

## 408 7 APPENDIX

### 409 7.1 Proofs

#### 410 7.1.1 Proof of Lemma 3.2

411 By Definition 3.1,  $\pi_\theta(s, a)$  should satisfy the following condition:

$$\begin{aligned} \int_{T_\epsilon} q(\epsilon) d\epsilon &= \int_{T_a} \pi_\theta(s, a) da \\ &= \int_{T_\epsilon} \pi_\theta(s, g_\theta(s, \epsilon)) \cdot \left| \det\left(\frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon}\right) \right| d\epsilon \quad (\text{by change of variable}). \end{aligned}$$

412 We denote the set of continuous distributions that satisfy this condition as  $P$ . That is,

$$P_\theta(s, \epsilon) \in P \text{ if } \int_{T_\epsilon} q(\epsilon) d\epsilon = \int_{T_\epsilon} P_\theta(s, \epsilon) \cdot \left| \det\left(\frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon}\right) \right| d\epsilon.$$

413 Clearly,  $\bar{P}_\theta(s, \epsilon) = q(\epsilon) \cdot \left| \det\left(\frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon}\right) \right|^{-1}$  is an element of  $P$ , and is well defined, because of

414 Equation 3. In fact,  $\bar{P}_\theta(s, \epsilon)$  is the only element of  $P$ , as the above condition should be met for any

415  $T_\epsilon \subset \mathbb{R}^n$ . That is, if there was another distribution  $\tilde{P}_\theta(s, \epsilon) \in P$ ,

$$\int_{T_\epsilon} (\bar{P}_\theta(s, \epsilon) - \tilde{P}_\theta(s, \epsilon)) \cdot \left| \det\left(\frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon}\right) \right| d\epsilon = 0, \forall T_\epsilon \subset \mathbb{R}^n.$$

416 Since  $\tilde{P}_\theta$  is different from  $\bar{P}_\theta$ , there exists an open set  $\hat{T}_\epsilon$  where  $\tilde{P}_\theta < \bar{P}_\theta$ . Then

$$\int_{\hat{T}_\epsilon} (\bar{P}_\theta(s, \epsilon) - \tilde{P}_\theta(s, \epsilon)) \cdot \left| \det\left(\frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon}\right) \right| d\epsilon > 0,$$

417 which is a contradiction. Therefore,  $\pi_\theta(s, a) = \pi_\theta(s, g_\theta(s, \epsilon)) = \bar{P}_\theta(s, \epsilon)$ . The inverse holds  
418 trivially.

#### 419 7.1.2 Proof of Definition 5

420 First, we show that  $f$  is an injective function by showing that the determinant of Jacobian  $\frac{\partial f}{\partial a}$  is  
421 always positive:

$$\begin{aligned} \det\left(\frac{\partial f}{\partial a}\right) &= \det(I + \alpha \cdot \nabla_a^2 A_{\pi_{\bar{\theta}}}(s, a)) \\ &= \Pi_{i=1}^n (1 + \alpha \cdot \lambda_i(s, a)) > 0 \quad (\because |\alpha| < \frac{1}{\max_{(s,a)} |\lambda_1(s, a)|}), \end{aligned}$$

422 where  $\lambda_i(s, a)$  are the eigenvalues of  $\nabla_a^2 A_{\pi_{\bar{\theta}}}(s, a)$  sorted in ascending order.

423 Then, for an arbitrary open set of action  $T \subset A$ ,  $\pi_\alpha$  selects  $\tilde{T} = f(T) \subset A$  with the same probability  
424 that the original policy  $\pi_{\bar{\theta}}$  selects  $T$ .

$$\begin{aligned} \int_{\tilde{T}} \pi_\alpha(s, \tilde{a}) d\tilde{a} &= \int_T \frac{\pi_{\bar{\theta}}(s, a)}{|\det(I + \alpha \nabla_a^2 A_{\pi_{\bar{\theta}}}(s, a))|} \cdot \left| \det\left(\frac{\partial f}{\partial a}\right) \right| da \\ &= \int_T \frac{\pi_{\bar{\theta}}(s, a)}{|\det(I + \alpha \nabla_a^2 A_{\pi_{\bar{\theta}}}(s, a))|} \cdot \left| \det(I + \alpha \cdot \nabla_a^2 A_{\pi_{\bar{\theta}}}(s, a)) \right| da \\ &= \int_T \pi_{\bar{\theta}}(s, a) da. \end{aligned}$$

425 Therefore,  $\pi_\alpha(s, \cdot)$  is a valid probability distribution as  $\pi_{\bar{\theta}}(s, \cdot)$ .

### 426 7.1.3 Proof of Proposition 4.2

427 For a given state  $s$ , we can estimate the (approximate) expected value of the state under  $\pi_\alpha$  as follows.

$$\begin{aligned}
& \int_{\tilde{a}} \pi_\alpha(s, \tilde{a}) A_{\pi_\theta}(s, \tilde{a}) d\tilde{a} \\
&= \int_a \frac{\pi_\theta(s, a)}{|\det(I + \alpha \nabla_a^2 A_{\pi_\theta}(s, a))|} A_{\pi_\theta}(s, a + \alpha \nabla_a A_{\pi_\theta}(s, a)) |\det(I + \alpha \nabla_a^2 A_{\pi_\theta}(s, a))| da \\
&\approx \int_a \pi_\theta(s, a) \left[ A_{\pi_\theta}(s, a) + \alpha \|\nabla_a A_{\pi_\theta}(s, a)\|^2 \right] da \\
&= \int_a \pi_\theta(s, a) A_{\pi_\theta}(s, a) da + \alpha \int_a \pi_\theta(s, a) \|\nabla_a A_{\pi_\theta}(s, a)\|^2 da.
\end{aligned}$$

428 Therefore, we can see the following holds:

$$\begin{aligned}
L_{\pi_\theta}(\pi_\alpha) &= \int_s \rho_{\pi_\theta}(s) \int_{\tilde{a}} \pi_\alpha(s, \tilde{a}) A_{\pi_\theta}(s, \tilde{a}) d\tilde{a} \\
&\approx \eta(\pi_\theta) + \alpha \int_s \rho_{\pi_\theta}(s) \int_a \pi_\theta(s, a) \|\nabla_a A_{\pi_\theta}(s, a)\|^2 da.
\end{aligned}$$

429 Since  $\pi_\theta(s, a)$  and  $\|\nabla_a A_{\pi_\theta}(s, a)\|^2$  are positive,  $L_{\pi_\theta}(\pi_\alpha)$  is greater or equal to  $\eta(\pi_\theta)$  when  $\alpha > 0$ .  
430 On the contrary, when  $\alpha < 0$ ,  $L_{\pi_\theta}(\pi_\alpha)$  is smaller or equal to  $\eta(\pi_\theta)$ . Since this argument builds upon  
431 local approximation, it holds only when  $|\alpha| \ll 1$ .

