
Improving Exploration in Deep Reinforcement Learning by State Planning Policies

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce an improvement for reinforcement learning (RL) algorithms for con-
2 tinuous setting called state planning policy RL (SPP-RL). In SPP-RL, the actor
3 plans for the next state provided the current state. To communicate the actor out-
4 put to the environment, we incorporate an inverse dynamics control model and
5 train it using supervised learning. We evaluate our improvement using the off-
6 policy state-of-the-art reinforcement learning algorithms: TD3 and SAC. The tar-
7 get states need to be physically relevant; the overall learning procedure is formu-
8 lated as a constrained optimization problem, solved via the classical Lagrangian
9 multipliers method. We benchmark the state planning RL approach using a set of
10 Safety-gym level 0 (no safety cost involved) environments and the AntPush env..
11 We find that SPP-RL significantly beats the baselines in terms of average return.
12 We assign the performance boost to the more efficient SPP-RL agent exploration,
13 performed in the target-state space rather than the action space. We report numer-
14 ical experiments confirming this finding.

15 1 Introduction

16 Research on reinforcement learning (RL) has brought many successful applications in diverse fields
17 of science and technology. RL application areas can be split into two classes: discrete (e.g., board
18 games) and continuous (e.g., robotic problems). Here, we are interested in continuous simulation
19 environments, mostly in robotics. The RL is concerned with training a policy governing agent
20 motion via interactions with the environment to maximize the expected total return.

21 Traditionally, RL is based on searching for the optimal policy within the space of state-action map-
22 pings; the policy is a function assigning an action to take depending on the current state. We propose
23 an improvement based on the principle of training an actor (a policy) operating entirely in the state
24 space (state-state mappings). We call such policies the state planning policies (SPP), whose ac-
25 tions determine desired trajectories in the state space. The task of training SPP may initially seem
26 infeasible due to a significantly larger dimension of states than actions. Nonetheless, quite surpris-
27 ingly, we show that the approach is feasible and often leads to significant improvements in average
28 performance and decreased sample efficiency for a class of robotic locomotion tasks.

29 We call our approach *State Planning Policy Reinforcement Learning (SPP-RL)*. It is a generic ap-
30 proach for problems specified using continuous environments. The main building block of SPP-RL –
31 the RL agent can be implemented using virtually any model-free RL algorithm. We chose to develop
32 our approach using the state-of-the-art off-policy DDPG [18], TD3 [8], and SAC [11] algorithms.
33 Note that, in SPP-RL we need another trainable model to communicate the policy output to the en-
34 vironment; as such, we incorporate a learnable inverse dynamics control model (IDM), see Fig. 1.
35 The overall algorithm optimizes the policy simultaneously with IDM. To ensure that the policy tar-
36 get states satisfy physical and under-actuation constraints, we formulate a constrained optimization

37 objective for policy training. Our work lies within the category of RL methods that have already
 38 implemented state-state policies, including work on hierarchical RL [21], the D3G algorithm [6],
 39 and behavioral cloning from observation [27]. In this work, we show the properties and advantages
 40 of state-planning policies that have not been demonstrated earlier.

41 **Summary of results.** Although the SPP-RL algorithm searches for the optimal policy within a
 42 much larger space, our performance benchmarks revealed that SPP-RL implementations often out-
 43 perform their vanilla RL counterparts. Experiments in Safety-Gym Level 0 environments [25] (with-
 44 out safety cost) demonstrate that SPP-TD3 and SPP-SAC outperform by a great margin TD3 and
 45 SAC, respectively. Experiments in AntPush task [21] show that SPP-TD3 outperforms hierarchical
 46 RL method HIRO [21] and provides some interpretability of the agent behavior.

47 We hypothesize that the superior performance of SPP-RL in the tested continuous environments
 48 originates in more efficient state-space exploration by state-state policies than traditional state-action
 49 policies; here noise is being added to target states rather than actions. To argue this, we performed
 50 series of experiments, including evaluation of a shadow agent utilizing experience from SPP and
 51 vanilla replay buffers (Sec. 5.2) and a study of the distributions of states gathered in different replay
 52 buffers (Sec. 5.4).

53 We implemented SPP-RL methods as a modular PyTorch library shared as open-source. SPP-RL
 54 algorithms are derived from their vanilla RL counterparts, making extending the library with new
 55 RL algorithms straightforward. We also share videos with test episodes of the trained agents to
 56 accompany benchmark plots [1].

57 1.1 Related work

58 We present a (non-exhaustive) list of related works; refer to Tab. 1 for a perspective on related
 59 work. The closest approach to ours is the D3G algorithm introduced by [6], which includes state
 60 planning policies, and introduces a novel form of the value function defined on state-next state pairs.
 61 There are two main ways our method is distinct. First, SPP employs the classical formulation of
 62 the value function. Also, we do not include a forward dynamics model nor the cycle loss. Instead,
 63 to guarantee consistency of the policy target-states in SPP, we formulate a constrained optimization
 64 problem (compare Fig. 2) solved via Lagrangian optimization.

65 Our work builds on the classical RL algorithms going back to REINFORCE [30], Asynchronous
 66 Actor-Critic [19], and especially the off-policy actor-critic algorithms including Q-Prop [9], DDPG
 67 [18], SAC [11, 12], and TD3 [8]. State planning policies have been used in hierarchical RL (HRL)
 68 methods like HIRO [21] and FuN [28]. Contrary to HRL SPP-RL approach does not employ a
 69 hierarchy of multiple policies nor state conditioned value functions.

