's distance and the corresponding optimization problem is known as the *optimal transport problem* (OTP).

### 4.1  Wasserstein Distance: Dual Formulation

We will use smoothed WDs to derive efficient regularizers for RL algorithms. To arrive at this goal, we first need to consider the dual form of Equation 1. Under the subspace topology [4] for $\mathcal{X}$ and $\mathcal{Y}$, let $\mathcal{C}(\mathcal{X})$ and $\mathcal{C}(\mathcal{Y})$ denote the space of continuous functions over $\mathcal{X}$ and $\mathcal{Y}$ respectively. The choice of the subspace topology ensures our discussion encompasses the discrete case.

Let $C : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ be a cost function, interpreted as the "ground cost" to move a unit of mass from $x$ to $y$. Define $\mathbb{I}$ as the function outputting values of its input predicates. Using Fenchel duality, we can obtain the following dual formulation of the problem in Eq. 1:

$$\mathrm{WD}_\gamma(\mu, \nu) = \max_{\lambda_\mu \in \mathcal{C}(\mathcal{X}), \lambda_\nu \in \mathcal{C}(\mathcal{Y})} \Psi(\lambda_\mu, \lambda_\nu), \tag{2}$$

**Algorithm 1** Random Features Wasserstein SGD
___
**Input:** kernels $\kappa, \ell$ over $\mathcal{X}, \mathcal{Y}$ respectively with corresponding random feature maps $\phi_\kappa, \phi_\ell$, smoothing parameter $\gamma$, gradient step size $\alpha$, number of optimization rounds $M$, initial dual vectors $\mathbf{p}_0^\mu, \mathbf{p}_0^\nu$.
**for** $t = 0, \cdots, M$ **do**
> 1. Sample $(x_t, y_t) \sim \mu \bigotimes \nu$.
> 2. Update: $\binom{\mathbf{p}_t^\mu}{\mathbf{p}_t^\nu}$ using Equation 4.

**Return:** $\mathbf{p}_M^\mu, \mathbf{p}_M^\nu$.
___

where $\Psi(\lambda_\mu, \lambda_\nu) = \int_{\mathcal{X}} \lambda_\mu(\mathbf{x}) d\mu(\mathbf{x}) - \int_{\mathcal{Y}} \lambda_\nu(\mathbf{y}) d\nu(\mathbf{y}) - E_C(\lambda_\mu, \lambda_\nu)$ and the damping term $E_C(\lambda_\mu, \lambda_\nu)$ equals:

$$E_C(\lambda_\mu, \lambda_\nu) = \mathbb{I}(\gamma > 0) \int_{\mathcal{X} \times \mathcal{Y}} \rho(\mathbf{x}, \mathbf{y}) d\xi(\mathbf{x}, \mathbf{y}) + \mathbb{I}(\gamma = 0) \mathbb{I}(\mathcal{A}) \tag{3}$$

for $\rho(\mathbf{x}, \mathbf{y}) = \gamma \exp(\frac{\lambda_\mu(\mathbf{x}) - \lambda_\nu(\mathbf{y}) - C(\mathbf{x}, \mathbf{y})}{\gamma})$ and $\mathcal{A} = [(\lambda_\mu, \lambda_\nu) \in \{(u, v) \text{ s.t. } \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y} : u(\mathbf{x}) - v(\mathbf{y}) \leq C(\mathbf{x}, \mathbf{y})\}]$.

We will set the damping distribution $d\xi(\mathbf{x}, \mathbf{y}) \propto 1$ for discrete domains and $d\xi(\mathbf{x}, \mathbf{y}) = d\mu(\mathbf{x}) d\nu(\mathbf{y})$ otherwise.

If $\lambda_\mu^*, \lambda_\nu^*$ are the functions achieving the maximum in Eq. 2, and $\gamma$ is sufficiently small then $\mathrm{WD}_\gamma(\mu, \nu) \approx \mathbb{E}_\mu [\lambda_\mu^*(\mathbf{x})] - \mathbb{E}_\nu [\lambda_\nu^*(\mathbf{y})]$, with equality when $\gamma = 0$. When for example $\gamma = 0$, $\mathcal{X} = \mathcal{Y}$, and $C(x, x) = 0$ for all $x \in \mathcal{X}$, it is easy to see $\lambda_\mu^*(x) = \lambda_\nu^*(x) = \lambda^*(x)$ for all $x \in \mathcal{X}$. In this case the difference between $\mathbb{E}_\mu [\lambda^*(\mathbf{x})]$ and $\mathbb{E}_\mu [\lambda^*(\mathbf{y})]$ equals the WD. In other words, the function $\lambda^*$ gives higher scores to regions of the space $\mathcal{X}$ where $\mu$ has more mass. This observation is key to the success of our algorithms in guiding optimization towards desired behaviors.

## 4.2 Computing $\lambda_\mu^*$ and $\lambda_\nu^*$

We combine several techniques to make the optimization of objective from Eq. 2 tractable. First, we replace $\mathcal{X}$ and $\mathcal{Y}$ with the functions from a RKHS corresponding to universal kernels [22]. This is justified since those function classes are dense in the set of continuous functions of their ambient spaces. In this paper we choose the RBF kernel and approximate it using random Fourier feature maps [26] to increase efficiency. Consequently, the functions $\lambda$ learned by our algorithms have the following form: $\lambda(\mathbf{x}) = (\mathbf{p}^\lambda)^\top \phi(\mathbf{x})$, where $\phi$ is a random feature map with $m$ standing for the number of random features and $\mathbf{p}^\lambda \in \mathbb{R}^m$. For the RBF kernel, $\phi$ is defined as follows: $\phi(\mathbf{z}) = \frac{1}{\sqrt{m}} \cos(\mathbf{G}\mathbf{z} + \mathbf{b})$ for $\mathbf{z} \in \mathbb{R}^d$, where $\mathbf{G} \in \mathbb{R}^{m \times d}$ is Gaussian with iid entries taken from $\mathcal{N}(0, 1), b \in \mathbb{R}^m$ with iid $b_i$s such that $b_i \sim \mathrm{Unif}[0, 2\pi]$ and the cos function acts elementwise.
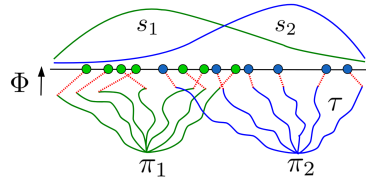


Figure 2: Two policies $\pi_1$ (green) and $\pi_2$ (blue) whose BEMs map trajectories to points in the real line.

Henceforth, when we refer to optimization over $\lambda$, we mean optimizing over corresponding dual vectors $\mathbf{p}^\lambda$ associated with $\lambda$. We can solve for the optimal dual functions by running Stochastic Gradient Descent (SGD) over the dual objective in Eq. 2. Algorithm 1 is the random features equivalent of Algorithm 3 in [12]. Given input kernels $\kappa, \ell$ and a fresh sample $(x_t, y_t) \sim \mu \bigotimes \nu$ the SGD step w.r.t. the current iterates $\mathbf{p}_{t-1}^\mu, \mathbf{p}_{t-1}^\nu$ satisfies:

$$F(\mathbf{p}_1, \mathbf{p}_2, x, y) = \exp \left( \frac{(\mathbf{p}_1)^\top \phi_\kappa(x) - (\mathbf{p}_2)^\top \phi_\ell(x) - C(x, y)}{\gamma} \right)$$

$$\binom{\mathbf{p}_{t+1}^\mu}{\mathbf{p}_{t+1}^\nu} = \binom{\mathbf{p}_t^\mu}{\mathbf{p}_t^\nu} + (1 - F(\mathbf{p}_t^\mu, \mathbf{p}_t^\nu, x_t, y_t)) v_t, \tag{4}$$

where $v_t = \frac{\alpha}{\sqrt{t}}(\phi_\kappa(x_t), -\phi_\ell(y_t))^\top$. An explanation and proof of these formulae is in Lemma C.2 in the Appendix. If $\mathbf{p}^\mu_*, \mathbf{p}^\nu_*$ are the optimal dual vectors, $p_* = (\mathbf{p}^\mu_*, \mathbf{p}^\nu_*)^\top$, $(x_1, y_1), \cdots, (x_k, y_k) \overset{\text{i.i.d}}{\sim} \mu \bigotimes \nu$, $\mathbf{v}^{\kappa,\ell}_i = (\phi_\kappa(x_i), -\phi_\ell(y_i))^\top$ for all $i$, and $\hat{\mathbb{E}}$ denotes the empirical expectation over the $k$ samples $\{(x_i, y_i)\}^k_{i=1}$, Algorithm 1 can be used to get an estimator of $\text{WD}_\gamma(\mu, \nu)$ as:

$$\widehat{\text{WD}}_\gamma(\mu, \nu) = \hat{\mathbb{E}}\left[\langle \mathbf{p}_*, \mathbf{v}^{\kappa,\ell}_i \rangle - \frac{F(\mathbf{p}^\mu_*, \mathbf{p}^\nu_*, x_i, y_i)}{\gamma}\right] \quad (5)$$

## 5  Behavior-Guided Reinforcement Learning

We explain now how to get practical algorithms based on the presented methods. Denote by $\pi_\theta$ a policy parameterized by $\theta \in \mathbb{R}^d$. The goal of policy optimization algorithms is to find a policy maximizing, as a function of the policy parameters, the expected total reward $\mathcal{L}(\theta) := \mathbb{E}_{\tau \sim \mathbb{P}_{\pi_\theta}}[\mathcal{R}(\tau)]$.

### 5.1  Behavioral Test Functions

If $C : \mathcal{E} \times \mathcal{E} \to \mathbb{R}$ is a cost function defined over behavior space $\mathcal{E}$, and $\pi_1, \pi_2$ are two policies, then in the case of **On-Policy** Embeddings:

$$\text{WD}_\gamma(\mathbb{P}^\Phi_{\pi_1}, \mathbb{P}^\Phi_{\pi_2}) \approx \mathbb{E}_{\tau \sim \mathbb{P}_{\pi_1}}[\lambda^*_1(\Phi(\tau))] - \mathbb{E}_{\tau \sim \mathbb{P}_{\pi_2}}[\lambda^*_2(\Phi(\tau))],$$

where $\lambda^*_1, \lambda^*_2$ are the optimal dual functions. The maps $s_1 := \lambda^*_1 \circ \Phi : \Gamma \to \mathbb{R}$ and $s_2 := \lambda^*_2 \circ \Phi : \Gamma \to \mathbb{R}$ define score functions over the space of trajectories. If $\gamma$ is close to zero, the score function $s_i$ gives higher scores to trajectories from $\pi_i$ whose behavioral embedding is common under $\pi_i$ but rarely appears under $\pi_j$ for $j \neq i$ (Fig. 2). In the case of **Off-Policy** Embeddings:

$$\text{WD}_\gamma(\mathbb{P}^{\Phi_S}_{\pi_1}, \mathbb{P}^{\Phi'_S}_{\pi_2}) \approx \mathbb{E}_{S \sim \mathbb{P}_S}[\lambda^*_1(S, \pi_1(S))] - \mathbb{E}_{S \sim \mathbb{P}'_S}[\lambda^*_2(S, \pi_2(S))],$$

where $\lambda^*_1, \lambda^*_2$ are maps from state policy pairs $(S, \pi_1(S))$ to scores, and $\mathbb{P}_S, \mathbb{P}'_S$ are probing distributions.