### 432 7.1.4 Proof of Lemma 4.3

433 When  $a = g_\theta(s, \epsilon)$ ,

$$\begin{aligned}
& g_\alpha(s, \epsilon) = g_\theta(s, \epsilon) + \alpha \cdot \nabla_{g_\theta(s, \epsilon)} A_{\pi_\theta}(s, g_\theta(s, \epsilon)) \\
&\Rightarrow \frac{\partial g_\alpha(s, \epsilon)}{\partial \epsilon} = \frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon} + \alpha \cdot \nabla_{g_\theta(s, \epsilon)}^2 A_{\pi_\theta}(s, g_\theta(s, \epsilon)) \cdot \frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon} \\
&= \left( I + \alpha \cdot \nabla_{g_\theta(s, \epsilon)}^2 A_{\pi_\theta}(s, g_\theta(s, \epsilon)) \right) \cdot \frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon} \\
&\Rightarrow \det \left( \frac{\partial g_\alpha(s, \epsilon)}{\partial \epsilon} \right) = \det \left( I + \alpha \cdot \nabla_{g_\theta(s, \epsilon)}^2 A_{\pi_\theta}(s, g_\theta(s, \epsilon)) \right) \cdot \det \left( \frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon} \right) \\
&= \det \left( I + \alpha \cdot \nabla_a^2 A_{\pi_\theta}(s, a) \right) \cdot \det \left( \frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon} \right)
\end{aligned}$$

434 Since  $\left| \det \left( \frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon} \right) \right| > 0$  (Equation 3), we can divide both sides of the above equation with  
435  $\det \left( \frac{\partial g_\theta(s, \epsilon)}{\partial \epsilon} \right)$  and prove Lemma 4.3.

### 436 7.1.5 Proof of Lemma 4.4

437 Here we assume  $\alpha = 1$ , and define  $\hat{A}$  for some state-action pair  $(s_t, a_t)$  at timestep  $t$  as follows,

$$\hat{A}_{\pi_\theta}(s_t, a_t) = \frac{1}{2} \mathbb{E}_{s_t, a_t, \dots \sim \pi_\theta} \left[ \sum_{k=t}^{\infty} \gamma^k r(s_k, a_k) \right],$$

438 Then

$$\nabla_a \hat{A}_{\pi_\theta}(s_t, a_t) = \frac{1}{2} \mathbb{E}_{s_t, a_t, \dots \sim \pi_\theta} \left[ \sum_{k=t}^{\infty} \gamma^k \frac{\partial r(s_k, a_k)}{\partial a_t} \right],$$

439 which can be rewritten using  $g_\theta$  as follows.

$$\nabla_a \hat{A}_{\pi_\theta}(s_t, a_t) = \frac{1}{2} \mathbb{E}_{s_t, \epsilon_t, \dots \sim q} \left[ \sum_{k=t}^{\infty} \gamma^k \frac{\partial r(s_k, g_\theta(s_t, \epsilon_t))}{\partial g_\theta(s_t, \epsilon_t)} \right], \quad (10)$$

Then, by differentiating Equation 9 with respect to  $\theta$  after plugging in  $\hat{A}$ , we can get following relationship:

$$\begin{aligned}
\frac{\partial L(\theta)}{\partial \theta} &= \mathbb{E}_{s_0, \epsilon_0, \dots, \sim q} \left[ \sum_{t=0}^{\infty} 2 \cdot \frac{\partial g_{\theta}(s_t, \epsilon_t)}{\partial \theta} \cdot (g_{\theta}(s_t, \epsilon_t) - g_{\alpha}(s_t, \epsilon_t)) \right] \\
&= \mathbb{E}_{s_0, \epsilon_0, \dots, \sim q} \left[ \sum_{t=0}^{\infty} 2 \cdot \frac{\partial g_{\theta}(s_t, \epsilon_t)}{\partial \theta} \cdot (g_{\theta}(s_t, \epsilon_t) - (g_{\theta}(s_t, \epsilon_t) + \nabla_a \hat{A}_{\pi_{\theta}}(s_t, a_t))) \right] \\
&\quad (\because \alpha = 1 \text{ by above assumption}) \\
&= \mathbb{E}_{s_0, \epsilon_0, \dots, \sim q} \left[ \sum_{t=0}^{\infty} -2 \cdot \frac{\partial g_{\theta}(s_t, \epsilon_t)}{\partial \theta} \cdot \nabla_a \hat{A}_{\pi_{\theta}}(s_t, a_t) \right] \\
&= \mathbb{E}_{s_0, \epsilon_0, \dots, \sim q} \left[ \sum_{t=0}^{\infty} -2 \cdot \frac{\partial g_{\theta}(s_t, \epsilon_t)}{\partial \theta} \cdot \frac{1}{2} \cdot \sum_{k=t}^{\infty} \gamma^k \frac{\partial r(s_k, g_{\theta}(s_t, \epsilon_t))}{\partial g_{\theta}(s_t, \epsilon_t)} \right] \\
&\quad (\because \text{Equation 10}) \\
&= \mathbb{E}_{s_0, \epsilon_0, \dots, \sim q} \left[ \sum_{t=0}^{\infty} \sum_{k=t}^{\infty} -\gamma^k \frac{\partial g_{\theta}(s_t, \epsilon_t)}{\partial \theta} \cdot \frac{\partial r(s_k, g_{\theta}(s_t, \epsilon_t))}{\partial g_{\theta}(s_t, \epsilon_t)} \right],
\end{aligned}$$

which equals to RP gradient in Appendix 7.2.2, but with different sign. However, they turn out to be the same because we minimize  $L(\theta)$ , but maximize  $\eta(\pi_{\theta})$ . Therefore, we can say that we gain RP gradient as the first gradient when we minimize Equation 9 for particular advantage function  $\hat{A}$  and  $\alpha = 1$ .

### 7.1.6 Proof of Proposition 4.5

By Definition 5 and Lemma 4.3,

$$\begin{aligned}
\frac{\pi_{\alpha}(s, \tilde{\alpha})}{\pi_{\bar{\theta}}(s, a)} &= \frac{1}{|\det(I + \alpha \cdot \nabla_a^2 A_{\pi_{\bar{\theta}}}(s, a))|} \\
&= |\det(\frac{\partial g_{\bar{\theta}}(s, \epsilon)}{\partial \epsilon})| \cdot |\det(\frac{\partial g_{\alpha}(s, \epsilon)}{\partial \epsilon})|^{-1},
\end{aligned}$$

where  $a = g_{\theta}(s, \epsilon)$ , and thus  $\tilde{a} = g_{\alpha}(s, \epsilon)$ . Since  $\pi_{\bar{\theta}} \triangleq g_{\bar{\theta}}$ , following holds:

$$\begin{aligned}
\pi_{\alpha}(s, \tilde{\alpha}) &= \pi_{\bar{\theta}}(s, a) \cdot |\det(\frac{\partial g_{\bar{\theta}}(s, \epsilon)}{\partial \epsilon})| \cdot |\det(\frac{\partial g_{\alpha}(s, \epsilon)}{\partial \epsilon})|^{-1} \\
&= q(\epsilon) \cdot |\det(\frac{\partial g_{\bar{\theta}}(s, \epsilon)}{\partial \epsilon})|^{-1} \cdot |\det(\frac{\partial g_{\bar{\theta}}(s, \epsilon)}{\partial \epsilon})| \cdot |\det(\frac{\partial g_{\alpha}(s, \epsilon)}{\partial \epsilon})|^{-1} \\
&\quad (\because \text{Lemma 3.2}) \\
&= q(\epsilon) \cdot |\det(\frac{\partial g_{\alpha}(s, \epsilon)}{\partial \epsilon})|^{-1},
\end{aligned}$$

which implies  $\pi_{\alpha} \triangleq g_{\alpha}$ .