70 Training predictive models (like IDMs) is fundamental for the model-based RL approach including
 71 algorithms: a locally linear latent dynamics model [29], model-based planning for discrete and
 72 continuous actions [14], model-predictive control [5], and model based policy optimization [16].
 73 We deployed IDMs for mapping current-target states to actions; other applications of IDMs in RL
 74 include the context of planning: search on replay buffer [7], episodic memory graph [31], and
 75 topological memory for navigation [26]. Existing many other applications of IDM in context of RL
 76 including: policy adaptation during deployment [13], sim to real transfer [4], adversarial exploration
 [15], curiosity-driven exploration [23], and video pre-training for minecraft [3].

Technique	State-state policy	Inverse/forward model	State cond. Q funct.	Policy Hierarchy	Planning horizon
SPP (ours)	yes	inverse	no	single policy	single step
D3G	yes	inverse & forward	yes	single policy	single step
HRL	yes(upper level)	inverse	yes(upper level)	multiple policies	multiple steps
Planning	yes	inverse	yes	N/A	multiple steps
Model based	no	forward	N/A	single policy	single step

Table 1: A Perspective on Related Work

78 **1.2 Background**

79 Following the standard setting used in RL literature, we work with infinite horizon *Markov decision*
 80 *process* (MDP) formalism $(\mathcal{S}, \mathcal{A}, P, r, \rho_0, \gamma)$, where \mathcal{S} is a *state space*, \mathcal{A} is a *action space*, $P: \mathcal{S} \times$
 81 $\mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a *transition probability distribution*, $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a *reward function*, ρ_0 is an
 82 *initial state distribution*, and $\gamma \in (0, 1)$ is a *discount factor*. From now on we assume that the MDP
 83 is fixed. In RL the agent interacts with E in discrete steps by selecting an action a_t for the state s_t
 84 at time t , causing the state transition $s_{t+1} = E(s_t, a_t)$, as a result the agent collects a scalar reward
 85 $r_{t+1}(s_t, a_t)$, the return is defined as the sum of discounted future reward $R_t = \sum_{i=t}^T \gamma^{(i-t)} r(s_i, a_i)$.
 86 The goal in RL is to learn a policy that maximizes the expected return from the start distribution.

87 **2 State Planning Policy Reinforcement Learning Approach**

88 Our SPP approach is rooted in state-state reinforcement learning, by which we mean setting in which
 89 RL agent is trained to plan goals in the state-space, the approach already employed e.g., in HRL,
 90 planning, D3G RL algorithms (see Tab. 1). In SPP a state planning policy π given the current state
 91 s_t outputs z_t – the desired target state to be reached by the environment in the next step. Forcing
 92 the environment to reach the desired state requires translating the target state to a suitable action a_t .
 93 Hence, we employ an additional model capable of mapping the current state-target state pair (s_t, z_t)
 94 to the action a_t – a (trainable) IDM model. Ideally, we like to have consistency $z_t(s_t) \approx s_{t+1}$.
 95 The consistency cannot be guaranteed a-priori, is rather achieved in SPP setting by employing a
 96 constrained optimization approach. A diagram illustrating SPP approach is presented in Fig. 1.
 97 We have freedom of choice of the particular RL algorithm (RL agent) and IDMs implementations.
 98 Currently, we use feed-forward neural networks, and RL Agent using implementations of the state-
 99 of-the-art off-policy RL algorithms: DDPG [18], TD3 [8] and SAC [11]. We present details of
 100 SPP-RL implementation in Sec. 4 using as the example SPP-DDPG. The encountered experiences
 101 during the execution of an off-policy RL algorithm are stored in replay buffer \mathcal{D} .

102 The main building block of SPP-RL are the state planning policies, intuitively a state planning policy
 103 selects a desired trajectory in the state-space of the environment.

104 **Definition 1.** We call a *state planning policy* a map $\pi_\theta: \mathcal{S} \rightarrow \mathcal{P}(\mathcal{S})$ parametrized using a vector of
 105 parameters $\theta \in \mathbb{R}^{n(\theta)}$, and we denote $\pi_\theta(z|s)$ a probability of the desired *target state* $z \in \mathcal{S}$ for the
 106 given current state $s \in \mathcal{S}$.

107 We call a *deterministic state planning policy* a parametrized map $\pi_\theta: \mathcal{S} \rightarrow \mathcal{S}$, and we denote
 108 $\pi_\theta(s) = z$.

109 We assume that π has continuous and bounded derivatives with respect to θ . We will call state
 110 planning policy whenever it is clear from the context deterministic/stochastic and omit the parameter
 111 subscript $\pi = \pi_\theta$.

112 Besides the state planning policy (Def. 1) the second main building block of the overall SPP agent
 113 is a model for mapping the current state-target pair (s_t, z_t) to suitable action a_t . Following the
 114 existing literature, we call such model the *inverse dynamics control model* (IDM), or simply the
 115 *control model*.

116 **Definition 2.** For a given MDP $(\mathcal{S}, \mathcal{A}, P, r, \rho_0, \gamma)$. Let $s, z \in \mathcal{S}$. We define the *control model*:

$$\text{CM}: \mathcal{S} \times \mathcal{S} \rightarrow \mathcal{A}, \quad \text{CM}(s, z) = a,$$

117 i.e. for the given two states CM computes the action a . We call s, z the initial state and the target state
 118 respectively, where a informally satisfies $\text{argmax}_{b \in \mathcal{S}} P(s, a, b) \sim z$ for stochastic E , or $z \approx E(s, a)$
 119 for deterministic E .