### 5.2  Repulsion and Imitation Learning

To illustrate the intuition behind behavioral test functions and on policy embeddings, we introduce an algorithm for multi-policy repulsion learning based on our framework. Algorithm 2 maintains two policies $\pi^{\mathbf{a}}$ and $\pi^{\mathbf{b}}$.

---
**Algorithm 2** Behvaior-Guided Repulsion Learning

---
**Input:** $\beta, \eta > 0$, $M \in \mathbb{N}$
**Initialize:** Initial stochastic policies $\pi^{\mathbf{a}}_0, \pi^{\mathbf{b}}_0$, parametrized by $\theta^{\mathbf{a}}_0, \theta^{\mathbf{b}}_0$ respectively, Behavioral Test Functions $\lambda^{\mathbf{a}}_1, \lambda^{\mathbf{b}}_2$
**for** $t = 1, \ldots, T$ **do**
    1. Collect $\{\tau^{\mathbf{a}}_i\}^M_{i=1} \sim \mathbb{P}_{\pi^{\mathbf{a}}_{t-1}}$ and $\{\tau^{\mathbf{b}}_i\}^M_{i=1} \sim \mathbb{P}_{\pi^{\mathbf{b}}_{t-1}}$.
    2. Form $\tilde{R}_{\mathbf{c}}(\tau_1, \tau_2)$ for $\mathbf{c} \in \{\mathbf{a}, \mathbf{b}\}$ using Equation 6.
    3. For $\mathbf{c} \in \{\mathbf{a}, \mathbf{b}\}$ and $(\tau_1, \tau_2) \sim \{\tau^{\mathbf{a}}_i\}^M_{i=1} \times \{\tau^{\mathbf{b}}_i\}^M_{i=1}$ use REINFORCE [36] to perform update:

$$\theta^{\mathbf{c}}_t = \theta^{\mathbf{c}}_{t-1} + \eta \nabla_\theta \tilde{R}_{\mathbf{c}}(\tau_1, \tau_2)$$

    5. Update $\lambda^{\mathbf{a}}_1, \lambda^{\mathbf{b}}_2$ with $\{\tau^{\mathbf{a}}_i, \tau^{\mathbf{b}}_i\}^M_{i=1}$ via Algorithm 1.

---

Each policy is optimized by taking a policy gradient step (using the REINFORCE gradient estimator [36]) to optimize surrogate rewards $\tilde{\mathcal{R}}_{\mathbf{a}}$ and $\tilde{\mathcal{R}}_{\mathbf{b}}$.
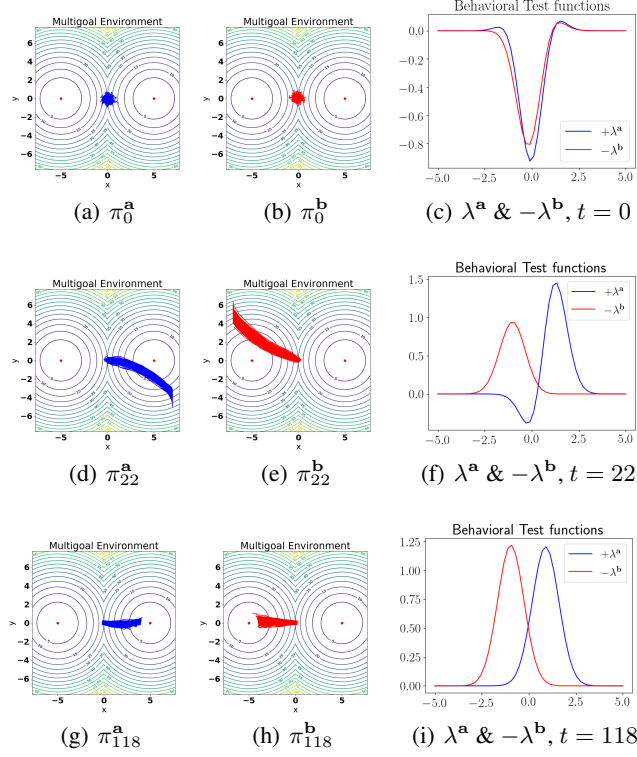
Figure 3: a) and b) Initial state of policies $\pi^{\mathbf{a}}, \pi^{\mathbf{b}}$ and Test functions $\lambda^{\mathbf{a}}, \lambda^{\mathbf{b}}$. d)-i) Policy evolution and Test Functions.

These combine the signal from the task's reward function $\mathcal{R}$ and the repulsion score encoded by the input BEM $\Phi$ and behavioral test functions $\lambda^{\mathbf{a}}$ and $\lambda^{\mathbf{b}}$:

$$\tilde{R}_{\mathbf{c}}(\tau_{\mathbf{a}}, \tau_{\mathbf{b}}) = \mathcal{R}(\tau_{\mathbf{c}}) + \beta \widehat{\mathrm{WD}}_\gamma(\mathbb{P}^{\Phi}_{\pi^{\mathbf{a}}}, \mathbb{P}^{\Phi}_{\pi^{\mathbf{a}}}), \mathbf{c} \in \{\mathbf{a}, \mathbf{b}\} \tag{6}$$

We test Algorithm 2 on an environment consisting of a particle that needs to reach one of two goals on the plane. Policies outputs a velocity vector and stochasticity is achieved by adding Gaussian noise to it. The embedding $\Phi$ maps trajectories $\tau$ to their mean displacement along the $x-$axis. Fig. 3 shows how the policies' behavior evolves throughout optimization and how the Test Functions guide the optimization by favouring the two policies to be far apart. The experiment details are in the Appendix (Section B.4). A related guided trajectory scoring approach to imitation learning is explored in Appendix B.3.

## 5.3 Algorithms

We propose to solve a WD-regularized objective to tackle behavior-guided policy optimization. All of our algorithms hinge on trying to maximize an objective of the form:

$$F(\theta) = \mathcal{L}(\theta) + \beta \mathrm{WD}_\gamma(\mathbb{P}^{\Phi}_{\pi_\theta}, \mathbb{P}^{\Phi}_{\mathbf{b}}), \tag{7}$$

where $\mathbb{P}^{\Phi}_{\mathbf{b}}$ is a base distribution[*] over behavioral embeddings (possibly dependent on $\theta$) and $\beta \in \mathbb{R}$ could be positive or negative. Although the base distribution $\mathbb{P}^{\Phi}_{\mathbf{b}}$ could be arbitrary, our algorithms will instantiate $\mathbb{P}^{\Phi}_{\mathbf{b}} = \frac{1}{|\mathcal{S}|} \cup_{\pi' \in \mathcal{S}} \mathbb{P}^{\Phi}_{\pi'}$ for some family of policies $\mathcal{S}$ (possibly satisfying $|\mathcal{S}| = 1$) we want the optimization to attract to / repel from.

In order to compute approximate gradients for $F$, we rely on the dual formulation of the WD. After substituting the composition maps resulting from Eq. 5.1 into Eq. 7, we obtain, for **on-policy** embeddings:

$$F(\theta) \approx \mathbb{E}_{\tau \sim \mathbb{P}_{\pi_\theta}} \left[ \mathcal{R}(\tau) + \beta s_1(\tau) \right] - \beta \mathbb{E}_{\phi \sim \mathbb{P}^{\Phi}_{\mathbf{b}}} \left[ \lambda_2^*(\phi) \right], \tag{8}$$

---

[*]Possibly using off policy embeddings.

7

where $s_1 : \Gamma \to \mathbb{R}$ equals $s_1 = \lambda_1^* \circ \Phi$, the Behavioral Test Function of policy $\pi_\theta$ and $\lambda_2^*$ is the optimal dual function of embedding distribution $\mathbb{P}_b^\Phi$. Consequently $\nabla_\theta F(\theta) \approx \nabla_\theta \mathbb{E}_{\tau \sim \mathbb{P}_{\pi_\theta}} [\mathcal{R}(\tau) + \beta s_1(\tau)]$. We learn a score function $s_1$ over trajectories that can guide our optimization by favoring those trajectories that show desired global behaviors. For **off-policy** embeddings, with state probing distributions $\mathbb{P}_S$ and $\mathbb{P}_S^b$ the analogous to Equation 9 is:

$$F(\theta) \approx \mathbb{E}_{\tau \sim \mathbb{P}_{\pi_\theta}} [\mathcal{R}(\tau)] + \beta \mathbb{E}_{(S, \pi_\theta(S)) \sim \mathbb{P}^{\Phi_S}} [\lambda_1^*(S, \pi_\theta(S))]$$
$$- \beta \mathbb{E}_{(S, \pi_b(S)) \sim \mathbb{P}_b^{\Phi_S}} [\lambda_2^*(S, \pi_b(S))], \tag{9}$$

Consequently, if $\mathbb{P}_b^{\Phi_S}$ is independent from $\theta$:

$$\nabla_\theta F(\theta) \approx \nabla_\theta \mathbb{E}_{\tau \sim \mathbb{P}_{\pi_\theta}} [\mathcal{R}(\tau)] + \beta \mathbb{E}_{s \sim \mathbb{P}_S} [\nabla_\theta \lambda_1^*(s, \pi_\theta(s))].$$

Eq. 8 and 9 are approximations to the true objective from Eq. 7 whenever $\gamma > 0$. In practice, the entropy regularization requires a damping term $E_C(\lambda_1^*, \lambda_2^*)$ as defined in Equation 3. If $\xi(\mathbb{P}_{\pi_\theta}^\Phi, \mathbb{P}_b^\Phi)$ is the damping joint distribution of choice and $\rho(\phi_1, \phi_2) = \gamma \exp\left(\frac{\lambda_{\pi_\theta}(\phi_1) - \lambda_b(\phi_2) - C(\phi_1, \phi_2)}{\gamma}\right)$ (for off policy embeddings $\phi$ is a state policy pair $(S, \pi(S))$), the damping term equals: $\mathbb{E}_{\phi_1, \phi_2 \sim \xi(\mathbb{P}_{\pi_\theta}^\Phi, \mathbb{P}_b^\Phi)} [\rho(\phi_1, \phi_2)]$. Gradients $\nabla_\theta$ through $E_C$ can be derived using a similar logic as the gradients above. When the embedding space $\mathcal{E}$ is not discrete and $\mathbb{P}_b^\Phi = \mathbb{P}_\pi^\Phi$ for some policy $\pi$, we let $\xi(\mathbb{P}_{\pi_\theta}^\Phi, \mathbb{P}_b^\Phi) = \mathbb{P}_{\pi_\theta}^\Phi \bigotimes \mathbb{P}_\pi^\Phi$, otherwise $\xi(\mathbb{P}_{\pi_\theta}^\Phi, \mathbb{P}_b^\Phi) = \frac{1}{|\mathcal{E}|^2} \mathbb{1}$, a uniform distribution over $\mathcal{E} \times \mathcal{E}$.

All of our methods perform a version of alternating SGD optimization: we take certain number of SGD steps over the internal dual Wasserstein objective, followed by more SGD steps over the outer objective having fixed the test functions.

We consider two approaches to optimizing this objective. Behavior-Guided Policy Gradient (BGPG) explores in the action space as in policy gradient methods [31, 32], while Behavior-Guided Evolution Strategies (BGES) considers a black-box optimization problem as in Evolution Strategies (ES, [30]).

## 5.4 Behavior-Guided Policy Gradient (BGPG)

Here we present the Behavior-Guided Policy Gradient (BGPG) algorithm (Alg. 3). Specifically, we maintain a stochastic policy $\pi_\theta$ and compute policy gradients as in prior work [31].