## 7.2 Formulations

### 7.2.1 Analytic Gradient of Generalized Advantage Estimator (GAE)

GAE [Schulman et al., 2015b] has been widely used in many RL implementations [Schulman et al., 2017, Raffin et al., 2021, Makoviichuk and Makovychuk, 2022] to estimate advantages. GAE finds a balance between variance and bias of the advantage estimation by computing the exponentially-weighted average of the TD residual terms ( $\delta_t^V$ ) Sutton et al. [1998], which are defined as follows:

$$\delta_t^V = r_t + \gamma V(s_{t+1}) - V(s_t).$$

456 Then, GAE can be formulated as follows:

$$A_t^{GAE}(\gamma, \lambda) = \sum_{k=0}^{\infty} (\gamma\lambda)^k \delta_{t+k}^V. \quad (11)$$

457 We can compute the gradients for these terms as:

$$\begin{aligned} \frac{\partial A_t^{GAE}}{\partial a_t} &= \sum_{k=0}^{\infty} (\gamma\lambda)^k \frac{\partial \delta_{t+k}^V}{\partial a_t} \\ &= \left(\frac{1}{\gamma\lambda}\right)^t \sum_{k=0}^{\infty} (\gamma\lambda)^{t+k} \frac{\partial \delta_{t+k}^V}{\partial a_t} \\ &= \left(\frac{1}{\gamma\lambda}\right)^t \left[ \sum_{k=0}^{\infty} (\gamma\lambda)^{t+k} \frac{\partial \delta_{t+k}^V}{\partial a_t} + \sum_{k=0}^{t-1} (\gamma\lambda)^k \frac{\partial \delta_k^V}{\partial a_t} \right] \\ &(\because \frac{\partial \delta_k^V}{\partial a_t} = 0 \text{ for } k < t) \\ &= \left(\frac{1}{\gamma\lambda}\right)^t \sum_{k=0}^{\infty} (\gamma\lambda)^k \frac{\partial \delta_k^V}{\partial a_t} \\ &= \left(\frac{1}{\gamma\lambda}\right)^t \frac{\partial A_0^{GAE}}{\partial a_t}. \end{aligned} \quad (12)$$

458 Based on this relationship, we can compute every  $\frac{\partial A_t^{GAE}}{\partial a_t}$  with only one backpropagation of  $A_0^{GAE}$ ,  
459 rather than backpropagating for every  $A_t^{GAE}$ .

### 460 7.2.2 RP Gradient Formulation

461 To differentiate Equation 1 with respect to  $\theta$ , we first rewrite Equation 1 as follows.

$$\begin{aligned} \eta(\pi_\theta) &= \mathbb{E}_{s_0, a_0, \dots \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \\ &= \mathbb{E}_{s_0, \epsilon_0, \dots \sim q} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, g_\theta(s_t, \epsilon_t)) \right]. \end{aligned}$$

462 Then, we can compute RP gradient as follows.

$$\begin{aligned} \frac{\partial \eta(\pi_\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \mathbb{E}_{s_0, \epsilon_0, \dots \sim q} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, g_\theta(s_t, \epsilon_t)) \right] \\ &= \mathbb{E}_{s_0, \epsilon_0, \dots \sim q} \left[ \frac{\partial}{\partial \theta} \sum_{t=0}^{\infty} \gamma^t r(s_t, g_\theta(s_t, \epsilon_t)) \right] \\ &= \mathbb{E}_{s_0, \epsilon_0, \dots \sim q} \left[ \sum_{t=0}^{\infty} \frac{\partial g_\theta(s_t, \epsilon_t)}{\partial \theta} \sum_{k=t}^{\infty} \gamma^k \frac{\partial r(s_k, g_\theta(s_k, \epsilon_k))}{\partial g_\theta(s_t, \epsilon_t)} \right] \\ &(\because \frac{\partial r(s_k, g_\theta(s_k, \epsilon_k))}{\partial g_\theta(s_t, \epsilon_t)} \neq 0 \text{ only when } k \geq t.) \\ &= \mathbb{E}_{s_0, \epsilon_0, \dots \sim q} \left[ \sum_{t=0}^{\infty} \sum_{k=t}^{\infty} \gamma^k \frac{\partial g_\theta(s_t, \epsilon_t)}{\partial \theta} \frac{\partial r(s_k, g_\theta(s_k, \epsilon_k))}{\partial g_\theta(s_t, \epsilon_t)} \right]. \end{aligned}$$

463 We can estimate this RP gradient using Monte Carlo sampling, and use it for gradient ascent to  
464 maximize Equation 1.

### 7.2.3 Estimator for Equation 4

Since  $\eta(\pi_{\bar{\theta}})$  does not depend on  $\theta$ , we can ignore the term and rewrite our loss function in Equation 4 using expectation as follows:

$$\begin{aligned} L_{\pi_{\bar{\theta}}}(\pi_{\theta}) &= \int_s \rho_{\pi_{\bar{\theta}}}(s) \int_a \pi_{\theta}(s, a) A_{\pi_{\bar{\theta}}}(s, a) \\ &= \mathbb{E}_{s \sim \rho_{\pi_{\bar{\theta}}}, a \sim \pi_{\theta}(s, \cdot)} [A_{\pi_{\bar{\theta}}}(s, a)]. \end{aligned}$$

However, note that we did not collect trajectories, or experience buffer, using  $\pi_{\theta}$ . Therefore, we use another importance sampling function  $q(s, a)$  [Schulman et al., 2015a].

$$L_{\pi_{\bar{\theta}}}(\pi_{\theta}) = \mathbb{E}_{s \sim \rho_{\pi_{\bar{\theta}}}, a \sim q(s, \cdot)} \left[ \frac{\pi_{\theta}(s, a)}{q(s, a)} A_{\pi_{\bar{\theta}}}(s, a) \right].$$

Then, we can estimate this value using Monte Carlo estimation as follows, where  $B$  is the experience buffer of size  $N$ , which stores state ( $s_i$ ), action ( $a_i$ ), and their corresponding advantage value ( $A_{\pi_{\bar{\theta}}}(s_i, a_i)$ ) obtained by following policy  $\pi_{\bar{\theta}}$ :

$$L_{\pi_{\bar{\theta}}}(\pi_{\theta}) \approx \frac{1}{N} \sum_{i=1}^N \frac{\pi_{\theta}(s_i, a_i)}{\pi_{\bar{\theta}}(s_i, a_i)} A_{\pi_{\bar{\theta}}}(s_i, a_i). \quad (13)$$

Note that we used  $\pi_{\bar{\theta}}$  in the place of  $q$ , as it is the most natural importance sampling function we can use [Schulman et al., 2015a, 2017].

## 7.3 Algorithm

Here we present the entire pseudocode of our algorithm. Before that, we'd like to discuss additional aspects of our algorithm other than those discussed in Section 4.3.

### 7.3.1 Out-of-range-ratio

As discussed in Section 2, PPO achieves its stable learning by updating policy only in the trust region near the current policy. Note that TRPO [Schulman et al., 2015a] achieves this by computing the KL divergence between the current policy and the original policy, and restraining it below certain threshold. However, PPO [Schulman et al., 2017] uses a much simpler approach, to restrict the ratio of probabilities  $\pi_{\theta}(s_i, a_i)$  and  $\pi_{\bar{\theta}}(s_i, a_i)$  for every state-action pair  $(s_i, a_i)$  in our buffer as follows,

$$1 - \epsilon_{clip} < \frac{\pi_{\theta}(s_i, a_i)}{\pi_{\bar{\theta}}(s_i, a_i)} < 1 + \epsilon_{clip},$$

where  $\epsilon_{clip}$  is a constant. To be specific, PPO uses following surrogate loss function to mandate this:

$$L_{PPO}(\theta) = \frac{1}{N} \sum_{i=1}^N \min\left(\frac{\pi_{\theta}(s_i, a_i)}{\pi_{\bar{\theta}}(s_i, a_i)} A_{\pi_{\bar{\theta}}}(s_i, a_i), \text{clip}\left(\frac{\pi_{\theta}(s_i, a_i)}{\pi_{\bar{\theta}}(s_i, a_i)}, 1 - \epsilon_{clip}, 1 + \epsilon_{clip}\right) A_{\pi_{\bar{\theta}}}(s_i, a_i)\right) \quad (14)$$

Note that some of the ratios could go out of bounds while optimizing this loss function (see Schulman et al. [2017] for details). However, the spirit of PPO lies in restricting this ratio.