120 Obviously, in order to work, SPP requires consistency of the target states generated by the policy
 121 with the actual next-states of the environment. We call this property the state consistency property
 122 (or simply consistency) of π , refer to Fig. 2. As it may be intractable to verify SPP for all possible
 123 interactions in continuous environments, we are interested in guaranteeing the state consistency for
 124 the experiences stored in the replay buffer. It is analogous to the behavioral cloning from observation
 125 loss [27].

126 **Property 1.** Let \mathcal{D} be a replay buffer, CM be an IDM and π be a (SPP) policy. We say that π has
 127 the *state consistency* property with threshold $d > 0$ if it holds that

$$\mathbb{E}_{\substack{(s_t, s_{t+1}) \in \mathcal{D} \\ z_t \sim \pi(s_t)}}} [\|s_{t+1} - z_t\|_2^2] \leq d,$$

128 for deterministic π we have $z_t = \pi(s_t)$. Given (z_t, s_{t+1}) , we call distance $\|z_t - s_{t+1}\|_2^2$ the *state-*
 129 *consistency distance* (refer Fig. 2). We will often assume that d is known from context and omit
 130 ‘with threshold $d > 0$ ’.

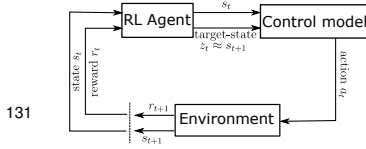


Figure 1: Diagram presenting
 our SPP method

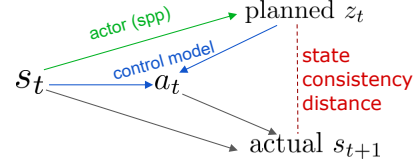


Figure 2: Diagram presenting the idea
 of state consistency property, ultimately
 we want to achieve $z_t \approx s_{t+1}$.

3 Ensuring State Consistency Property by Constrained Optimization

134 Our SPP algorithm is utilizing three parametrized models: IDM CM_ψ , policy model π_θ , and Q -
 135 function model(s) Q_ϕ . In our current implementation, all of the models are feed-forward neural
 136 networks. One way of ensuring the state consistency property (Prop. 1) is to modify the policy
 137 training loss, such that the expected values are maximized under fixed state consistency distance
 138 penalty. However, such an approach has many disadvantages, e.g., choosing appropriate learning
 139 temperature (λ) for the state consistency penalty is a very delicate issue, see the ablation study in
 140 Appendix 5.5. It is easy to notice that setting its value too low would result in π biased towards
 141 nonphysical target-states. On the other hand, setting its value too high would make π overly con-
 142 servative for off-policy states. Hence, we find a solution relying on constrained optimization more
 143 appealing for the studied problem. Namely, the objective for policy training is to maximize the sum
 144 of discounted rewards assuming a fixed threshold for the state consistency distance.

145 **Definition 3.** Let π be a state planning policy (Def. 1), CM be a control model (Def. 2), \mathcal{D} be the
 146 replay buffer with experience generated by executing an off-policy RL algorithm. In particular z_t 's
 147 are target states evaluated on-the-fly by π (susceptible to be changed during the course of algorithm),
 148 and s_{t+1} 's are E next-states. Let $d > 0$ be a fixed hyperparameter.

149 We define the *constrained objective for state planning policy* π as follows

$$\max_{\pi} \mathbb{E}_{\substack{\tau \sim \pi, CM \\ a_i = CM(s_i, z_i)}} \left[\sum_{i=0}^T \gamma^i r(s_i, a_i) \right], \quad (1a)$$

$$\text{s.t.} \quad \mathbb{E}_{\substack{(s_t, s_{t+1}) \in \mathcal{D} \\ z_t \sim \pi(s_t)}}} [\|s_{t+1} - z_t\|_2^2] \leq d, \quad (1b)$$

150 where d is a hyperparameter for determining the allowed threshold for the expected divergence
 151 of predictions from actual next-states, (1a) is an expectation over the policy trajectories generated
 152 by both of the state planning policy and IDM (trajectory is composed out of tuples (s_t, z_t, a_t)).
 153 The whole optimization process of (1a) is being performed off-policy, see Sec. 4. The con-
 154 strained objective in (1) is being solved using the standard Lagrange multiplier method. The
 155 max-min Lagrangian objective $\mathcal{L}(\pi, \lambda)$ for the constrained optimization problem takes the form
 156 $\max_{\pi} \min_{\lambda \geq 0} \mathcal{L}(\pi, \lambda) = \mathbb{E}_{\tau \sim \pi} [R_0(\pi)] - \lambda \left(\mathbb{E}_{\mathcal{D}} \left[\|s_{t+1} - z_t\|_2^2 \Big|_{z_t \sim \pi(s_t)} \right] - d \right)$. For more details
 157 refer to App. B.1.

4 Algorithm Implementation

159 We briefly present here details of the SPP Algorithm implementation, more detailed discussion
 160 can be found in Appendix B. We implemented SPP-DDPG, SPP-TD3, and SPP-SAC as a modu-

161 lar Python library within PyTorch framework [22]. All gradient optimization steps were performed
 162 using the Adam optimizer by [17]. Many of the algorithmic choices were motivated by the Spin-
 163 ning Up RL on-line resource [2]. We publish the modular SPP-RL software package as open-source
 164 [1]. For illustrative purposes, we present the pseudo-code of the full SPP-DDPG algorithm in Algo-
 165 rithm 1. SPP-SAC and SPP-TD3 algorithms are presented in Appendix B. We emphasize that SPP
 166 algorithms are not using any extra samples, i.e. the samples utilized for ICM training are added to
 167 the buffer and then reutilized for RL training, and if the buffer is full new samples are not being
 168 added anymore. An important caveat of our policy implementation in SPP-DDPG, not present in
 169 the vanilla DDPG, is that the output of π is being normalized in order to reflect the physical bounds.
 170 Contrary to vanilla DDPG where π outputs actions within well-defined uniform bounds, in SPP a
 171 suitable normalization of target state π output is being computed online – depends on the past E ob-
 172 servations. Implementation of π is a feed-forward neural network (refer to Appendix C for details)
 173 with tangential outputs bounded within $[-1, 1]$. Hence, we normalize the output of neural network
 174 π by utilizing the current mean and min/max values of the past observations in replay buffer \mathcal{D} . We
 recompute mean and min/max values after each episode of the algorithm. Our base implementation