---

**Algorithm 3** Behavior-Guided Policy Gradient

---

**Input:** Initialize stochastic policy $\pi_0$ parametrized by $\theta_0$, $\beta < 0, \eta > 0, M \in \mathbb{N}$
**for** $t = 1, \ldots, T$ **do**

    1. Run $\pi_{t-1}$ in the environment to get advantage values $A^{\pi_{t-1}}(s, a)$ and trajectories $\{\tau_i^{(t)}\}_{i=1}^M$
    2. Update policy and test functions via several alternating policy gradient steps over $F(\theta)$.
    3. Use samples from $\mathbb{P}_{\pi_{t-1}} \bigotimes \mathbb{P}_{\pi_\theta}$ and Algorithm 1 to update $\lambda_1, \lambda_2$ and take SGA step $\theta_t = \theta_{t-1} + \eta \hat{\nabla}_\theta \hat{F}(\theta_{t-1})$

---

For **on-policy** embeddings the objective function $F(\theta)$ takes the form:
$$F(\theta) = \mathbb{E}_{\tau_1, \tau_2 \sim \mathbb{P}_{\pi_{t-1}} \bigotimes \mathbb{P}_{\pi_\theta}} \left[ \hat{R}(\tau_1, \tau_2) \right], \tag{10}$$

where $\hat{R}(\tau_1, \tau_2) = \sum A^{\pi_{t-1}}(s_i, a_i) \frac{\pi_\theta(a_i|s_i)}{\pi_{t-1}(a_i|s_i)} + \widehat{\mathrm{WD}}_\gamma(\mathbb{P}_{\pi_{t-1}}^\Phi, \mathbb{P}_{\pi_\theta}^\Phi)$. To optimize the Wasserstein distance we use Algorithm 1. Importantly, stochastic gradients of $F(\theta)$ can be approximated by samples from $\pi_\theta$. In its simplest form, the gradient $\hat{\nabla}_\theta \hat{F}$ can be computed by the vanilla policy gradient over the advantage component and using the REINFORCE estimator through the components involving Test Functions acting on trajectories from $\mathbb{P}_{\pi_\theta}$. For **off-policy** embeddings, $\hat{\nabla}_\theta \hat{F}$ can be computed by sampling from the product of the state probing distributions. Gradients through the differentiable test functions can be computed by the chain rule: $\nabla_\theta \lambda(S, \pi_\theta(S)) = (\nabla_\phi \lambda(\phi))^\top \nabla_\theta \phi$ for $\phi = (S, \pi_\theta(S))$.

BGPG can be thought of as a variant of Trust Region Policy Optimization with a Wasserstein penalty. As opposed to vanilla TRPO, the optimization path of BGPG flows through policy parameter space while encouraging it to follow a smooth trajectory through the geometry of the behavioral manifold. We proceed to show that given the right embedding and cost function, we can prove a monotonic

improvement theorem for BGPG, showing that our methods satisfy at least similar guarantees as TRPO.

Furthermore, Let $V(\pi)$ be the expected reward of policy $\pi$ and $\rho_\pi(s) = \mathbb{E}_{\tau \sim \mathbb{P}_\pi} \left[ \sum_{t=0}^T \mathbf{1}(s_t = s) \right]$ be the visitation measure.

Two distinct policies $\pi$ and $\tilde{\pi}$ can be related via the equation (see: [33]) $V(\tilde{\pi}) = V(\pi) + \int_{\mathcal{S}} \rho_{\tilde{\pi}}(s) \left( \int_{\mathcal{A}} \tilde{\pi}(a|s) A^\pi(s,a) da \right) ds$ and the linear approximations to $V$ around $\pi$ via: $L(\tilde{\pi}) = V(\pi) + \int_{\mathcal{S}} \rho_\pi(s) \left( \int_{\mathcal{A}} \tilde{\pi}(a|s) A^\pi(s,a) da \right) ds$ (see: [14]). Let $\mathcal{S}$ be a finite set. Consider the following embedding $\Phi^s : \Gamma \to \mathbb{R}^{|\mathcal{S}|}$ defined by $(\Phi(\tau))_s = \sum_{t=0}^T \mathbf{1}(s_t = s)$ and related cost function defined as: $C(\mathbf{v}, \mathbf{w}) = \|\mathbf{v} - \mathbf{w}\|_1$. Then $\mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_\pi^{\Phi^s})$ is related to visitation frequencies since $\mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_\pi^{\Phi^s}) \geq \sum_{s \in \mathcal{S}} |\rho_\pi(s) - \rho_{\tilde{\pi}}(s)|$. These observations enable us to prove an analogue of Theorem 1 from [31] (see Section C.2 for the proof), namely:

**Theorem 5.1.** *If* $\mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_\pi^{\Phi^s}) \leq \delta$ *and* $\epsilon = \max_{s,a} |A^\pi(s,a)|$, *then* $V(\tilde{\pi}) \geq L(\tilde{\theta}) - \delta\epsilon$.

As in [31], Theorem 5.1 implies a policy improvement guarantee for BGPG.

### 5.5 Behavior Guided Evolution Strategies (BGES)

ES takes a black-box optimization approach to RL, by considering a rollout of a policy, parameterized by $\theta$ as a black-box function $F$. This approach has gained in popularity recently [30, 20, 6].

---

**Algorithm 4** Behavior-Guided Evolution Strategies

---

**Input:** learning rate $\eta$, noise standard deviation $\sigma$, iterations $T$, BEM $\Phi$, $\beta$ ($> 0$ for repulsion, $< 0$ for imitation).
**Initialize:** Initial policy $\pi_0$ parametrized by $\theta_0$, Behavioral Test Functions $\lambda_1, \lambda_2$. Evaluate policy $\pi_0$ to return trajectory $\tau_0$
**for** $t = 1, \ldots, T-1$ **do**

    1. Sample $\epsilon_1, \cdots, \epsilon_n$ independently from $\mathcal{N}(0, I)$.
    2. Evaluate policies $\{\pi_t^k\}_{k=1}^n$ parameterized by $\{\theta_t + \sigma\epsilon_k\}_{k=1}^n$, get rewards $R_k$ and trajectories $\tau_k$ for all $k$.
    3. Update $\lambda_1$ and $\lambda_2$ using Algorithm 1.
    4. Approximate $\widehat{\mathrm{WD}}\gamma(\mathbb{P}_{\pi_t^k}^\Phi, \mathbb{P}_{\pi_t}^\Phi)$ plugging in $\lambda_1, \lambda_2$ into Eq. 5 for each perturbed policy $\pi_k$
    5. Update Policy: $\theta_{t+1} = \theta_t + \eta \nabla_{ES} F$, where:
$$\nabla_{ES} F = \frac{1}{\sigma} \sum_{k=1}^n [(1-\beta)(R_k - R_t) + \beta \widehat{\mathrm{WD}}\gamma(\mathbb{P}_{\pi_t^k}^\Phi, \mathbb{P}_{\pi_t}^\Phi)] \epsilon_k$$

---

If we take this approach to optimizing the objective in Eq. 7, the result is a black-box optimization algorithm which seeks to maximize the reward and simultaneously maximizes or minimizes the difference in behavior from the base embedding distribution $\mathbb{P}_b^\Phi$. We call it Behavior-Guided Evolution Strategies (BGES) algorithm (see: Alg. 4).

When $\beta > 0$, and we take $\mathbb{P}_b^\Phi = \mathbb{P}_{\pi_{t-1}}^\Phi$, BGES resembles the NSR-ES algorithm from [7], an instantiation of *novelty search* [19]. The positive weight on the WD-term enforces newly constructed policies to be behaviorally different from the previous ones while the $\mathcal{R}-$term drives the optimization to maximize the reward. The key difference in our approach is the probabilistic embedding map, with WD rather than Euclidean distance. We show in Section 7.2 that BGES outperforms NSR-ES for challenging exploration tasks.

## 6 Related Work

Our work is related to research in multiple areas in neuroevolution and machine learning:

**Behavior Characterizations:** The idea of directly optimizing for behavioral diversity was introduced by [19] and [18], who proposed to search directly for *novelty*, rather than simply assuming it would naturally arise in the process of optimizing an objective function. This approach has been applied to deep RL [7] and meta-learning [11]. In all of this work, the policy is represented via a behavioral characterization (BC), which requires domain knowledge. In our setting, we move from

deterministic BCs to stochastic behavioral embeddings, thus requiring the use of metrics capable of comparing probabilistic distributions.

**Distance Metrics:**   WDs have been used in many applications in machine learning where guarantees based on distributional similarity are required [13, 1]. We make use of WDs in our setting for a variety of reasons. First and foremost, the dual formulation of the WD allows us to recover Behavioral Test Functions, providing us with behavior-driven trajectory scores. In contrast to KL divergences, WDs are sensitive to user-defined costs between pairs of samples instead of relying only on likelihood ratios. Furthermore, as opposed to KL divergences, it is possible to take SGD steps using entropy-regularized Wasserstein objectives. Computing an estimator of the KL divergence is hard without a density model. Since in our framework multiple unknown trajectories may map to the same behavioral embedding, the likelihood ratio between two embedding distributions may be ill-defined.

**WDs for RL:**   We are not the first to propose using WDs in RL. [37] have recently introduced Wasserstein Gradient Flows (WGFs), which casts policy optimization as gradient descent flow on the manifold of corresponding probability measures, where geodesic lengths are given as second-order WDs. We note that computing WGFs is a nontrivial task. In [37] this is done via particle approximation methods, which we show in Section 7 is substantially slower than our methods. The WD has also been employed to replace KL terms in standard Trust Region Policy Optimization [27]. This is a very special case of our more generic framework (cf. Section 5.3). In [27] it is suggested to solve the corresponding RL problems via Fokker-Planck equations and diffusion processes, yet no empirical evidence of the feasibility of this approach is provided. We propose general practical algorithms and provide extensive empirical evaluation.

**Distributional RL**   Distributional RL (DRL, [2]) expands on traditional off-policy methods [23] by attempting to learn a distribution of the return from a given state, rather than just the expected value. These approaches have impressive experimental results [2, 8], with a growing body of theory [28, 25, 3, 29]. Superficially it may seem that learning a distribution of returns is similar to our approach to PPEs, when the BEM is a distribution over rewards. Indeed, reward-driven embeddings used in DRL can be thought of as special cases of the general class of BEMs. We note two key differences: 1) DRL methods are off-policy whereas our BGES and BGPG algorithms are on-policy, and 2) DRL is typically designed for discrete domains, since Q-Learning with continuous action spaces is generally much harder. Furthermore, we note that while the WD is used in DRL, it is only for the convergence analysis of the DRL algorithm [2].

# 7   Experiments

Here we seek to test whether our approach to RL translates to performance gains for by evaluating BGPG and BGES, versus their respective baselines for a range of tasks. For each subsection we provide additional details in the Appendix.