As we discussed in Section 4.3, we first update our policy towards  $\pi_{\alpha}$  by minimizing Equation 9 before performing PPO update. If  $\alpha$  is sufficiently small, most of the ratios  $\frac{\pi_{\theta}(s_i, a_i)}{\pi_{\bar{\theta}}(s_i, a_i)}$  will stay near 1, which does not harm the PPO's assumption. However, if  $\alpha$  is big, many of the ratios could go out of bounds. Therefore, we compute the out-of-range-ratio after updating our policy to  $\pi_{\alpha}$  as follows,

$$\text{out-of-range-ratio} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\left(\left|\frac{\pi_{\theta}(s_i, a_i)}{\pi_{\bar{\theta}}(s_i, a_i)} - 1\right| > \epsilon_{clip}\right), \quad (15)$$

where  $N$  is the size of the buffer, and  $\mathbb{I}$  is the indicator function. Note that this ratio depends on  $\epsilon_{clip}$ . If this out-of-range-ratio is larger than a pre-defined threshold, we decrease  $\alpha$ .

### 7.3.2 Importance sampling function

After updating our policy to  $\pi_\alpha$ , if we do not take additional measure and maximize surrogate loss function in Equation 14, the updated policy can totally ignore  $\pi_\alpha$ , which is undesirable. To avoid such situation, we use another importance sampling function  $\pi_h$  defined as follows:

$$\pi_h(s, a) = \frac{1}{2}(\pi_\theta(s, a) + \pi_\alpha(s, a)).$$

By using this importance sampling function and making our updated function to stay near  $\pi_h$ , we can guarantee that the updated function stays close to  $\pi_\alpha$ . Then our final surrogate loss function is

$$L_{GIPPO}(\theta) = \frac{1}{N} \sum_{i=1}^N \min\left(\frac{\pi_\theta(s_i, a_i)}{\pi_h(s_i, a_i)} A_{\pi_\theta}(s_i, a_i), \text{clip}\left(\frac{\pi_\theta(s_i, a_i)}{\pi_h(s_i, a_i)}, 1 - \epsilon_{clip}, 1 + \epsilon_{clip}\right) A_{\pi_\theta}(s_i, a_i)\right). \quad (16)$$

See how it differs from the original loss function in Equation 14.

### 7.3.3 Pseudocode

In Algorithm 1, we present pseudocode that illustrates the outline of our algorithm, GI-PPO. Note that we used three criteria (variance, bias, out-of-range-ratio) to control  $\alpha$  after we update our policy to the  $\alpha$ -policy. If there is no reason to reduce  $\alpha$ , we increase it, because we can expect higher expected return with bigger  $\alpha$  as shown in Proposition 4.2. Also note that there are five hyperparameters here, which are  $\alpha_0, \beta, \delta_{det}, \delta_{oorr}$ , and  $\max(\alpha)$ . We present specific values for these hyperparameters for our experiments in Appendix 7.4.

---

#### Algorithm 1 GI-PPO

---

```

 $\alpha \leftarrow \alpha_0$ , Initial value
 $\beta \leftarrow$  Constant multiplier larger than 1 for  $\alpha$ 
 $\delta_{det}, \delta_{oorr} \leftarrow$  Constant thresholds
 $B \leftarrow$  Experience buffer
while Training not ended do
    Clear  $B$ 
    while Not collected enough experience do
        | Collect experience  $\{s_t, \epsilon_t, a_t, r_t, s_{t+1}\} \rightarrow B$ 
    end
    Estimate advantage  $A$  for every  $(s_i, a_i)$  in  $B$  using Eq. 11
    Estimate advantage gradient  $\frac{\partial A}{\partial a}$  for every  $(s_i, a_i)$  in  $B$  using Eq. 12
    For current  $\alpha$ , approximate  $\alpha$ -policy by minimizing loss in Equation 9
    // Variance
    For each  $\epsilon_i$  and its corresponding state-action pair  $(s_i, a_i)$ , estimate  $\det(I + \alpha \cdot \nabla_a^2 A_{\pi_\theta}(s, a))$  by
    Lemma 4.3 and get its sample minimum ( $\psi_{min}$ ) and maximum ( $\psi_{max}$ )
    // Bias
    Evaluate expected additional return in Equation 13 with our current policy to get  $R_\alpha$ 
    // Out-of-range-ratio
    Evaluate out-of-range-ratio in Equation 15 with our current policy to get  $R_{oorr}$ 
    if  $\psi_{min} < 1 - \delta_{det}$  or  $\psi_{max} > 1 + \delta_{det}$  or  $R_\alpha < 0$  or  $R_{oorr} > \delta_{oorr}$  then
        |  $\alpha = \alpha / \beta$ 
    end
    else
        |  $\alpha = \alpha \times \beta$ 
    end
     $\alpha = \text{clip}(\alpha, 0, \max(\alpha))$ 
    Do PPO update by maximizing the surrogate loss in Equation 16
end

```

---

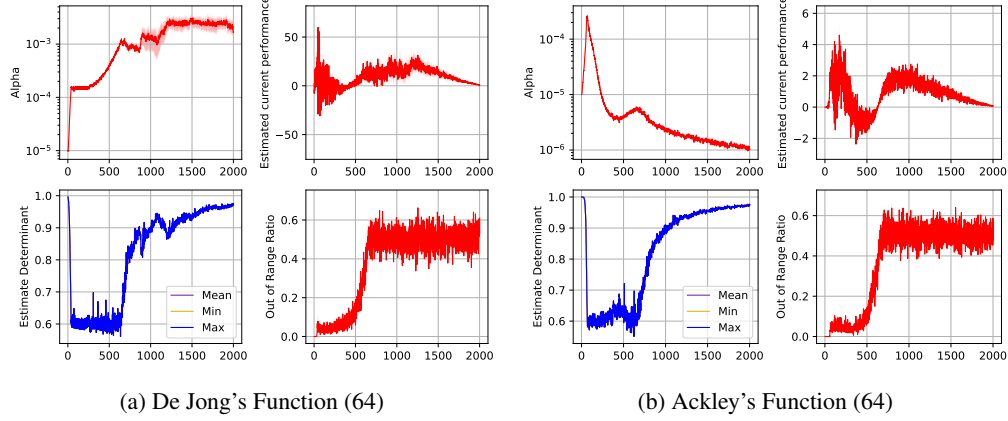


Figure 6: Change of  $\alpha$  during solving function optimization problems in Section 5.1. Since these problems are one-step problems, and observations are all the same, there is no difference between mean, min, and max of estimated determinant.

### 7.3.4 Example: Function Optimization Problems (Section 5.1)

Here we present how  $\alpha$  changes as we solve function optimization problems in Section 5.1 with Figure 6. There are 4 plots for each of the problems.

- Alpha: Shows change of  $\alpha$  over training epoch.
- Estimated current performance: Shows expected additional return, which corresponds to Equation 13 ( $R_\alpha$  in Algorithm 1).
- Estimate Determinant: Shows statistics of estimated  $\det(I + \alpha \cdot \nabla_a^2 A_{\pi_\theta}(s, a))$ .
- Out of Range Ratio: Shows out-of-range-ratio in Equation 15.