Algorithm 1: SPP-DDPG Algorithm

input : environment E ; initial model parameters θ, ϕ, ψ ; state planning distance threshold d ; empty
 replay buffer \mathcal{D} ; the DDPG algorithm hyperparameters
output: trained model parameters θ, ϕ, ψ ; total return
repeat
 | Sample random action $a \sim \mathcal{U}$;
 | Store experience $(s_t, a_t, z_t = s_{t+1}, r_{t+1}, s_{t+1})$ in replay buffer \mathcal{D} ; (use next-state as the initial actor
 | actions)
until random exploration is done;
repeat
 | **if** buffer \mathcal{D} is not full **then**
 | | Compute actor prediction $z_t = \pi(s_t) + \varepsilon$, where $\varepsilon \sim \mathcal{N}$;
 | | Compute action $a_t = \text{CM}(s_t, z_t)$ and observe reward r_{t+1} and next state s_{t+1} ;
 | | Store experience $(s_t, z_t, a_t, s_{t+1}, r_{t+1})$ in \mathcal{D} ;
 | **end**
 | **if** it's time to update CM **then**
 | | Sample $\{b_i = \{(s_t, s_{t+1}), a\}\}_{i=1}^b$ batches of samples from replay buffer \mathcal{D} ;
 | | SGD train CM using the batches and MSE loss;
 | **end**
 | **if** it's time to update actor and critic **then**
 | | **for** update steps **do**
 | | | Randomly sample $\mathcal{B} = \{(s_t, z_t, a_t, s_{t+1}, r_{t+1})\}$ set of batches from \mathcal{D} ;
 | | | Compute $\tilde{a}_{k+1} = \text{CM}(s_{t+1}, \pi(s_{t+1}))$;
 | | | Compute targets $y = r_{t+1} + \gamma Q_{\phi_{\text{target}}}^{\pi, \text{CM}}(s_{t+1}, \tilde{a}_{k+1})$ (using target parameters ϕ_{target});
 | | | Update $\phi = \phi - \frac{l_\phi}{|\mathcal{B}|} \cdot \nabla_\phi \sum_{\mathcal{B}} (y - Q_\phi^{\pi, \text{CM}}(s_t, a_t))^2$;
 | | | Update policy parameters (ascent w.r.t θ of max-min Lagrangian obj.) $\theta =$
 | | | $\theta + l_\theta \left(\frac{1}{|\mathcal{B}|} \cdot \nabla_\theta \sum_{\mathcal{B}} Q_\phi^{\pi, \text{CM}}(s_t, a_t) \Big|_{a_t = \text{CM}(s_t, \pi_\theta(s_t))} - \frac{\lambda}{|\mathcal{B}|} \cdot \nabla_\theta \sum_{\mathcal{B}} \|s_{t+1} - \pi_\theta(s_t)\|_2^2 \right)$;
 | | | Update (descent w.r.t. λ of max-min Lagrangian obj.)
 | | | $\lambda = \lambda + l_\lambda \left(\frac{1}{|\mathcal{B}|} \sum_{\mathcal{B}} \|s_{t+1} - \pi_\theta(s_t)\|_2^2 - d \right)$;
 | | | Update actor & critic $\phi_{\text{target}} = (1 - \tau)\phi_{\text{target}} + \tau\phi$; $\theta_{\text{target}} = (1 - \tau)\theta_{\text{target}} + \tau\theta$;
 | | **end**
 | **end**
until convergence;

175 of DDPG algorithm is parametrized by the usual hyper-parameters including episode length, update
 176 batch size, Polyak averaging parameter (τ), actor and critic learning rate l , maximal episode length,
 177 number of test episodes, γ . To ensure that we perform minimization (2) within the domain of posi-
 178 tive λ values, we optimize the parameter of the softplus function. All relevant hyper-parameters are
 179 provided in Appendix C.
 180

181 5 Experimental Evaluation

182 To show the feasibility of our method, we performed experiments on a set of benchmarks using
183 continuous environments, most of them having large space dimensions. We performed all of the
184 reported experiments using the default vector state input. We show that SPP-RL implementations
185 compare favorably to their vanilla RL counterparts. As our SPP approach differs considerably from
186 the vanilla off-policy RL, we performed a thorough hyper-parameter sweep from scratch. We pro-
187 vide the hyper-parameter values from the actual SPP implementations in Appendix C. For the sake
188 of presentation we share videos with example test episodes rendered using the trained actors and
189 high-resolution benchmark plots. All of our experiments are reproducible, we share the sources,
190 training, evaluation, and trained models online [1]. We also evaluated SPP-RL in classical MuJoCo
191 tasks, and the performance is comparable to vanilla (see [1]). All of the reported experiments were
192 run using CPU only, and a single experiment was always run on a single CPU core, i.e., we have not
193 performed collecting experience in parallel. The experiments were performed on an example ma-
194 chine: AMD Ryzen Tr. 1920X, 64 Gb RAM, Ubuntu OS 18.04. Example average time of execution
195 of 10^6 steps stands at SPP-DDPG 5hrs 58', DDPG 3hrs, 7', (SPP-)SAC 19hrs 20', SAC 11hrs.