## 7.1   Behavior-Guided Policy Gradient

Our key question is whether our techniques lead to outperformance for BGPG vs. baseline TRPO methods using KL divergence, which are widely used in the reinforcement learning community. For the BEM, we use the concatenation-of-actions, as used already in TRPO. We consider a variety of challenging problems from the DeepMind Control Suite [34] and Roboschool (RS). In Fig. 4 we see that BGPG does indeed outperform KL-based TRPO methods, with gains across all six environments. We also confirm results from [31] that a trust region typically improves performance.
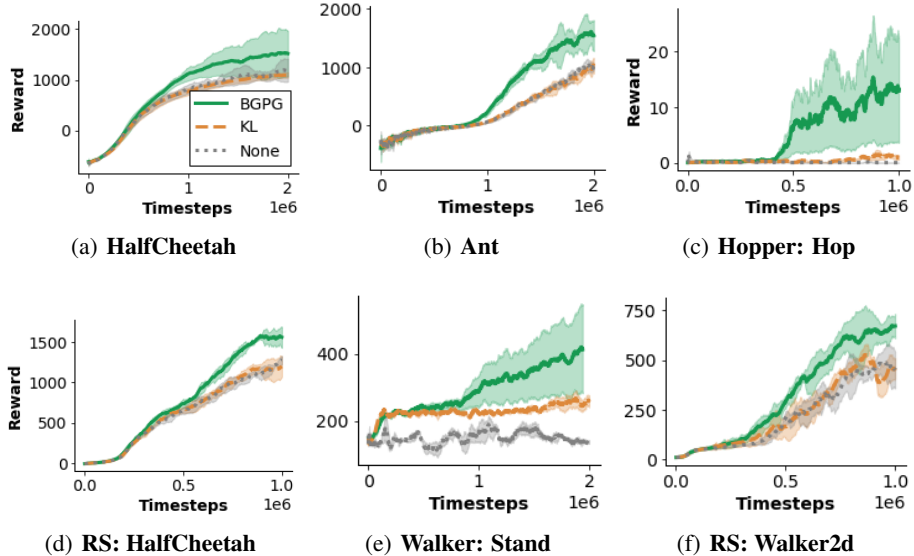
Figure 4: **BGPG vs. TRPO:** We compare BGPG and TRPO (KL divergence) on several continuous control tasks. As a baseline we also include results without a trust region ($\beta = 0$ in Algorithm 3). Plots show the mean $\pm$ std across 5 random seeds.

**Wall Clock Time:**   To illustrate computational benefits of alternating optimization of WD in BGPG, we compare it to the method introduced in [37]. In practice, the WD across different state samples can be optimized in a batched manner, details of which are in the Appendix. In Table 7.1 we see that BGPG is substantially faster.

Table 1: Clock time (s) to achieve a normalized reward of $90\%$ of the best achieved. All experiments were run on the same CPU.

|               | [37]  | BGPG  |
| ------------- | ----- | ----- |
| Pendulum      | 3720  | 777   |
| Hopper: Stand | 26908 | 10817 |
| Hopper: Hop   | 23542 | 12820 |
| Walker: Stand | 13497 | 4082  |

## 7.2   Behavior-Guided Evolution Strategies

Next we seek to evaluate the ability for BGES to use its behavioral repulsion for exploration.

**Deceptive Rewards**   A common challenge in RL is *deceptive* rewards. These arise since agents can only learn from data gathered via experience in the environment. To test BGES in this setting, we created two intentionally deceptive environments. In both cases the agent is penalized at each time step for its distance from a goal. The deception comes from a barrier, which means initially positive rewards from moving directly forward will lead to a suboptimal policy.

We consider two agents—a two-dimensional point and a larger quadruped. Details are provided in the Appendix (Section B). We compare with state-of-the-art **on-policy** methods for exploration: NSR-ES [7], which assumes the BEM is deterministic and uses the Euclidean distance to compare policies, and NoisyNet-TRPO [10]. We used the reward-to-go and final state BEMs for the quadruped and point respectively.
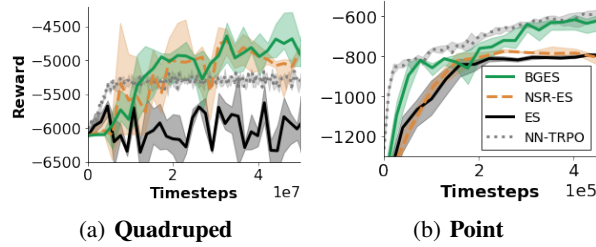
(a) **Quadruped**  (b) **Point**

Figure 5: **Deceptive Rewards.** Plots show the mean $\pm$ std across 5 random seeds for two environments: Quadruped and Point.

Policies avoiding the wall correspond to rewards: $R > -5000$ and $R > -800$ for the quadruped and point respectively. In the prior case an agent needs to first learn how to walk and the presence of the wall is enough to prohibit vanilla ES from even learning forward locomotion. As we see in Fig. 5, BGES is the only method that drives the agent to the goal in *both* settings.

## 8 Conclusion and Future Work

In this paper we proposed a new paradigm for on-policy learning in RL, where policies are embedded into expressive latent behavioral spaces and the optimization is conducted by utilizing the repelling/attraction signals in the corresponding probabilistic distribution spaces. The use of Wasserstein distances (WDs) guarantees flexibility in choosing cost funtions between embedded policy trajectories, enables stochastic gradient steps through corresponding regularized objectives (as opposed to KL divergence methods) and provides an elegant method, via their dual formulations, to quantify behaviorial difference of policies through the behavioral test functions. Furthermore, the dual formulations give rise to efficient algorithms optimizing RL objectives regularized with WDs.

We also believe the presented methods shed new light on several other challenging problems of modern RL, including: learning with safety guarantees (a repelling signal can be used to enforce behaviors away from dangerous ones) or anomaly detection for reinforcement learning agents (via the above score functions). Finally, we are interested in extending our method to the off policy setting.

## References

[1] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 214–223, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.

[2] M. G. Bellemare, W. Dabney, and R. Munos. A distributional perspective on reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, pages 449–458. JMLR.org, 2017.

[3] M. G. Bellemare, N. L. Roux, P. S. Castro, and S. Moitra. Distributional reinforcement learning with linear function approximation. In K. Chaudhuri and M. Sugiyama, editors, *Proceedings of Machine Learning Research*, volume 89 of *Proceedings of Machine Learning Research*, pages 2203–2211. PMLR, 16–18 Apr 2019.

[4] N. Bourbaki. *General Topology: Elements of Mathematics*. Addison-Wesley, 1966.

[5] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

[6] K. Choromanski, A. Pacchiano, J. Parker-Holder, Y. Tang, D. Jain, Y. Yang, A. Iscen, J. Hsu, and V. Sindhwani. Provably robust blackbox optimization for reinforcement learning. *CoRR*, abs/1903.02993, 2019.

[7] E. Conti, V. Madhavan, F. P. Such, J. Lehman, K. O. Stanley, and J. Clune. Improving exploration in evolution strategies for deep reinforcement learning via a population of novelty-seeking

agents. In *Advances in Neural Information Processing Systems 31, NeurIPS*, pages 5032–5043, 2018.

[8] W. Dabney, G. Ostrovski, D. Silver, and R. Munos. Implicit quantile networks for distributional reinforcement learning. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1096–1105, Stockholmsmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR.

[9] P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, Y. Wu, and P. Zhokhov. Openai baselines. `https://github.com/openai/baselines`, 2017.

[10] M. Fortunato, M. G. Azar, B. Piot, J. Menick, M. Hessel, I. Osband, A. Graves, V. Mnih, R. Munos, D. Hassabis, O. Pietquin, C. Blundell, and S. Legg. Noisy networks for exploration. In *International Conference on Learning Representations, ICLR*, 2018.

[11] A. Gajewski, J. Clune, K. O. Stanley, and J. Lehman. Evolvability es: Scalable and direct optimization of evolvability. In *Proceedings of the Genetic and Evolutionary Computation Conference*, GECCO '19, pages 107–115, New York, NY, USA, 2019. ACM.

[12] A. Genevay, M. Cuturi, G. Peyré, and F. Bach. Stochastic optimization for large-scale optimal transport. In *Advances in neural information processing systems*, pages 3440–3448, 2016.

[13] R. Jiang, A. Pacchiano, T. Stepleton, H. Jiang, and S. Chiappa. Wasserstein fair classification. *arXiv preprint arXiv:1907.12059*, 2019.

[14] S. Kakade and J. Langford. Approximately optimal approximate reinforcement learning. 2:267–274, 2002.

[15] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[16] S. Kullback and R. A. Leibler. On information and sufficiency. *Ann. Math. Statist.*, 22(1):79–86, 03 1951.

[17] S. J. Lee and Z. Popovic. Learning behavior styles with inverse reinforcement learning. *ACM Trans. Graph.*, 29(4):122:1–122:7, 2010.

[18] J. Lehman. *Evolution through the Search for Novelty*. PhD thesis, University of Central Florida, 2012.

[19] J. Lehman and K. O. Stanley. Exploiting open-endedness to solve problems through the search for novelty. In *Proceedings of the Eleventh International Conference on Artificial Life (Alife XI*. MIT Press, 2008.

[20] H. Mania, A. Guy, and B. Recht. Simple random search of static linear policies is competitive for reinforcement learning. *Advances in Neural Information Processing Systems 31, NeurIPS*, 2018.

[21] E. Meyerson, J. Lehman, and R. Miikkulainen. Learning behavior characterizations for novelty search. In *Proceedings of the 2016 on Genetic and Evolutionary Computation Conference, Denver, CO, USA, July 20 - 24, 2016*, pages 149–156, 2016.

[22] C. A. Micchelli, Y. Xu, and H. Zhang. Universal kernels. *Journal of Machine Learning Research*, 7:2651–2667, 2006.

[23] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. A. Riedmiller. Playing atari with deep reinforcement learning. *ArXiv*, abs/1312.5602, 2013.

[24] J. K. Pugh, L. B. Soros, and K. O. Stanley. Quality diversity: A new frontier for evolutionary computation. *Frontiers in Robotics and AI*, 2016.

[25] C. Qu, S. Mannor, and H. Xu. Nonlinear distributional gradient temporal-difference learning. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5251–5260, Long Beach, California, USA, 09–15 Jun 2019. PMLR.

[26] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 1177–1184. Curran Associates, Inc., 2008.

[27] P. H. Richemond and B. Maginnis. On wasserstein reinforcement learning and the fokker-planck equation. *CoRR*, abs/1712.07185, 2017.

[28] M. Rowland, M. Bellemare, W. Dabney, R. Munos, and Y. W. Teh. An analysis of categorical distributional reinforcement learning. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 29–37, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR.

[29] M. Rowland, R. Dadashi, S. Kumar, R. Munos, M. G. Bellemare, and W. Dabney. Statistics and samples in distributional reinforcement learning. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5528–5536, Long Beach, California, USA, 09–15 Jun 2019. PMLR.

[30] T. Salimans, J. Ho, X. Chen, S. Sidor, and I. Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *CoRR, abs/1703.03864*, 2017.

[31] J. Schulman, S. Levine, P. Abbeel, M. I. Jordan, and P. Moritz. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 1889–1897, 2015.

[32] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[33] R. S. Sutton, A. G. Barto, et al. *Introduction to Reinforcement Learning*, volume 135. MIT press Cambridge, 1998.

[34] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.

[35] C. Villani. Optimal transport: Old and new. *Springer*, 2008.

[36] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

[37] R. Zhang, C. Chen, C. Li, and L. Carin. Policy optimization as wasserstein gradient flows. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pages 5741–5750, 2018.

# Appendix: Behavior-Guided Reinforcement Learning
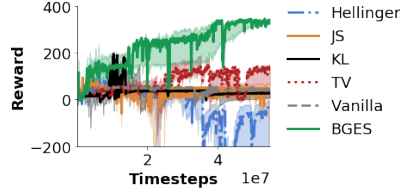
## A    Additional Experiments



Figure 6: **Escaping Local Maxima** A comparison of BGES with those using different distances on Policy Embeddings.