For these problems, we have set  $\alpha_0 = 10^{-5}$ ,  $\delta_{det} = 0.4$ , and  $\delta_{orr} = 0.5$ .

For De Jong's function (Figure 6a), we can observe that  $\alpha$  shows steady increase over time, but it is mainly upper bounded by variance criterion and out-of-range-ratio. In the end,  $\alpha$  stabilizes around  $10^{-3}$ . In contrast, for Ackley's function (Figure 6b) we can see that  $\alpha$  increases rapidly at first, but it soon decreases due to the variance and bias criteria. Compare these graphs with optimization curve in Figure 3d. Then we would be able to observe that this large  $\alpha$  contributes to the faster convergence of GI-PPO compared to PPO at the early stage. However, after that,  $\alpha$  decreases slowly due to the out-of-range-ratio. In the end,  $\alpha$  reaches approximately  $10^{-6}$ , which is much lower than  $10^{-3}$  of De Jong's function. These observations align well with our intuition that RP gradient would play bigger role for De Jong's function than Ackley's function, because De Jong's function exhibits RP gradients with much lower variance than those of Ackley's.

## 7.4 Experimental Details

### 7.4.1 Baseline Methods

Here we explain implementation details of the baseline methods that were used for comparisons in Section 5. First, we'd like to point out that we used relatively a short time horizon to collect experience, and used the critic network for bootstrapping, instead of collecting experience for the whole trajectory for all the methods. Therefore, the collected experience could be biased. However, we chose to implement in this way, because it allows faster learning than collecting the whole trajectory in terms of number of simulation steps.

534 **LR** We can estimate original LR gradient as follows, using log-derivative trick [Williams and Peng,  
535 1989, Glynn, 1990]

$$\begin{aligned}\frac{\partial \eta(\pi_\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \mathbb{E}_{s_0, a_0, \dots \sim \pi_\theta} \left[ L(s_0, a_0, \dots) \sum_{t=0}^{\infty} \log \pi_\theta(s_t, a_t) \right] \\ &= \frac{\partial}{\partial \theta} \mathbb{E}_{s_0, a_0, \dots \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} L(s_t, a_t, \dots) \log \pi_\theta(s_t, a_t) \right],\end{aligned}$$

536 where  $L(s_0, a_0, \dots) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$  is the discounted accumulated return. However, this LR  
537 gradient often suffers from high variance, because  $L$  varies a lot from trajectory to trajectory. To  
538 reduce such variance and faithfully compare our method to LR gradient-based method, we decided to  
539 use advantage function in the place of  $L$ , and particularly used GAE [Schulman et al., 2015b] as the  
540 advantage function as we used for PPO. After calculating this LR gradient, we take a gradient ascent  
541 step with pre-defined learning rate.

542 **RP** By using RP gradient in Appendix 7.2.2, we can perform gradient ascent as we did in LR. Also,  
543 because we use short time horizon and critic network, this method is very similar to SHAC [Xu et al.,  
544 2022]. Please refer to it for more details.

545 **PPO** Our PPO implementation is based on that of RL Games [Makoviichuk and Makoviychuk,  
546 2021]. However, to be as fair as possible, we used the same critic network, instead of using that  
547 implemented already. Also, we omitted several other additional losses that are not vital to PPO’s  
548 formulation, such as entropy loss.

549 **LR+RP** After we compute LR ( $\nabla_{LR}$ ) and RP ( $\nabla_{RP}$ ) gradient as shown above, we can interpolate  
550 them using their sample variance as follows [Parmas et al., 2018].

$$\begin{aligned}\nabla_{LR+RP} &= \nabla_{LR} \cdot \kappa_{LR} + \nabla_{RP} \cdot (1 - \kappa_{LR}), \\ \kappa_{LR} &= \frac{\sigma_{RP}^2}{\sigma_{RP}^2 + \sigma_{LR}^2},\end{aligned}$$

551 where  $\sigma_{RP}$  and  $\sigma_{LR}$  are sample standard deviation of RP and LR gradients. We gain these terms by  
552 computing trace of covariance matrix of the sample gradients from different trajectories.

553 Since we have to compute this sample statistics, we have to do multiple different backpropagations  
554 for different trajectories, which incur a lot of computation time. Also, we found out that computing  
555 covariance matrix is also time consuming when controller has a large number of parameters. There-  
556 fore, we decided to use only limited number of sample gradients (16) to compute sample variance,  
557 and also truncate the gradient to smaller length (512) to facilitate computation.

558 **PE** We tried to faithfully re-implement policy enhancement scheme of [Qiao et al., 2021].

#### 559 7.4.2 Network architecture and Hyperparameters

560 In this section, we provide network architectures and hyperparameters that we used for experiments  
561 in Section 5. For each of the experiments, we used the same network architectures, the same length  
562 of time horizons before policy update, and the same optimization procedure for critic updates, etc.  
563 We present these common settings first for each of the problems.

564 For GI-PPO, there are hyperparameters for update towards  $\alpha$ -policy, and those for PPO update. We  
565 denote the hyperparameters for  $\alpha$ -policy update by appending ( $\alpha$ ), and those for PPO update by  
566 appending (PPO). For the definition of hyperparameters for  $\alpha$ -policy update, please see Algorithm 1  
567 for details.

568 **Function Optimization Problems (Section 5.1)** For these problems, common settings are as  
569 follows.

- 570 • Actor Network: MLP with [32, 32] layers and ELU activation function



- 571 • Critic Network: MLP with [32, 32] layers and ELU activation function
  - 572 • Critic Hyperparameters: Learning rate =  $10^{-3}$ , Iterations = 16, Batch Size = 4
  - 573 • Number of parallel environments: 64
  - 574 • Horizon Length: 1
  - 575 •  $\gamma$  (Discount factor): 0.99
  - 576 •  $\tau$  (GAE): 0.95
- 577 Hyperparameters for **LR** are as follows.
- 578 • Learning Rate: Dejong(1), Dejong(64) =  $10^{-3}$ , Ackley(1) =  $10^{-4}$ , Ackley(64) =  $3 \cdot 10^{-4}$
  - 579 • Learning Rate Scheduler: Linear<sup>1</sup>
- 580 Hyperparameters for **RP** are as follows.
- 581 • Learning Rate: Dejong(1), Dejong(64) =  $10^{-2}$ , Ackley(1), Ackley(64) =  $10^{-3}$
  - 582 • Learning Rate Scheduler: Linear
- 583 Hyperparameters for **LR+RP** are as follows.
- 584 • Learning Rate: Dejong(1), Dejong(64) =  $10^{-3}$ , Ackley(1) =  $10^{-4}$ , Ackley(64) =  $3 \cdot 10^{-4}$
  - 585 • Learning Rate Scheduler: Linear
- 586 Hyperparameters for **PPO** are as follows.
- 587 • Learning Rate: Dejong(1), Ackley(1) =  $10^{-4}$ , Dejong(64), Ackley(64) =  $10^{-2}$
  - 588 • Learning Rate Scheduler: Constant
  - 589 • Batch Size for Actor Update: 64
  - 590 • Number of Epochs for Actor Update: 5
  - 591 •  $\epsilon_{clip}$ : 0.2
- 592 Hyperparameters for **GI-PPO** are as follows.
- 593 • ( $\alpha$ ) Learning Rate:  $10^{-3}$
  - 594 • ( $\alpha$ ) Batch Size for Actor Update: 64
  - 595 • ( $\alpha$ ) Number of Epochs for Actor Update: 16
  - 596 • ( $\alpha$ )  $\alpha_0$ :  $10^{-5}$
  - 597 • ( $\alpha$ )  $\max(\alpha)$ : 1.0
  - 598 • ( $\alpha$ )  $\beta$ : 1.1
  - 599 • ( $\alpha$ )  $\delta_{det}$ : 0.4
  - 600 • ( $\alpha$ )  $\delta_{orr}$ : 0.5
  - 601 • (PPO) Learning Rate: Dejong(1), Ackley(1) =  $10^{-4}$ , Dejong(64), Ackley(64) =  $10^{-2}$
  - 602 • (PPO) Learning Rate Scheduler: Constant
  - 603 • (PPO) Batch Size for Actor Update: 64
  - 604 • (PPO) Number of Epochs for Actor Update: 5
  - 605 • (PPO)  $\epsilon_{clip}$ : 0.2

---

<sup>1</sup>Learning rate decreases linearly to the minimum value as learning progresses.