196 5.1 Safety-Gym (Locomotion Tasks)

197 We use environments from the safety-gym suite by [25]. Currently, we employed only Level 0 envi-
198 ronments (which does not involve the cost function for violating the safety). We find Level 0 tasks
199 from the safety-gym suite as the perfect ground to study the performance of SPP-RL in robotic lo-
200 comotion environments, the goal being to steer agents (robots) to solve planar goal-reaching tasks.
201 Moreover, we concentrate on difficulties arising from higher dimensionality of the state (and actions)
202 space rather than maximizing returns under safety constraints. The experiments were performed us-
203 ing solely the vector state input. We leave investigating the higher-level environments considering
204 the cost function as a topic of future research. We chose a subset of the most challenging Level
205 0 tasks, including Car-Push, Doggo-Goal, Doggo-Button environments. We also create a custom
206 environment (termed Doggo-Columns) based on Doggo-Goal with additional 10 fixed pillars placed
207 in the arena, obscuring the paths toward the goal. We evaluate our SPP-TD3 implementation against
208 state-of-the-art off-policy algorithms like TD3 and SAC. The results presented in Fig. 3 clearly
209 show that SPP-TD3 is superior to vanilla off-policy algorithms within the studied safety-gym envi-
210 ronments. Also, there is a noticeable difference in the learned behavior of the trained agents. The
211 agents trained using the SPP-RL approach show smarter and more efficient behavior; for instance,
212 the trained using SPP doggo robot learned an efficient gait of moving backward to mark a goal or
213 press a button. For comparison, we publish videos of the trained agents online [1]. We argue that
214 the performance boost exhibited by SPP-RL algorithms over vanilla counterparts is due to improved
215 exploration. In Sec. 5.2 we show results from evaluating a TD3 shadow agent, i.e. vanilla TD3
216 agent utilizing for training some portion of experience from SPP-TD3 replay buffer. In Sec 5.4 we
217 investigate differences in distribution of states collected by both of the methods.

218 5.2 More Efficient Exploration in SPP-RL

219 Our experimental evaluation using the safety-gym environments show that SPP-RL implementations
220 outperform by a great margin their vanilla off-policy RL counterparts (TD3 and SAC) in terms of the
221 average returns (See Fig.4). In this section, we argue the performance boost of SPP-RL compared
222 to the vanilla RL counterparts. Our intuition is that exploration performed in the target-state space
223 rather than in the action space may be more efficient in some cases. Exploration using SPP policies
224 results in more viable experience being collected in the replay buffer, leading to more efficient
225 Actor & Critic training. It is also possible that the constrained optimization induces some kind of
226 curriculum. To confirm the mentioned intuition, we evaluated the performance of a TD3 shadow
227 agent i.e., a vanilla TD3 Actor&Critic trained using (partially) experience collected by the SPP-TD3
228 agent. Both of the agents were trained in parallel. The TD3 shadow agent updates were performed
229 using samples drawn from two of the replay buffers. The replay buffers of the SPP-TD3 and TD3
230 agent were used according to a 50/50 ratio. The results are presented in Fig. 4. Such TD3 shadow
231 agent outperforms vanilla TD3, and eventually, its performance matches SPP-TD3 agent's in all of
232 the studied safety-gym environments, excluding Doggo-Button. We present plots of the discretized
233 distributions of states encoded using a random encoder and cross-entropy of two distributions w.r.t.
234 the quantity of gathered experience.

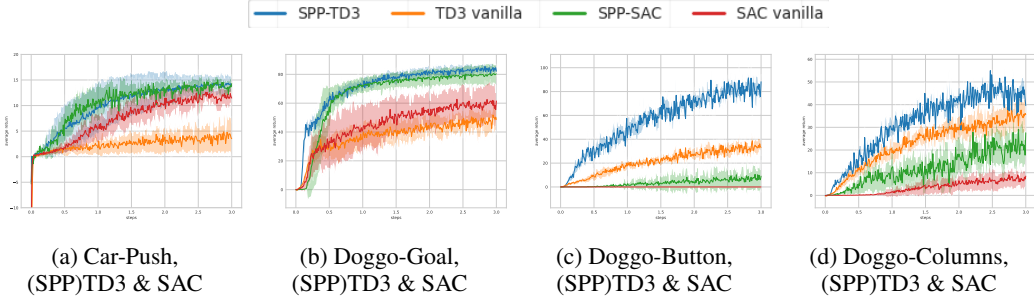


Figure 3: Experimental comparison of SPP-TD3 with corresponding vanilla off-policy RL on set of safety-gym level 0 environments. Figures show test return computed every 5k frames averaged over 10 different seeds. The continuous curve is the mean, whereas the faded color regions std. deviation. D3G did not converge (return oscillated around zero, or it diverged in CarPush to a large negative score - removed from the plot for clarity). Fine-tuning D3G to make it work in this setting is beyond scope of the research. Refer to Appendix C for exact hyperparameters that we used to perform those experiments.