### A.1    Escaping Local Maxima.

In Fig. 6 we compare our methods with methods using regularizers based on other distances or divergences (specifically, Hellinger, Jensen-Shannon (JS), KL and Total Variation (TV) distances), as well as vanilla ES (i.e., with no distance regularizer). Experiments were performed on a **Swimmer** environment from $\mathrm{OpenAI\,Gym}$ [5], where the number of samples of the ES optimizer was drastically reduced. BGES is the only one that manages to obtain good policies which also proves that the benefits come here not just from introducing the regularizer, but from its particular form.

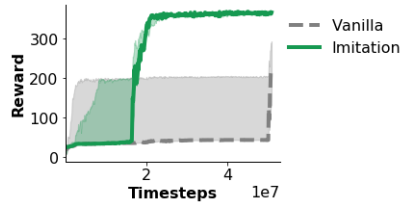### A.2    Imitation Learning



Figure 7: **Imitation Learning.**

As discussed in Section 5.4, we can also utilize the BGES algorithm for imitation learning, by setting $\beta < 0$, and using an expert's trajectories for the Policy Embedding. For this experiment we use the reward-to-go BEM (Section 5). In Fig. 7, we show that this approach significantly outperforms vanilla ES on the $\mathrm{Swimmer}$ task. Although conceptually simple, we believe this could be a powerful approach with potential extensions, for example in designing safer algorithms.

# B Further Experimental Details

## B.1 BGPG

Here we reproduce a full version of Algorithm 3:

---

**Algorithm 5** Behvaior-Guided Policy Gradient with **On-Policy** Embeddings

---

**Input:** Initialize stochastic policy $\pi_0$ parametrized by $\theta_0$, $\beta < 0, \eta > 0, M, L \in \mathbb{N}$

**for** $t = 1, \ldots, T$ **do**

   1. Run $\pi_{t-1}$ in the environment to get advantage values $A^{\pi_{t-1}}(s, a)$ and trajectories $\{\tau_i^{(t)}\}_{i=1}^M$

   2. Update policy and test functions via several alternating gradient steps over the objective:

$$
F(\theta) = \mathop{\mathbb{E}}_{\tau_1, \tau_2 \sim \mathbb{P}_{\pi_{t-1}} \otimes \mathbb{P}_{\pi_\theta}} \Big[ \sum_{i=1}^H A^{\pi_{t-1}}(s_i, a_i) \frac{\pi_\theta(a_i|s_i)}{\pi_{t-1}(a_i|s_i)}
$$
$$
+ \beta \lambda_1(\Phi(\tau_1)) - \beta \lambda_2(\Phi(\tau_2)) + \beta \gamma \exp\left( \frac{\lambda_1(\Phi(\tau_1)) - \lambda_2(\Phi(\tau_2)) - C(\Phi(\tau_1)), \Phi(\tau_2))}{\gamma} \right) \Big]
$$

   Where $\tau_1 = s_0, a_0, r_0, \cdots, s_H, a_H, r_H$. Let $\theta_{t-1}^{(0)} = \theta_{t-1}$.

   **for** $\ell = 1, \cdots, L$ **do**

      a. Approximate $\mathbb{P}_{\pi_{t-1}} \otimes \mathbb{P}_{\pi_\theta}$ via $\frac{1}{M}\{\tau_i^{(t)}\}_{i=1}^M \otimes \frac{1}{M}\{\tau_i^\theta\}_{i=1}^M := \hat{P}_{\pi_t, \pi_\theta}$ where $\tau_i^\theta \overset{\text{i.i.d}}{\sim} \mathbb{P}_{\pi_\theta}$

      b. Take SGA step $\theta_{t-1}^{(\ell)} = \theta_{t-1}^{(\ell-1)} + \eta \hat{\nabla}_\theta \hat{F}(\theta_{t-1}^{(\ell-1)})$ using samples from $\hat{P}_{\pi_{t-1}, \pi_\theta}$.

      c. Use samples from $\hat{P}_{\pi_{t-1}, \pi_\theta}$ and Algorithm 1 to update $\lambda_1, \lambda_2$.

Set $\theta_t = \theta_{t-1}^{(M)}$.

---

---

**Algorithm 6** Behvaior-Guided Policy Gradient with **Off-Policy** Embeddings

---

**Input:** Initialize stochastic policy $\pi_0$ parametrized by $\theta_0$, $\beta < 0, \eta > 0, M, L \in \mathbb{N}$, state probing distribution $\mathbb{P}_{\mathcal{S}_0}$.

**for** $t = 1, \ldots, T$ **do**

   1. Run $\pi_{t-1}$ in the environment to get advantage values $A^{\pi_{t-1}}(s, a)$ . 2. Update policy and test functions via several alternating gradient steps over the objective:

$$
F(\theta) = \mathop{\mathbb{E}}_{\tau_1 \sim \mathbb{P}_{\pi_{t-1}}} \Big[ \sum_{i=1}^H A^{\pi_{t-1}}(s_i, a_i) \frac{\pi_\theta(a_i|s_i)}{\pi_{t-1}(a_i|s_i)} \Big]
$$
$$
+ \mathbb{E}_{s_1, s_2 \sim \mathbb{P}_{\mathcal{S}_\theta} \otimes \mathbb{P}_{\mathcal{S}_{t-1}}} \Big[ \beta \lambda_1(s_1, \pi_\theta(s_1)) - \beta \lambda_2(s_2, \pi_{t-1}(s_2))
$$
$$
+ \beta \gamma \exp\left( \frac{\lambda_1(s_1, \pi_\theta(s_1)) - \lambda_2(s_2, \pi_{t-1}(s_2)) - C((s_1, \pi_\theta(s_1)), (s_2, \pi_{t-1}(s_2)))}{\gamma} \right) \Big]
$$

   Where $\tau_1 = s_0, a_0, r_0, \cdots, s_H, a_H, r_H$. Let $\theta_{t-1}^{(0)} = \theta_{t-1}$.

   **for** $\ell = 1, \cdots, L$ **do**

      a. Approximate the expectation $\mathbb{E}_{s_1, s_2 \sim \mathbb{P}_{\mathcal{S}_\theta} \otimes \mathbb{P}_{\mathcal{S}_{t-1}}}$ via $2M$ samples.

      b. Take SGA step $\theta_{t-1}^{(\ell)} = \theta_{t-1}^{(\ell-1)} + \eta \hat{\nabla}_\theta \hat{F}(\theta_{t-1}^{(\ell-1)})$ using samples from a. and trajectories from current $\pi_\theta$.

      c. Use samples from $a.$ and Algorithm 1 to update $\lambda_1, \lambda_2$.

Set $\theta_t = \theta_{t-1}^{(M)}$.

---

**A Lower-variance Gradient Estimator via Off-Policy embeddings:** As explained in Section 5.2, the BGPG considers an objective which involves two parts: the conventional surrogate loss function for policy optimization [32], and a loss function that involves the Behavior Test Functions.

Though we could apply vanilla reinforced gradients on both parts, it is straightforward to notice that the second part can be optimized with reparameterized gradients [15], which arguably have lower variance compared to the reinforced gradients. In particular, we note that under random feature approximation (5), as well as the action-concatenation embedding, the Wasserstein distance loss $\widehat{\text{WD}}_\gamma(P_{\pi_\theta}^\Phi, P_b^\Phi)$ is a differentiable function of $\theta$. To see this more clearly, notice that under a Gaussian policy $a \sim \mathcal{N}(\mu_\theta(s), \sigma_\theta(s)^2)$ the actions $a = \mu_\theta(s) + \sigma_\theta(s) \cdot \epsilon$ are reparametrizable for $\epsilon$ being standard Gaussian noises. We can directly apply the reparametrization trick to this second objective to obtain a gradient estimator with potentially much lower variance. In our experiments, we applied this lower-variance gradient estimator. In Algorithm 6 we allow the state probing distribution to evolve with the iteration index of the algorithm $t$.

**Trust Region Policy Optimization:**  Though the original TRPO [31] construct the trust region based on KL-divergence, we propose to construct the trust region with WD. For convenience, we adopt a dual formulation of the trust region method and aim to optimize the augmented objective $\mathbb{E}_{\tau \sim \pi_\theta}[R(\tau)] - \beta \text{WD}_\gamma(\mathbb{P}_{\pi'}^\Phi, \mathbb{P}_{\pi_\theta}^\Phi)$. We apply the concatenation-of-actions embedding and random feature maps to calculate the trust region. We identify several important hyperparameters: the RKHS (for the test function) is produced by RBF kernel $k(x, y) = \exp(\|x - y\|_2^2/\sigma^2)$ with $\sigma = 0.1$; the number of random features is $D = 100$; recall the embedding is $\Phi(\tau) = [a_1, a_2...a_H]$ where $H$ is the horizon of the trajectory, here we take 10 actions per state and embed them together, this is equivalent to reducing the variance of the gradient estimator by increasing the sample size; the regularized entropy coefficient in the WD definition as $\gamma = 0.1$; the trust region trade-off constant $\beta \in \{0.1, 1, 10\}$. The alternate gradient descent is carried out with $T = 100$ alternating steps and test function coefficients $\mathbf{p} \in \mathbb{R}^D$ are updated with learning rate $\alpha_{\mathbf{p}} = 0.01$.

The baseline algorithms are: No trust region, and trust region with KL-divergence. The KL-divergence is identified by a maximum KL-divergence threshold per update, which we set to $\epsilon = 0.01$.

Across all algorithms, we adopt the open source implementation [9]. Hyper-parameters such as number of time steps per update as well as implementation techniques such as state normalization are default in the original code base.

The additional experiment results can be found in Figure 8 where we show comparison on additional continuous control benchmarks: Tasks with DM are from DeepMind Contol Suites [34]. We see that the trust region constructed from the WD consistently outperforms other baselines (importantly, trust region methods are always better than the baseline without trust region, this confirms that trust region methods are critical in stabilizing the updates).
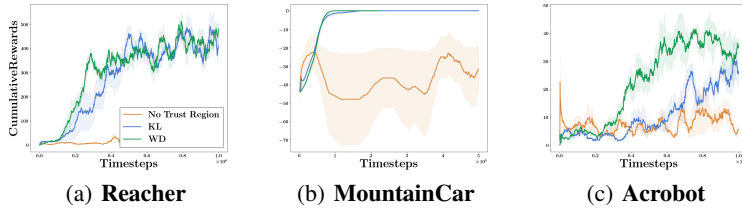


|   (a) **Reacher**   |   (b) **MountainCar**   |   (c) **Acrobot**   |

Figure 8: Additional Experiment on TRPO. We compare No Trust Region with two alternative trust region constructions: KL-divergence and Wassertein distance (ours).