606 **Differentiable Physics Problems (Section 5.2)** For these problems, common settings are as  
607 follows.

- 608 • Actor Network: MLP with [64, 64] (Cartpole), and [128, 64, 32] (Ant, Hopper) layers and  
609 ELU activation function
- 610 • Critic Network: MLP with [64, 64] layers and ELU activation function
- 611 • Critic Hyperparameters: Learning rate =  $10^{-3}$ (Cartpole),  $2 \cdot 10^{-3}$ (Ant), and  $2 \cdot 10^{-4}$ (Hopper),  
612 Iterations = 16, Batch Size = 4
- 613 • Number of parallel environments: 64(Cartpole, Ant), 256(Hopper)
- 614 • Horizon Length: 32
- 615 •  $\gamma$  (Discount factor): 0.99
- 616 •  $\tau$  (GAE): 0.95

617 Hyperparameters for **LR** are as follows.

- 618 • Learning Rate:  $10^{-4}$
- 619 • Learning Rate Scheduler: Linear

620 Hyperparameters for **RP** are as follows.

- 621 • Learning Rate: Cartpole =  $10^{-2}$ , Ant, Hopper =  $2 \cdot 10^{-3}$
- 622 • Learning Rate Scheduler: Linear

623 Hyperparameters for **PPO** are as follows.

- 624 • Learning Rate: Cartpole =  $3 \cdot 10^{-4}$ , Ant, Hopper =  $10^{-4}$
- 625 • Learning Rate Scheduler: Cartpole = Adaptive<sup>2</sup>, Ant, Hopper = Constant
- 626 • Batch Size for Actor Update: 2048
- 627 • Number of Epochs for Actor Update: 5
- 628 •  $\epsilon_{clip}$ : 0.2

629 Hyperparameters for **GI-PPO** are as follows.

- 630 • ( $\alpha$ ) Learning Rate: Cartpole =  $10^{-2}$ , Ant =  $5 \cdot 10^{-4}$ , Hopper =  $5 \cdot 10^{-3}$
- 631 • ( $\alpha$ ) Batch Size for Actor Update: Cartpole, Ant = 2048, Hopper = 8192
- 632 • ( $\alpha$ ) Number of Epochs for Actor Update: 16
- 633 • ( $\alpha$ )  $\alpha_0$ : Cartpole, Ant =  $5 \cdot 10^{-1}$ , Hopper =  $5 \cdot 10^{-3}$
- 634 • ( $\alpha$ )  $\max(\alpha)$ : 1.0
- 635 • ( $\alpha$ )  $\beta$ : 1.02
- 636 • ( $\alpha$ )  $\delta_{det}$ : 0.4
- 637 • ( $\alpha$ )  $\delta_{orr}$ : Cartpole, Hopper = 0.75, Ant = 0.5
- 638 • (PPO) Learning Rate: Cartpole =  $3 \cdot 10^{-4}$ , Ant, Hopper =  $10^{-4}$
- 639 • (PPO) Learning Rate Scheduler: Cartpole = Adaptive, Ant, Hopper = Constant
- 640 • (PPO) Batch Size for Actor Update: 2048
- 641 • (PPO) Number of Epochs for Actor Update: 5
- 642 • (PPO)  $\epsilon_{clip}$ : 0.2

---

<sup>2</sup>Learning rate is adaptively controlled, so that the KL divergence between the updated policy and the original policy is maintained at certain value, 0.008 in this case.

643 **Traffic Problems (Section 5.3)** For these problems, common settings are as follows.

- 644 • Actor Network: MLP with [512, 64, 64] layers and ELU activation function
- 645 • Critic Network: MLP with [64, 64] layers and ELU activation function
- 646 • Critic Hyperparameters: Learning rate =  $10^{-3}$ , Iterations = 16, Batch Size = 4
- 647 • Number of parallel environments: 64
- 648 • Horizon Length: 32
- 649 •  $\gamma$  (Discount factor): 0.99
- 650 •  $\tau$  (GAE): 0.95

651 Hyperparameters for **LR** are as follows.

- 652 • Learning Rate:  $3 \cdot 10^{-4}$
- 653 • Learning Rate Scheduler: Linear

654 Hyperparameters for **RP** are as follows.

- 655 • Learning Rate:  $10^{-3}$
- 656 • Learning Rate Scheduler: Linear

657 Hyperparameters for **LR+RP** are as follows.

- 658 • Learning Rate:  $3 \cdot 10^{-4}$
- 659 • Learning Rate Scheduler: Linear

660 Hyperparameters for **PPO** are as follows.

- 661 • Learning Rate:  $3 \cdot 10^{-4}$
- 662 • Learning Rate Scheduler: Constant
- 663 • Batch Size for Actor Update: 2048
- 664 • Number of Epochs for Actor Update: 5
- 665 •  $\epsilon_{clip}$ : 0.2

666 Hyperparameters for **GI-PPO** are as follows.

- 667 • ( $\alpha$ ) Learning Rate:  $10^{-5}$
- 668 • ( $\alpha$ ) Batch Size for Actor Update: 2048
- 669 • ( $\alpha$ ) Number of Epochs for Actor Update: 16
- 670 • ( $\alpha$ )  $\alpha_0$ :  $10^{-1}$
- 671 • ( $\alpha$ )  $\max(\alpha)$ : 1.0
- 672 • ( $\alpha$ )  $\beta$ : 1.1
- 673 • ( $\alpha$ )  $\delta_{det}$ : 0.4
- 674 • ( $\alpha$ )  $\delta_{orr}$ : 0.5
- 675 • (PPO) Learning Rate:  $3 \cdot 10^{-4}$
- 676 • (PPO) Learning Rate Scheduler: Constant
- 677 • (PPO) Batch Size for Actor Update: 2048
- 678 • (PPO) Number of Epochs for Actor Update: 5
- 679 • (PPO)  $\epsilon_{clip}$ : 0.2