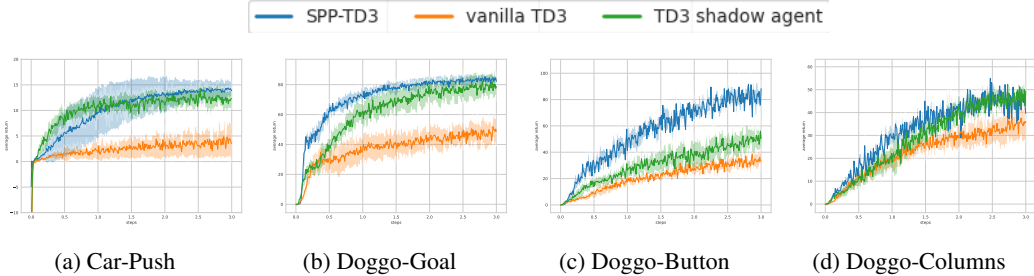


Figure 4: Experimental evaluation of (SPP)TD3 agents and a TD3 shadow agent, i.e. vanilla TD3 agent trained utilizing experience collected by a SPP-TD3 agent, on a set of safety-gym level 0 environments. Presented metrics are same as in Fig. 3.

235 5.3 Harder Exploratory Task AntPush

236 We describe an experiment in AntPush environment from [21]. The experiments were performed
 237 using solely the vector state input. The task is to control the ant such that it reaches the goal. The
 238 goal is hidden within a chamber behind a block. Therefore, Ant needs to learn to walk around
 239 the block and push it to the right first before eventually reaching the goal. Success is defined as
 240 finishing the episode within a radius 5 from the goal. We benchmark SPP-TD3 against the state-
 241 of-the-art hierarchical RL HIRO method by [21]. Specifically, we used the implementation [24].
 242 Instead of reporting the achieved success rate of a single training run like in [21], which may be
 243 spurious if a lucky seed is chosen, we report the mean and std.dev. of the AntPush success rate using
 244 10 random seeded training runs. Our experiments revealed that HIRO is highly susceptible to the
 245 random seed used. Only a single HIRO agent out of 10 trained using random seeds in total achieved
 246 a positive success rate, comparing to 7 out of 10 SPP-TD3 agents successfully learned to solve the
 247 task. The performance reported in Fig. 5a shows SPP-TD3 is eventually superior to HIRO. Example
 248 two solution paths are marked on Fig. 5b (blue curves). The right path is suboptimal as Ant blocks
 249 the entrance to the chamber where the goal is. The left path is optimal, Ant traverses to the left to
 250 push the red block away and open the passage towards the goal (green arrow).

251 Finally, Figs 5c, 5d show the obtained paths in the state space (blue dashed), and the policy target
 252 states z_t 's (orange solid), only the first two coordinates corresponding to the position of the Ant
 253 body in x, y coordinates are illustrated. Observe that in the case of the suboptimal path in Fig. 5c,
 254 the planned path diverts to the left from the actual path (blue dashed), which indicates that the agent
 255 learned and attempted the correct behavior of pushing the red brick away and successfully open the

256 entrance to the goal. In this case, however, the block is not movable; it is stuck as it was pushed
 257 forward before, hence as we see, the actual path in the state-space diverts in the middle. Figs 5c,5d
 258 show that apparently, the policy target states path being more erratic than the actual path in the state
 259 space. Erratic behavior can be mitigated by adjusting the hyperparameter d in (1b) (the smaller d ,
 260 the closer the paths will be). Nonetheless, the policy target paths (orange) in Figs 5c,5d could be
 261 potentially used to cluster agent behavior, qualitatively differentiating two example agents executing
 (sub)optimal path. This information could then be used to pick appropriate agents for deployment.

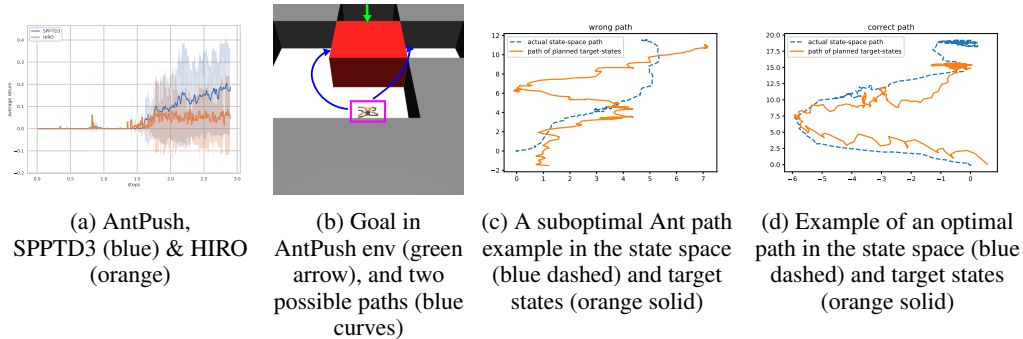


Figure 5: Experiment in AntPush environment from [21]. Fig. 5a shows success rate for SPP-TD3 & HIRO computed every 5k frames averaged over 10 different seeds (10 independent training runs were used). The continuous curve is the mean, whereas the faded color regions std. dev. Fig. 5b show two possible solution Ant paths. Figs 5c,5d show the paths in the state space (blue dashed) and target states (orange solid), the coordinates describing the position of the Ant body (x, y) are used.