**Wasserstein AO vs. Particle Approximation:**  To calculate the regularized Wasserstein distance, we propose a gradient descent method that iteratively updates the test function. The alternting optimization (AO) scheme consists of updating both the test function and the distribution parameters such that the regularized Wasserstein distance of the trainable distribution against the reference distribution is minimized. Alternatively, we can also adopt a particle approximation method to calculate the Wasserstein distance and update the distribution parameters using an approximate gradient descent method [37]. We see the benefit in clock time in Fig 9.
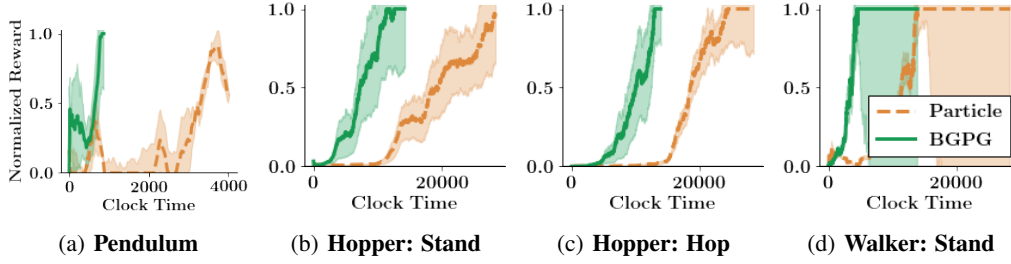
Figure 9: The clock-time comparison (in sec) of BGPG (alternating optimization) with particle approximation.

One major advantage of BGPG against particle approximation is its ease of parallelization. In particular, when using the concatenation-of-actions embedding, the aggregate Wasserstein distance can be decomposeed into an average of a set of Wasserstein distances over states. To calculate this aggregated gradient, BGPG can easily leverage the matrix multiplication; on the other hand, particle approximation requires that the dual optimal variables of each subproblem be computed, which is not straightforward to parallelize.

We test both methods in the context of trust region policy search, in which we explicitly calculate the Wasserstein distance of consecutive policies and enforce the constraints using a line search as in [31]. Both methods require the trust region trade-off parameter $\beta \in \{0.1, 1, 10\}$. We adopt the particle method in [37] where for each state there are $M = 16$ particles. The gradients are derived based a RKHS where we adaptively adjust the coefficient of the RBF kernel based on the mean distance between particles. For the AO, we find that it suffices to carry out $T \in \{1, 5, 10\}$ gradient descents to approximate the regularized Wasserstein distance.

## B.2  BGES

Here we reproduce a detailed version of Algorithm 4:

---

**Algorithm 7** Behavior-Guided Evolution Strategies with **On-Policy** Embeddings

---

**Input:** learning rate $\eta$, noise standard deviation $\sigma$, iterations $T$, BEM $\Phi$, $\beta$
**Initialize:** Initial policy $\pi_0$ parametrized by $\theta_0$, Behavioral Test Functions $\lambda_1, \lambda_2$. Evaluate policy $\pi_0$ to return trajectory $\tau_0$ and subsequently use the BEM to produce an initial $\hat{\mathbb{P}}^\Phi_{\pi_0}$.
**for** $t = 1, \ldots, T - 1$ **do**

    1. Sample $\epsilon_1, \cdots, \epsilon_n$ independently from $\mathcal{N}(0, I)$.

    2. Evaluate policies $\{\pi_t^k\}_{k=1}^n$ parameterized by $\{\theta_t + \sigma\epsilon_k\}_{k=1}^n$ to return rewards $R_k$ and trajectories $\tau_k$ for all $k$.

    3. Use BEM to map trajectories $\tau_k$ to produce empirical $\hat{\mathbb{P}}^\Phi_{\pi_t^k}$ for all $k$.

    4. Update $\lambda_1$ and $\lambda_2$ using Algorithm 1, where $\mu = \frac{1}{n} \cup_{k=1}^n \hat{\mathbb{P}}^\Phi_{\pi_{t-1}^k}$ and $\nu = \frac{1}{n} \cup_{k=1}^n \hat{\mathbb{P}}^\Phi_{\pi_t^k}$ are the uniform distribution over the set of from 3 for $t - 1$ and $t$.

    5. Approximate $\widehat{\mathrm{WD}}\gamma(\mathbb{P}^\Phi_{\pi_t^k}, \mathbb{P}^\Phi_{\pi_t})$ plugging in $\lambda_1, \lambda_2$ into Eq. 5 for each perturbed policy $\pi_k$

    6. Update Policy: $\theta_{t+1} = \theta_t + \eta\nabla_{ES}F$, where:

$$\nabla_{ES}F = \frac{1}{\sigma}\sum_{k=1}^n [(1-\beta)(R_k - R_t) + \beta\widehat{\mathrm{WD}}\gamma(\mathbb{P}^\Phi_{\pi_t^k}, \mathbb{P}^\Phi_{\pi_t})]\epsilon_k$$

---

**Efficient Exploration:** To demonstrate the effectiveness of our method in exploring deceptive environments, we constructed two new environments using the MuJoCo simulator. For the point environment, we have a 6 dimensional state and 2 dimensional action, with the reward at each timestep calculated as the distance between the agent and the goal. We use a horizon of 50 which is sufficient to reach the goal. The quadruped environment is based on Ant from the Open AI Gym [5], and has a similar reward structure to the point environment but a much larger state space (113) and action space (8). For the quadruped, we use a horizon length of 400.

18

To leverage the trivially parallelizable nature of ES algorithms, we use the ray library, and distribute the rollouts across 72 workers using AWS. Since we are sampling from an isotropic Gaussian, we are able to pass only the seed to the workers, as in [30]. However we do need to return trajectory information to the master worker.

For both the point and quadruped agents, we use random features with dimensionality $m = 1000$, and 100 warm-start updates for the WD at each iteration. For point, we use the final state embedding, learning rate $\eta = 0.1$ and $\sigma = 0.01$. For the quadruped, we use the reward-to-go embedding, as we found this was needed to learn locomotion, as well as a learning rate of $\eta = 0.02$ and $\sigma = 0.02$. The hyper-parameters were the same for all ES algorithms. When computing the WD, we used the previous 2 policies, $\theta_{t-1}$ and $\theta_{t-2}$.

Our approach includes several new hyperparameters, such as the kernel for the Behavioral Test Functions and the choice of BEM. For our experiments we did not perform any hyperparameter optimization. We only considered the rbf kernel, and only varied the BEM for BGES. For BGES, we demonstrated several different BEMs, and we show an ablation study for the point agent in Fig. 12 where we see that both the reward-to-go (RTG) and Final State (SF) worked, but the vector of all states (SV) did not (for 5 seeds). We leave learned BEMs as exciting future work.
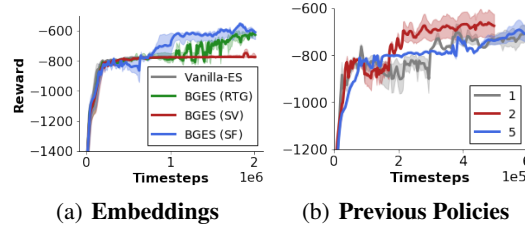


(a) **Embeddings**    (b) **Previous Policies**

Figure 10: A sensitivity analysis investigating a) the impact of the embedding and b) the number of previous policies $\theta_{t-i}, i \in 1, 2, 5$

For embeddings, we compare the reward-to-go (RTG), concatenation of states (SV) and final state (SF). In both the RTG and SF case the agent learns to navigate past the wall ($> -800$). For the number of previous policies, we use the SF embedding, and using 2 appears to work best, but both 1 and 5 do learn the correct behavior.

**Escaping Local Maxima:**  We also demonstrated that our method leads to faster training even in more standard settings, where exploration is not that crucial, but the optimization can be trapped in local maxima. To show it, we compared baseline ES algorithm for ES optimization from [30] with its enhancements, where regularizers using different metrics on the space of probabilistic distributions corresponding to policy embeddings were used, as in the previous paragraph. We noticed that adding Wasserstein regularizers drastically improved optimization, whereas regularizers based on other distances/divergencies, namely: Hellinger, Jensen-Shannon, KL and TV did not have any impact. We considered Swimmer task from OpenAI Gym and to make it challenging, reduced the number of perturbations per iteration to 80. In that setting our method was the only one that was not trapped in local maxima and managed to learn effective policies.

### B.3  Imitation Learning:

For the Imitation Learning experiment we used the reward-to-go embedding, with learning rate $\eta = 0.1$ and $\sigma = 0.01$. We use one oracle policy, which achieves $> 360$ on the environment. The only information provided to the algorithm is the embedded trajectory, used to compute the WD. This has exciting future applications since no additional information about the oracle is required in order to significantly improve learning.

### B.4  Repulsion Learning and Attraction learning

Here we reproduce a full version of Algorithm 2:

**Algorithm 8** Behvaior-Guided Repulsion (and Attraction) Learning with **On-Policy** Embeddings

**Input:** $\beta, \eta > 0, M \in \mathbb{N}$

**Initialize:** Initial stochastic policies $\pi_0^{\mathbf{a}}, \pi_0^{\mathbf{b}}$, parametrized by $\theta_0^{\mathbf{a}}, \theta_0^{\mathbf{b}}$ respectively, Behavioral Test Functions $\lambda_1^{\mathbf{a}}, \lambda_2^{\mathbf{b}}$

**for** $t = 1, \ldots, T$ **do**

   1.  Collect $M$ trajectories $\{\tau_i^{\mathbf{a}}\}_{i=1}^M$ from $\mathbb{P}_{\pi_{t-1}^{\mathbf{a}}}$ and $M$ trajectories $\{\tau_i^{\mathbf{b}}\}_{i=1}^M$ from $\mathbb{P}_{\pi_{t-1}^{\mathbf{b}}}$.

      Approximate $\mathbb{P}_{\pi_{t-1}^{\mathbf{a}}} \bigotimes \mathbb{P}_{\pi_{t-1}^{\mathbf{b}}}$ via $\frac{1}{M}\{\tau_i^{\mathbf{a}}\}_{i=1}^M \bigotimes \frac{1}{M}\{\tau_i^{\mathbf{b}}\}_{i=1}^M := \hat{P}_{\pi_{t-1}^{\mathbf{a}}, \pi_{t-1}^{\mathbf{b}}}$

   2. Form two distinct surrogate rewards for joint trajectories of agents **a** and **b**:

$$\tilde{R}_{\mathbf{a}}(\tau_1, \tau_2) = \mathcal{R}(\tau_1) + \beta\lambda_1^{\mathbf{a}}(\Phi(\tau_1)) + \beta\gamma \exp\left(\frac{\lambda_1^{\mathbf{a}}(\Phi(\tau_1)) - \lambda_2^{\mathbf{b}}(\Phi(\tau_2)) - C(\Phi(\tau_1)), \Phi(\tau_2))}{\gamma}\right)$$

$$\tilde{R}_{\mathbf{b}}(\tau_1, \tau_2) = \mathcal{R}(\tau_2) - \beta\lambda_2^{\mathbf{b}}(\Phi(\tau_2)) + \beta\gamma \exp\left(\frac{\lambda_1^{\mathbf{a}}(\Phi(\tau_1)) - \lambda_2^{\mathbf{b}}(\Phi(\tau_2)) - C(\Phi(\tau_1)), \Phi(\tau_2))}{\gamma}\right)$$

   3. For $\mathbf{c} \in \{\mathbf{a}, \mathbf{b}\}$ use the Reinforce estimator to take gradient steps:

$$\theta_t^{\mathbf{c}} = \theta_{t-1}^{\mathbf{c}} + \eta \mathop{\mathbb{E}}_{\tau^{\mathbf{a}}, \tau^{\mathbf{b}} \sim \hat{P}_{\pi_{t-1}^{\mathbf{a}}, \pi_{t-1}^{\mathbf{b}}}} \left[ \tilde{\mathcal{R}}_{\mathbf{c}}(\tau^{\mathbf{a}}, \tau^{\mathbf{b}}) \left( \sum_{i=0}^{H-1} \nabla_{\theta_{t-1}^{\mathbf{c}}} \log\left(\pi_{t-1}^{\mathbf{c}}(a_i^{\mathbf{c}}|s_i^{\mathbf{c}})\right) \right) \right]$$

     Where $\tau^{\mathbf{a}} = s_0^{\mathbf{a}}, a_0^{\mathbf{a}}, r_0^{\mathbf{a}}, \cdots, s_H^{\mathbf{a}}, a_H^{\mathbf{a}}, r_H^{\mathbf{a}}$ and $\tau^{\mathbf{b}} = s_0^{\mathbf{b}}, a_0^{\mathbf{b}}, r_0^{\mathbf{b}}, \cdots, s_H^{\mathbf{b}}, a_H^{\mathbf{b}}, r_H^{\mathbf{b}}$.