## 680 7.5 Problem Definitions

681 Here we present details about the problems we suggested in Section 5.

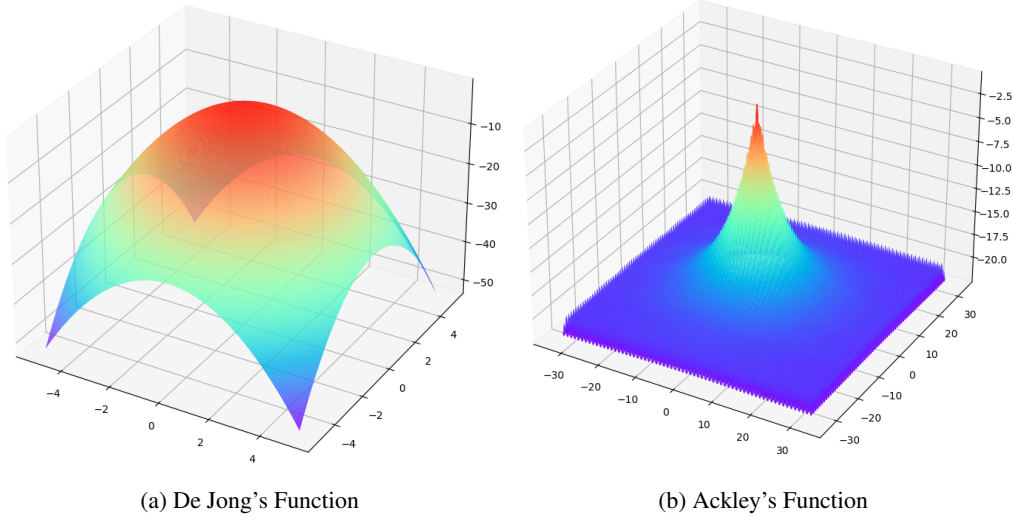


Figure 7: Landscape of target functions in 2 dimensions.

### 682 7.5.1 Function Optimization Problems (Section 5.1)

683 (N-dimensional) De Jong's function ( $F_D$ ) and Ackley's function ( $F_A$ ) are defined for  $n$ -dimensional  
 684 vector  $x$  and are formulated as follows [Molga and Smutnicki, 2005]:

$$F_D(x) = \sum_{i=1}^n x_i^2,$$

$$F_A(x) = -20 \cdot \exp\left(-0.2 \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}\right) - \exp\left(\frac{1}{n} \sum_{i=1}^n \cos(2\pi x_i)\right) + 20 + \exp(1).$$

685 As we mentioned in Section 5.1, we multiply -1 to these functions to make these problems maximiza-  
 686 tion problems. Also, even though these two functions have their optimum at  $x = 0$ , they exhibit very  
 687 different landscape as shown in Figure 7.

688 When it comes to the formal definition as RL environments, we can define these problems as follows.

- 689 • Episode Length: 1
- 690 • Observation:  $[0]$
- 691 • Action:  $n$ -dimensional vector  $\mathbf{x}$ , all of which element is in  $[-1, 1]$ .
- 692 • Reward: De Jong's Function =  $F_D(5.12\mathbf{x})$ , Ackley's Function =  $F_A(32.768\mathbf{x})$

### 693 7.5.2 Traffic Problems (Section 5.3)

694 Before introducing pace car problem, which we specifically discussed in this paper, we'd like to  
 695 briefly point out the traffic model that we used to simulate motions of individual vehicles in the traffic  
 696 environment.

697 **Traffic Model** In our traffic simulation, the state  $s \in \mathbb{R}^{2N}$  is defined as a concatenation of all  
 698 vehicle states,  $q_n \in \mathbb{R}^2$ , where  $1 \leq n \leq N$  stands for vehicle ID.  $q_n$  can be represented simply with  
 699 the vehicle's position ( $x_n$ ) and velocity ( $v_n$ ).

$$q_n = \begin{bmatrix} x_n \\ v_n \end{bmatrix}$$

$$s = [q_1^T, \dots, q_N^T]^T,$$

where  $n$  is the vehicle ID. Since this traffic state  $s$  changes over time, we represent the traffic state at timestep  $t$  as  $s_t$ . Then our differentiable traffic simulator is capable of providing the gradient of the simulation state at the next timestep  $t + 1$  with respect to the state at the current timestep  $t$ , or  $\frac{ds_{t+1}}{ds_t}$ .

This is possible because we use the Intelligent Driver Model (IDM) to model car-following behavior in our simulator Treiber et al. [2000], which is differentiable. IDM describes vehicle speed  $\dot{x}$  and acceleration  $\dot{v}$  as a function of desired velocity  $v_0$ , safe time headway in meters  $T$ , maximum acceleration  $a$ , comfortable deceleration  $b$ , the minimum distance between vehicles in meters  $\delta$ , vehicle length  $l$ , and difference in speed with the vehicle in front  $\Delta v_\alpha$  as follows:

$$\dot{x}_\alpha = \frac{dx_\alpha}{dt} = v_\alpha,$$

$$\dot{v}_\alpha = \frac{dv_\alpha}{dt} = a \left( 1 - \left( \frac{v_\alpha}{v_0} \right)^\delta - \left( \frac{s^*(v_\alpha, \Delta v_\alpha)}{s_\alpha} \right)^2 \right),$$

$$s_\alpha = x_{\alpha-1} - x_\alpha - l_{\alpha-1},$$

$$s^*(v_\alpha, \Delta v_\alpha) = s_0 + v_\alpha T + \frac{v_\alpha \Delta v_\alpha}{2\sqrt{a b}},$$

$$\Delta v_\alpha = v_\alpha - v_{\alpha-1},$$

where  $\alpha$  means the order of the vehicle in the lane. Therefore,  $(\alpha - 1)$ -th vehicle runs right in front of  $\alpha$ -th vehicle, and this relationship plays an important role in IDM. Also note that the computed acceleration term ( $\dot{v}$ ) is differentiable with respect to the other IDM variables.

IDM is just one of many car-following models used in traffic simulation literature. We choose IDM for its prevalence in previous mixed autonomy literature, however, any ODE-based car-following model will also work in our simulator as far as it is differentiable. In our simulator, automatic differentiation governs the gradient computation.

**Lane** Under IDM, lane membership is one of the most vital variables for simulation. This is because IDM builds on the leader-follower relationship. When a vehicle changes its lane with lateral movement, as shown as red arrows in Figure 9, such relationships could change, and our gradients do not convey any information about it, because lane membership is a discrete variable in nature. Note that this lateral movement also affects the longitudinal movement, rendered in green arrows in Figure 9, of a vehicle, and gradient from it is valid as far as the lane membership does not change. However, when the lane membership changes, even if our simulator gives gradient, it tells us only “partial” information in the sense that it only gives information about longitudinal behavior. Therefore, we could say that analytical gradients we get from this environment is “incomplete” (not useless at all, because it still gives us information about longitudinal movement), and thus biased.

**Pace Car Problem** In this problem, there is a single autonomous vehicle that we have to control to regulate other human driven vehicles, which follow IDM, to run at the given target speed. Since human driven vehicles control their speeds based on the relationship to their leading vehicles, our autonomous vehicle has to control itself to run in front of those human driven vehicles and adjust their speeds to the target speed. The number of lanes and human driven vehicles varies across different environments as follows.

- Single Lane: Number of lanes = 1, Number of vehicles per lane = 1
- 2 Lanes: Number of lanes = 2, Number of vehicles per lane = 2
- 4 Lanes: Number of lanes = 4, Number of vehicles per lane = 4
- 10 Lanes: Number of lanes = 10, Number of vehicles per lane = 1

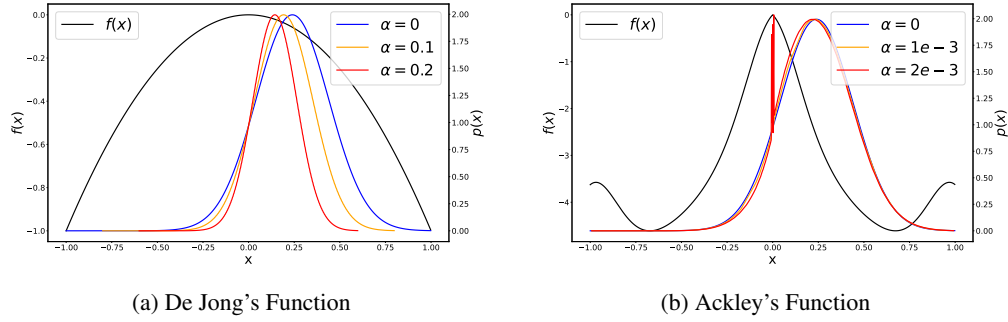


Figure 8:  $\alpha$ -policies for 1-dimensional (a) De Jong's function and (b) Ackley's function [Molga and Smutnicki, 2005].