262

263 5.4 SPP-RL vs Vanilla RL Replay Buffer

264 As argued in Sec. 5.2, the performance boost visible in SPP-RL vs. vanilla RL approaches is pre-
 265 sumably attributed to more efficient exploration performed by SPP-RL algorithms than vanilla RL.
 266 One empirical argument is given in Sec.5.2. Here we provide empirical evidence that the distribu-
 267 tion of observations in replay buffers gathered by SPP-RL and vanilla RL implementations differs
 268 considerably. In Fig. 6 we present an empirical study of distributions of states gathered in SPP-SAC
 269 and vanilla SAC replay buffers for the Doggo-Goal task. The state space in this task has 72 dimen-
 270 sions. Hence to make it amenable to visual investigation and entropy computation, we encode the
 271 state vectors using a random encoder. The encoder architecture that we used for this task is a simple
 272 architecture with random weights (not optimized): $72 \rightarrow 20 \tanh \rightarrow 10 \tanh \rightarrow 2$. Fig. 6 show
 273 plots of discretized state distributions gathered by example run of vanilla SAC and SPP-SAC respec-
 274 tively using the Doggo-Goal environment and encoded using the random encoder, observe that
 275 the state distributions are different, i.e., the distribution for vanilla SAC is visibly more concentrated
 276 than the one for SPP-SAC. We also compute cross-entropy to quantify the difference between those
 277 distributions as the training of both algorithms progresses and the replay buffer is filled up. Observe
 278 that the cross-entropy is increasing as the replay buffer is being filled up, suggesting that vanilla
 279 RL and SPP-RL algorithms gather different observation distributions in the replay buffer. We also
 280 performed an analogous analysis for the (SPP)TD3 approach, but it looked qualitatively similar and
 281 is not reported here.

282 5.5 Ablation Study of the Lagrangian Objective

283 Using the Doggo-Goal environment, we performed an ablation study of the SPP-TD3 the most im-
 284 portant feature. We investigate the impact of the Lagrangian objective (1) on the overall SPP-RL
 285 performance. We compare the implementation with Lagrangian objective to the implementation uti-
 286 lizing fixed λ values (parameter not trained using dual optimization), including $\lambda = 1, 0.5, 0.1, 0.01$.
 287 Performance varies greatly depending on this parameter, demonstrating how delicate the matter of
 288 choosing appropriate λ (when fixed) per given environment is. Observe that $\lambda = 1$ results in a lack
 289 of convergence, and $\lambda = 0.1$ or 0.01 results in even better performance than the SPP-TD3 imple-
 290 mentation with the Lagrangian multipliers. However, employing the Lagrange multipliers provides
 291 a natural way of solving the constrained objective optimization and avoids separate fine-tuning of λ

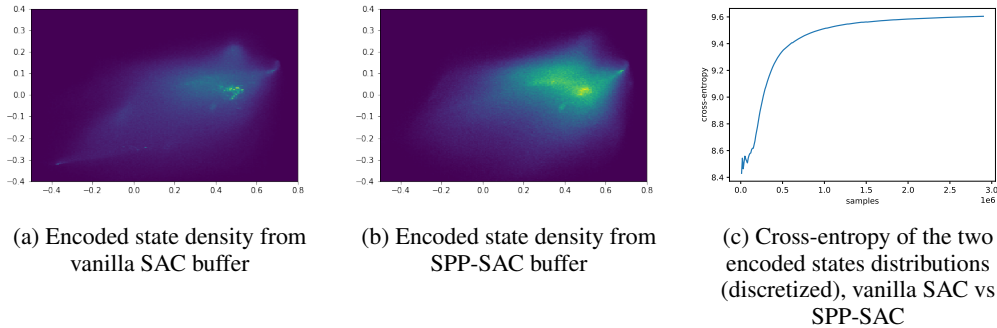


Figure 6: Visualizations of states distribution (density histogram plot) in the replay buffer at the end of training taken from single vanilla SAC, SPP-SAC runs. States are encoded in 2D using the random encoder.

292 per each environment. Moreover, the average target – next-state distance in case of $\lambda = 0.01$ is way
 293 above the set target for the Lagrangian method (0.2), whereas the Lagrangian objective successfully
 294 keeps it close to the target. We also present in Fig. 7. The study *SPP-TD3 less init. sample.*, when
 295 much less randomly generated experience is added to the replay buffer at the beginning of executing
 296 Alg. 3 (the first *repeat until* block). In some cases, like SPP-TD3 for Doggo-Goal, see Table 4, we
 297 choose to include a lot of random samples in the buffer (400k). However, this has no considerable
 298 effect on performance (Fig. 7a). The SPP-RL agent utilizing much fewer random samples (the same
 299 number as vanilla RL) has comparable performance. Other ablations that we tested included: SPP-
 300 TD3 $Q(s, s')$ critic (state-state critic), replacing the traditional target Q function computation using
 301 state-action pairs with state-next state pairs, and π not being normalized using the current mean and
 302 min/max values of the replay buffer observations. We do not show this ablations, as the algorithm
 did not converge for these settings.

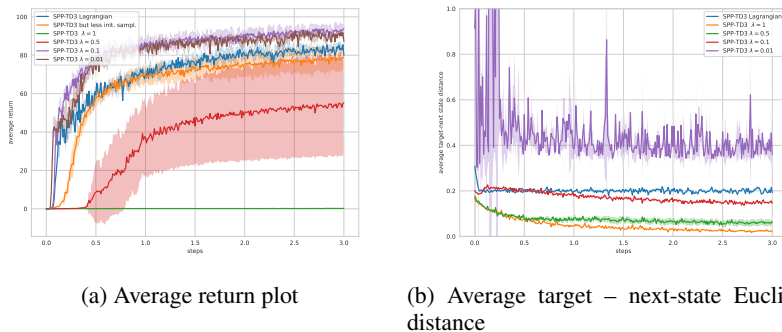


Figure 7: Ablation Study for the Lagrangian objective of SPP-TD3 algorithm, performed in the Doggo-Goal environment. The continuous curve is the mean, whereas the faded color regions std. deviation. computed from 5 independent runs.