   5. Use samples from $\hat{P}_{\pi_{t-1}^{\mathbf{a}}, \pi_{t-1}^{\mathbf{b}}}$ and Algorithm 1 to update the Behavioral Test Functions $\lambda_1^{\mathbf{a}}, \lambda_2^{\mathbf{b}}$.

Algorithm 8 is the de version of the repulsion algorithm from Section 5.2. The algorithm maintains two policies $\pi^{\mathbf{a}}$ and $\pi^{\mathbf{b}}$. Each policy is optimized by taking a policy gradient step (using the Reinforce gradient estimator) in the direction optimizing surrogate rewards $\tilde{\mathcal{R}}_{\mathbf{a}}$ and $\tilde{\mathcal{R}}_{\mathbf{b}}$ that combines the signal from the task's reward function $\mathcal{R}$ and the repulsion (or attraction) score encoded by the behavioral test functions $\lambda^{\mathbf{a}}$ and $\lambda^{\mathbf{b}}$.

We conducted experiments testing Algorithm 2 on a simple Mujoco environment consisting of a particle that moves on the plane and whose objective is to learn a policy that allows it to reach one of two goals. Each policy outputs a velocity vector and stochasticity is achieved by adding Gaussian noise to the mean velocity encoded by a neural network with two size 5 hidden layers and ReLu activations. If an agent performs action $a$ at state $s$, it moves to state $a + s$. The reward of an agent after performing action $a$ at state $s$ equals $-\|a\|^2 * 30 - \min(d(s, \text{Goal}_1), d(s, \text{Goal}_2))^2$ where $d(x, y)$ denotes the distance between $x$ and $y$ in $\mathbb{R}^2$. The initial state is chosen by sampling a Gaussian distribution with mean $\binom{0}{0}$ and diagonal variance 0.1. In each iteration step we sample 100 trajectories. In the following pictures we plot the policies' behavior by plotting 100 trajectories of each. The embedding $\Phi : \Gamma \to \mathbb{R}$ maps trajectories $\tau$ to their mean displacement in the $x-$axis. We use the squared absolute value difference as the cost function. When $\beta < 0$ we favour attraction and the agent are encouraged to learn a similar policy to solve the same task.
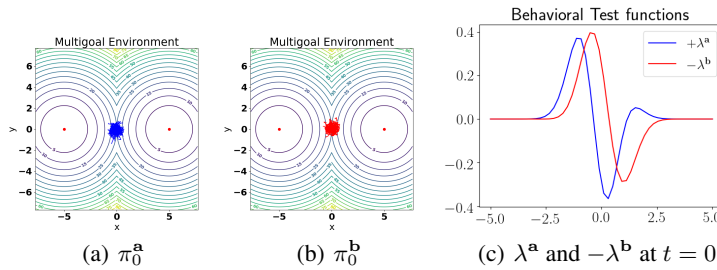


(a) $\pi_0^{\mathbf{a}}$       (b) $\pi_0^{\mathbf{b}}$       (c) $\lambda^{\mathbf{a}}$ and $-\lambda^{\mathbf{b}}$ at $t = 0$

Figure 11: Initial state of policies $\pi^{\mathbf{a}}, \pi^{\mathbf{b}}$ and Behavioral Test functions $\lambda^{\mathbf{a}}, \lambda^{\mathbf{b}}$ in the Multigoal environment.

There are two optimal policies, moving the particle to the left goal or moving it to the right goal. We now plot how the policies' behavior and evolves throughout optimization and how the Behavioral Test Functions guide the optimization by favouring the two policies to be close by or far apart.
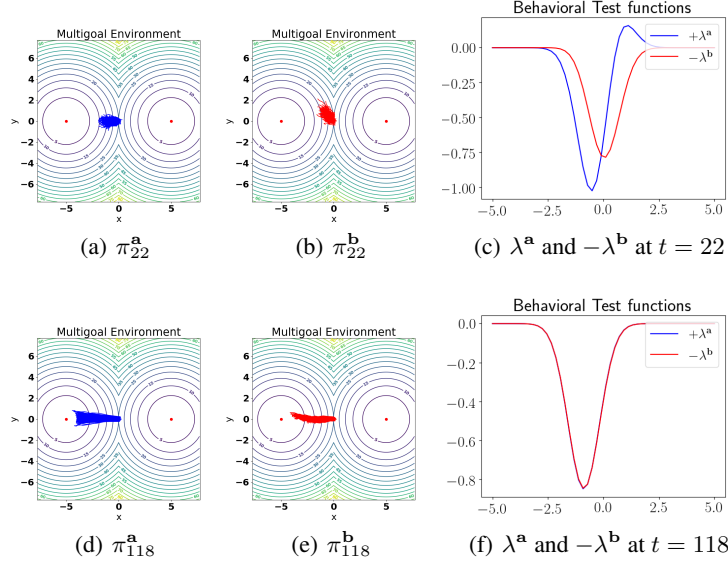


Figure 12: Evolution of the policies and Behavioral Test Functions throughout optimization.

Let $\mathcal{X}$ and $\mathcal{Y}$ be the domains of two measures $\mu$, and $\nu$. Recall that in case $\gamma = 0$, $\mathcal{X} = \mathcal{Y}$, and $C(x, x) = 0$ for all $x \in \mathcal{X}$, then $\lambda_\mu^*(x) = \lambda_\nu^*(x) = \lambda^*(x)$ for all $x \in \mathcal{X}$. In the case of regularized Wasserstein distances with $\gamma > 0$, this relationship may not hold true even if the cost satisfies the same diagonal assumption. For example when the regularizing measure is the product measure, and $\mu, \nu$ have disjoint supports, since the soft constraint $\gamma \exp\left(\frac{\lambda_\mu(\mathbf{x}) - \lambda_\nu(\mathbf{y}) - C(\mathbf{x}, \mathbf{y})}{\gamma}\right)$ is enforced in expectation over the product measure there may exist optimal solutions $\lambda_\mu^*, \lambda_\nu^*$ that do not satisfy $\lambda_\mu^* = \lambda_\nu^*$.

# C   Theoretical results

We start by exploring some properties of the Wasserstein distance and its interaction with some simple classes of embeddings. The first lemma we show has the intention to show conditions under which two policies can be shown to be equal provided the Wasserstein distance between its trajectory embeddings is zero. This result implies that our framework is capable of capturing equality of policies when the embedding space equals the space of trajectories.

**Lemma C.1.** *Let $\mathcal{S}$ and $\mathcal{A}$ be finite sets, the* MDP *be episodic (i.e. of finite horizon $H$), and* $\Phi(\tau) = \sum_{t=0}^{H} e_{s_t, a_t}$ *with* $e_{s,a} \in \mathbb{R}^{|\mathcal{S}| + |\mathcal{A}|}$ *the indicator vector for the state action pair $(s, a)$. Let* $C(\mathbf{v}, \mathbf{w}) = \|\mathbf{v} - \mathbf{w}\|_p^p$ *for $p \geq 1$. If $\gamma = 0$ and $\mathrm{WD}_\gamma(\mathbb{P}_\pi^\Phi, \mathbb{P}_{\pi'}^\Phi) = 0$ then $\pi = \pi'$.*

*Proof.* If $\mathrm{WD}_\gamma(\mathbb{P}_\pi^\Phi, \mathbb{P}_{\pi'}^\Phi) = 0$, there exists a coupling $\Pi$ between $\mathbb{P}_\pi^\Phi$ and $\mathbb{P}_{\pi'}^\Phi$ such that:

$$\mathbb{E}_{u,v \sim \Pi}\left[\|u - v\|_p^p\right] = 0$$

Consequently:

$$\mathbb{E}_{u,v \sim \Pi}\left[\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} |u_{s,a} - v_{s,a}|^p\right] = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathbb{E}_{u,v \sim \Pi}\left[|u_{s,a} - v_{s,a}|^p\right] = 0$$

Therefore for all $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\left|\mathbb{E}_{u \sim \mathbb{P}_\pi^\Phi}[u_{s,a}] - \mathbb{E}_{v \sim \mathbb{P}_{\pi'}^\Phi}[v_{s,a}]\right|^p \leq \mathbb{E}_{u,v \sim \Pi}\left[|u_{s,a} - v_{s,a}|^p\right] = 0$$

21

Where $u_{s,a}$ and $v_{s,a}$ denote the $(s,a)$ entries of $u$ and $v$ respectively. Notice that for all $(s,a) \in \mathcal{S} \times \mathcal{A}$:

$$\mathbb{P}^\Phi_\pi(s,a) = \mathbb{P}^\Phi_{\pi'}(s,a) \tag{11}$$

Since for all $s \in \mathcal{S}$ and $p \geq 1$:

$$\left| \sum_{a \in \mathcal{A}} u_{s,a} - v_{s,a} \right|^p \leq \sum_{a \in \mathcal{A}} |u_{s,a} - v_{s,a}|^p$$

Therefore for all $s \in \mathcal{S}$:

$$\left| \mathbb{E}_{u \sim \mathbb{P}^\Phi_\pi} \left[ \sum_{a \in \mathcal{A}} u_{s,a} \right] - \mathbb{E}_{v \sim \mathbb{P}^\Phi_{\pi'}} \left[ \sum_{a \in \mathcal{A}} v_{s,a} \right] \right|^p \leq \mathbb{E}_{u,v \sim \Pi} \left[ \sum_{a \in \mathcal{A}} |u_{s,a} - v_{s,a}|^p \right] = 0$$

Consequently $\mathbb{P}^\Phi_\pi(s) = \mathbb{P}^\Phi_{\pi'}(s)$ for all $s \in \mathcal{S}$. By Bayes rule, this plus equation 11 yields:

$$\mathbb{P}^\Phi_\pi(a|s) = \mathbb{P}^\Phi_{\pi'}(a|s)$$

And therefore: $\pi = \pi'$. $\qquad\square$

These results can be extended in the following ways:

1. In the case of a continuous state space, it is possible to define embeddings using Kernel density estimators. Under the appropriate smoothness conditions on the visitation frequencies, picking an adequate bandwidth and using the appropriate norm to compare different embeddings it is possible to derive similar results to those in Lemma C.1 for continuous state spaces.

2. For embeddings such as $\Phi_5$ in Section 3.1 or $\Phi(\tau) = \sum_{t=0}^H e_{s_t,a_t}$, when $\gamma = 0$, if $\mathrm{WD}_\gamma(\mathbb{P}^\Phi_\pi, \mathbb{P}^\Phi_{\pi'}) \leq \epsilon$ then $|V(\pi) - V(\pi')| \leq \epsilon R$ for $R = \max_{\tau \in \Gamma} \mathcal{R}(\tau)$ thus implying that a small Wasserstein distance between $\pi$ and $\pi'$s PPEs implies a small difference in their value functions.