See Figure 10 for 10-lane environment.

Then we can formally define RL environments as follows, where  $N$  is the number of human driven vehicles in total, and  $v_{tgt}$  is the target speed.

- Episode Length: 1000
- Observation:  $s \in \mathbb{R}^{2N}$
- Action:  $\mathbf{a} \in \mathbb{R}^2$  ( $\mathbf{a}_0$  = Acceleration,  $\mathbf{a}_1$  = Steering)
- Reward:  $1 - \frac{1}{N} \sum_{i=1}^N \min(\frac{|v_i - v_{tgt}|}{v_{tgt}}, 1)$

However, to facilitate learning, we additionally adopted following termination conditions to finish unpromising trials early.

- Terminate when autonomous vehicle collides with human driven vehicle
- Terminate when autonomous vehicle goes out of lane
- Terminate when autonomous vehicle runs too far away from the other vehicles
- Terminate when autonomous vehicle falls behind the other vehicles

When one of these conditions is met, it gets reward of  $-1$  and the trial terminates.

## 7.6 Renderings

### 7.6.1 $\alpha$ -policies

Here we illustrate  $\alpha$ -policies for De Jong and Ackley's function in Section 5.1. In Figure 8, the original policy is rendered in blue, and alpha policies for different  $\alpha$  values are rendered in orange and red. Those  $\alpha$ -policies are obtained from analytical gradients ( $\frac{\partial f(x)}{\partial x}$ ) of the given function  $f(x)$ , which is rendered in black. Note that singularities arise in Ackley's  $\alpha$ -policy near 0 when  $\alpha = 2 \cdot 10^{-3}$ , as it has a steep change of gradients there. This explains one of the reasons why we need to keep  $\alpha$  value sufficiently small.

### 7.6.2 Traffic Simulation

**Gradient flow** As described in Appendix 7.5.2, analytical gradients are only valid along longitudinal direction, as illustrated with green arrows in Figure 9. When a vehicle changes lane with lateral movement, as shown with red arrows, even though such change would incur different behaviors of following vehicles in multiple lanes, analytical gradients cannot convey such information. However, they still tell us information about longitudinal behavior - which is one of the motivations of our research, to use analytical gradients even in these biased environments.

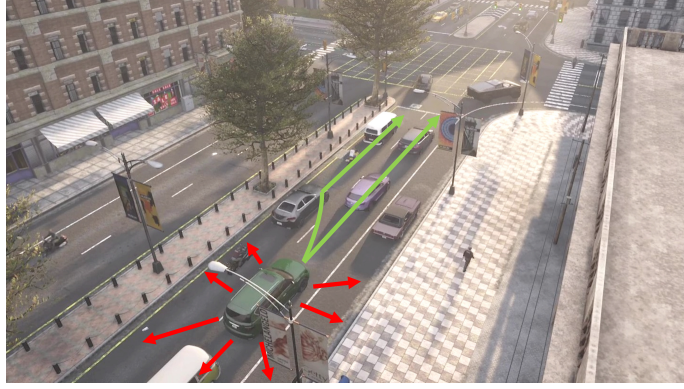


Figure 9: **Traffic Environment.** In traffic flows, only partial gradients of forward simulation are known; acceleration along a (longitudinal) traffic lane is continuous and admits gradient flow (green), while (lateral) lane-changing is a discrete event and thus prohibits gradient flow (red). Analytical gradients convey only partial information about traffic dynamics.

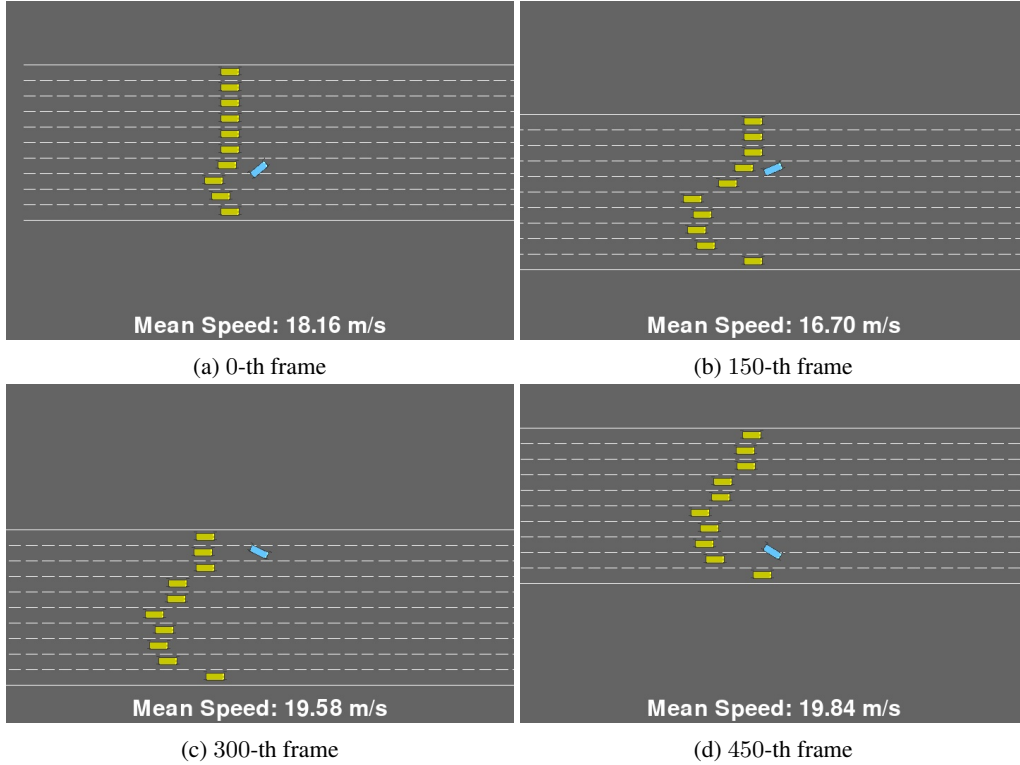


Figure 10: 10-Lane pace car environment. The autonomous vehicle (Blue) has to interfere paths of the other human driven vehicles (Yellow) to limit their speeds to  $10m/s$ . Even though it is a hard problem to achieve high score, our method achieved far better score than the baseline methods.

764 **Pace Car Problem** We additionally provide renderings of the 10-lane pace car problem defined in  
 765 Appendix 7.5.2. In Figure 10, we can observe 10 parallel lanes, 10 human driven vehicles rendered  
 766 in yellow, and the autonomous vehicle rendered in blue that we have to control. As shown there,  
 767 human driven vehicles adjust their speeds based on IDM when the autonomous vehicle blocks their  
 768 way. Therefore, the autonomous vehicle has to learn how to change their lanes, and also how to  
 769 run in appropriate speed to regulate the following vehicles. Our experimental results in Section 5.3  
 770 show that our method achieves better score than the baseline methods by adopting biased analytical  
 771 gradients in PPO.