303

304 6 Conclusions

305 We evaluated the state planning policies in reinforcement learning, where the policy selects target
 306 states for the environment. Experiments performed on continuous benchmark environments often
 307 show the superior performance of SPP-RL compared to state-of-the-art vanilla off-policy RL algo-
 308 rithms. There are various avenues for future work pertaining to this research. One path is to include
 309 in our approach physically informed control models. Another important work path is to implement a
 310 long-term policy planning method scheme and application in the safety RL setting of SPP approach.

References

- [1] Spp-rl supplementary material webpage. <https://sites.google.com/view/sppr1>, 2021.
- [2] J. Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- [3] B. Baker, I. Akkaya, P. Zhokhov, J. Huizinga, J. Tang, A. Ecoffet, B. Houghton, R. Sampedro, and J. Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos, 2022.
- [4] P. Christiano, Z. Shah, I. Mordatch, J. Schneider, T. Blackwell, J. Tobin, P. Abbeel, and W. Zaremba. Transfer from simulation to real world through learning deep inverse dynamics model, 2016.
- [5] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [6] A. Edwards, H. Sahni, R. Liu, J. Hung, A. Jain, R. Wang, A. Ecoffet, T. Miconi, C. Isbell, and J. Yosinski. Estimating $q(s,s')$ with deep deterministic dynamics gradients. In H. D. III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 2825–2835. PMLR, 13–18 Jul 2020.
- [7] B. Eysenbach, R. R. Salakhutdinov, and S. Levine. Search on the replay buffer: Bridging planning and reinforcement learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [8] S. Fujimoto, H. van Hoof, and D. Meger. Addressing Function Approximation Error in Actor-Critic Methods. *arXiv e-prints*, page arXiv:1802.09477, Feb. 2018.
- [9] S. Gu, T. Lillicrap, Z. Ghahramani, R. E. Turner, and S. Levine. Q-prop: Sample-efficient policy gradient with an off-policy critic. 2017.
- [10] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. Reinforcement learning with deep energy-based policies. 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 1352–1361, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- [11] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. 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 1861–1870, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [12] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- [13] N. Hansen, R. Jangir, Y. Sun, G. Alenyà, P. Abbeel, A. A. Efros, L. Pinto, and X. Wang. Self-supervised policy adaptation during deployment. In *International Conference on Learning Representations*, 2021.
- [14] M. Henaff, W. F. Whitney, and Y. LeCun. Model-based planning with discrete and continuous actions, 2017.
- [15] Z.-W. Hong, T.-J. Fu, T.-Y. Shann, and C.-Y. Lee. Adversarial active exploration for inverse dynamics model learning. In L. P. Kaelbling, D. Kragic, and K. Sugiura, editors, *Proceedings of the Conference on Robot Learning*, volume 100 of *Proceedings of Machine Learning Research*, pages 552–565. PMLR, 30 Oct–01 Nov 2020.

- 358 [16] M. Janner, J. Fu, M. Zhang, and S. Levine. When to trust your model: Model-based policy
359 optimization. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and
360 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran
361 Associates, Inc., 2019.
- 362 [17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2014. cite
363 arxiv:1412.6980Comment: Published as a conference paper at the 3rd International Confer-
364 ence for Learning Representations, San Diego, 2015.
- 365 [18] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
366 Continuous control with deep reinforcement learning. In Y. Bengio and Y. LeCun, editors,
367 *ICLR*, 2016.
- 368 [19] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and
369 K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In M. F. Balcan
370 and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine*
371 *Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1928–1937, New
372 York, New York, USA, 20–22 Jun 2016. PMLR.
- 373 [20] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. A. Ried-
374 miller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013.
- 375 [21] O. Nachum, S. S. Gu, H. Lee, and S. Levine. Data-efficient hierarchical reinforcement learn-
376 ing. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett,
377 editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates,
378 Inc., 2018.
- 379 [22] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
380 N. Gimeshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Te-
381 jani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative
382 style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer,
383 F. d' Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing*
384 *Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- 385 [23] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-
386 supervised prediction. In *ICML*, 2017.
- 387 [24] Z. Qin. Repository implementing hiro algorithm. [https://github.com/ziangqin-stu/
388 rl_hiro](https://github.com/ziangqin-stu/rl_hiro), 2021.
- 389 [25] A. Ray, J. Achiam, and D. Amodei. Benchmarking safe exploration in deepreinforcement
390 learning, 2019.
- 391 [26] N. Savinov, A. Dosovitskiy, and V. Koltun. Semi-parametric topological memory for naviga-
392 tion. In *International Conference on Learning Representations*, 2018.
- 393 [27] F. Torabi, G. Warnell, and P. Stone. Behavioral cloning from observation. In *Proceedings of*
394 *the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages
395 4950–4957. International Joint Conferences on Artificial Intelligence Organization, 7 2018.
- 396 [28] A. S. Vezhnevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and
397 K. Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *Proceedings*
398 *of the 34th International Conference on Machine Learning - Volume 70, ICML'17*, page
399 3540–3549. JMLR.org, 2017.
- 400 [29] M. Watter, J. Springenberg, J. Boedecker, and M. Riedmiller. Embed to control: A locally
401 linear latent dynamics model for control from raw images. In C. Cortes, N. D. Lawrence,
402 D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing*
403 *Systems 28*, pages 2746–2754. Curran Associates, Inc., 2015.
- 404 [30] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforce-
405 ment learning. *Mach. Learn.*, 8(3–4):229–256, May 1992.
- 406 [31] G. Yang, A. Zhang, A. S. Morcos, J. Pineau, P. Abbeel, and R. Calandra. Think, act and learn
407 on an episodic memory graph. In *ICLR 2020 workshop: Beyond tabula rasa in RL (BeTR-RL)*,
408 2020.