## C.1   Random features stochastic gradients

Let $\phi_\kappa$ and $\phi_\ell$ be two feature maps over $\mathcal{X}$ and $\mathcal{Y}$ and corresponding to kernels $\kappa$ and $\ell$ respectively. For this and the following sections we will make use of the following expression:

$$G(\mathbf{p}^\mu, \mathbf{p}^\nu) = \beta \int_\mathcal{X} (\mathbf{p}^\mu)^\top \phi_\kappa(\mathbf{x}) d\mu(\mathbf{x}, \theta) - \beta \int_\mathcal{Y} (\mathbf{p}^\nu)^\top \phi_\ell(\mathbf{y}) d\nu(\mathbf{y}) + \tag{12}$$

$$\gamma\beta \int_{\mathcal{X} \times \mathcal{Y}} \exp\left( \frac{(\mathbf{p}^\mu)^\top \phi_\kappa(\mathbf{x}) - (\mathbf{p}^\nu)^\top \phi_\ell(\mathbf{y}) - C(\mathbf{x}, \mathbf{y})}{\gamma} \right) d\mu(\mathbf{x}) d\nu(\mathbf{y})$$

We now show how to compute gradients with respect to the random feature maps:

**Lemma C.2.** *The gradient $\nabla_{\binom{\mathbf{p}^\mu}{\mathbf{p}^\nu}} G(\mathbf{p}^\mu, \mathbf{p}^\nu)$ of the objective function from Equation 12 with respect to the parameters $\binom{\mathbf{p}^\mu}{\mathbf{p}^\nu}$ satisfies:*

$$\nabla_{\binom{\mathbf{p}^\mu}{\mathbf{p}^\nu}} G(\mathbf{p}^\mu, \mathbf{p}^\nu) = \beta \mathbb{E}_{(\mathbf{x},\mathbf{y}) \sim \mu \otimes \nu} \left[ \left( 1 - \exp\left( \frac{(\mathbf{p}^\mu)^\top \phi_\kappa(\mathbf{x}) - (\mathbf{p}^\nu)^\top \phi_\ell - C(\mathbf{x}, \mathbf{y})}{\gamma} \right) \right) \binom{\phi_\kappa(\mathbf{x})}{-\phi_\ell(\mathbf{y})} \right]$$

*Proof.* A simple use of the chain rule, taking the gradients inside the expectation, and the fact that $\mathbf{p}^\mu$ and $\mathbf{p}^\nu$ are vectors yields the desired result. $\qquad\square$

The main consequence of this formulation is the stochastic gradients we use in Algorithm 1.

## C.2 Behavior Guided Policy Gradient and Wasserstein trust region

For a policy $\pi$, we denote as: $V^\pi$, $Q^\pi$ and $A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$ the: value function, $Q$-function and advantage function.

The chief goal of this section is to prove Theorem 5.1. We restate the section's definitions here for the reader's convenience: To ease the discussion we make the following assumptions:

- Finite horizon $T$.
- Undiscounted MDP.
- States are time indexed. In other words, states visited at time $t$ can't be visited at any other time.
- $\mathcal{S}$ and $\mathcal{A}$ are finite sets.

The third assumption is solely to avoid having to define a time indexed Value function. It can be completely avoided. We chose not to do this in the spirit of notational simplicity. These assumptions can be relaxed, most notably we can show similar results for the discounted and infinite horizon case. We chose to present the finite horizon proof because of the nature of our experimental results.

Let $\Phi = \mathrm{id}$ be the identity embedding so that $\mathcal{E} = \Gamma$. In this case $\mathbb{P}_\pi^\Phi$ denotes the distribution of trajectories corresponding to policy $\pi$. We define the value function $V^\pi : \mathcal{S} \to \mathbb{R}$ as

$$V^\pi(s_t = s) = \mathbb{E}_{\tau \sim \mathbb{P}_\pi^{\mathrm{id}}} \left[ \sum_{\ell=t}^{T} R(s_{\ell+1}, a_\ell, s_\ell) | s_t = s \right]$$

The Q-function $Q^\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ as:

$$Q^\pi(s_t, a_t = a) = \mathbb{E}_{\tau \sim \mathbb{P}_\pi^{\mathrm{id}}} \left[ \sum_{\ell=t}^{T} R(s_{\ell+1}, a_\ell, s_\ell) \right]$$

Similarly, the advantage function is defined as:

$$A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$$

We denote by $V(\pi) = \mathbb{E}_{\tau \sim \mathbb{P}_\pi^{\mathrm{id}}} \left[ \sum_{t=0}^{T} R(s_{t+1}, a_t, s_t) \right]$ the expected reward of policy $\pi$ and define the visitation frequency as:

$$\rho_\pi(s) = \mathbb{E}_{\tau \sim \mathbb{P}_\pi^{\mathrm{id}}} \left[ \sum_{t=0}^{T} \mathbf{1}(s_t = s) \right]$$

The first observation in this section is the following lemma:

**Lemma C.3.** *two distinct policies $\pi$ and $\tilde\pi$ can be related via the following equation :*

$$V(\tilde\pi) = V(\pi) + \sum_{s \in \mathcal{S}} \left( \rho_{\tilde\pi}(s) \left( \sum_{a \in \mathcal{A}} \tilde\pi(a|s) A^\pi(s,a) \right) \right)$$

*Proof.* Notice that $A^\pi(s,a) = \mathbb{E}_{s' \sim P(s'|a,s)} \left[ R(s',a,s) + V^\pi(s') - V^\pi(s) \right]$. Therefore:

$$\mathbb{E}_{\tau \sim \mathbb{P}_{\tilde\pi}^{\mathrm{id}}} \left[ \sum_{t=0}^{T} A_\pi(s_t, a_t) \right] = \mathbb{E}_{\tau \sim \mathbb{P}_{\tilde\pi}^{\mathrm{id}}} \left[ \sum_{t=0}^{T} R(s_{t+1}, a_t, s_t) + V^\pi(s_{t+1}) - V^\pi(s_t) \right]$$

$$= \mathbb{E}_{\tau \sim \mathbb{P}_{\tilde\pi}^{\mathrm{id}}} \left[ \sum_{t=0}^{T} R(s_{t+1}, a_t, s_t) \right] - \mathbb{E}_{s_0} \left[ V^\pi(s_0) \right]$$

$$= -V(\pi) + V(\tilde\pi)$$

The result follows. $\square$

See [33] for an alternative proof. We also consider the following linear approximation to $V$ around policy $\pi$ (see: [14]):

$$L(\tilde\pi) = V(\pi) + \sum_{s \in \mathcal{S}} \left( \rho_\pi(s) \left( \sum_{a \in \mathcal{A}} \tilde\pi(a|s) A^\pi(s,a) \right) \right)$$

Where the only difference is that $\rho_{\tilde{\pi}}$ was substituted by $\rho_{\pi}$. Consider the following embedding $\Phi^s : \Gamma \to \mathbb{R}^{|\mathcal{S}|}$ defined by $(\Phi(\tau))_s = \sum_{t=0}^{T} \mathbf{1}(s_t = s)$, and related cost function defined as: $C(\mathbf{v}, \mathbf{w}) = \|\mathbf{v} - \mathbf{w}\|_1$.

**Lemma C.4.** *The Wasserstein distance* $\mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_{\pi}^{\Phi^s})$ *is related to visit frequencies since:*

$$\mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_{\pi}^{\Phi^s}) \geq \sum_{s \in \mathcal{S}} |\rho_{\pi}(s) - \rho_{\tilde{\pi}}(s)|$$

*Proof.* Let $\Pi$ be the optimal coupling between $\mathbb{P}_{\tilde{\pi}}^{\Phi^s}$ and $\mathbb{P}_{\pi}^{\Phi^s}$. Then:

$$\mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_{\pi}^{\Phi^s}) = \mathbb{E}_{u,v \sim \Pi}\left[\|u - v\|_1\right]$$
$$= \sum_{s \in \mathcal{S}} \mathbb{E}_{u,v \sim \Pi}\left[|u_s - v_s|\right]$$

Where $u_s$ and $v_s$ denote the $s \in \mathcal{S}$ indexed entry of the $u$ and $v$ vectors respectively. Notice that for all $s \in \mathcal{S}$ the following is true:

$$\left| \underbrace{\mathbb{E}_{u \sim \mathbb{P}_{\pi}^{\Phi^s}}[u_s]}_{\rho_{\pi}(s)} - \underbrace{\mathbb{E}_{v \sim \mathbb{P}_{\pi}^{\Phi^s}}[v_s]}_{\rho_{\pi'}(s)} \right| \leq \mathbb{E}_{u,v \sim \Pi}\left[|u_s - v_s|\right]$$

The result follows.

$\square$

These observations enable us to prove an analogue of Theorem 1 from [31], namely:

**Theorem C.5.** *If* $\mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_{\pi}^{\Phi^s}) \leq \delta$ *and* $\epsilon = \max_{s,a} |A^{\pi}(s,a)|$, *then* $V(\tilde{\pi}) \geq L(\tilde{\theta}) - \delta\epsilon$.

As in [31], Theorem 5.1 implies a policy improvement guarantee for BGPG from Section 5.4.

*Proof.* Notice that:

$$V(\tilde{\pi}) - L(\tilde{\pi}) = \sum_{s \in \mathcal{S}}\left((\rho_{\tilde{\pi}}(s) - \rho_{\pi}(s))\left(\sum_{a \in \mathcal{A}} \tilde{\pi}(a|s)A^{\pi}(s,a)\right)\right)$$

Therefore by Holder inequality:

$$|V(\tilde{\pi}) - L(\tilde{\pi})| \leq \underbrace{\left(\sum_{s \in \mathcal{S}} |\rho_{\pi}(s) - \rho_{\tilde{\pi}}(s)|\right)}_{\leq \mathrm{WD}_0(\mathbb{P}_{\tilde{\pi}}^{\Phi^s}, \mathbb{P}_{\pi}^{\Phi^s}) \leq \delta} \underbrace{\left(\sup_{s \in \mathcal{S}} \left|\sum_{a \in \mathcal{A}} \tilde{\pi}(a|s)A^{\pi}(s,a)\right|\right)}_{\leq \epsilon}$$

The result follows.
$\square$

We can leverage the results of Theorem C.5 to show wasserstein trust regions methods with embedding $\Phi^s$ give a monotonically improving sequence of policies. The proof can be concluded by following the logic of Section 3 in [31].

### C.3 Off policy embeddings and their properties.

It is easy to see that if the cost function equals the $l_2$ norm between state-policy pairs and if $\mathrm{WD}_0(\mathbb{P}_{\pi}^{\Phi_S}, \mathbb{P}_{\pi'}^{\Phi_S}) = 0$ then $\mathbb{E}_{\mathbb{P}_S}[\mathbf{1}(\pi(S) \neq \pi'(S))] = 0$. If $\mathbb{P}_S$ has mass only in relevant areas of the state space, a value of zero implies the two policies behave similarly where it matters. In the case when the user may care only about the action of a policy within a set of states of interest, this notion applies.

When the sampling distribution can be identified with the stationary distribution over states of the current policy, we can recover trust region-type of results for BGPG.