
Relative Trajectory Balance is equivalent to Trust-PCL

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent progress in generative modeling has highlighted the importance of Rein-
2 forcement Learning (RL) for fine-tuning, with KL-regularized methods in particular
3 proving to be highly effective for both autoregressive and diffusion models. Com-
4 plementing this line of work, the Relative Trajectory Balance (RTB) objective
5 was recently introduced in the context of Generative Flow Networks (GFlowNets)
6 to serve the same role of improving fine-tuning in sequential generative models.
7 Building on prior work linking GFlowNets and maximum-entropy RL, we estab-
8 lish in this paper an equivalence between RTB and Trust-PCL, an off-policy RL
9 method with KL regularization. This equivalence situates RTB within the broader
10 theoretical landscape of KL-regularized RL, and clarifies its relationship to earlier
11 methods. Leveraging this insight, we revisit an illustrative example from the RTB
12 paper and show that KL-regularized RL methods achieve comparable performance,
13 offering an alternative perspective to what was previously reported.

14 1 Introduction

15 Generative models are playing an increasingly prominent role in artificial intelligence, ranging from
16 language and reasoning [Brown et al., 2020, Wei et al., 2022], to vision [Ho et al., 2020, Rombach
17 et al., 2022] and scientific discovery [Abramson et al., 2024]. At the heart of this rapid progress,
18 fine-tuning is central to aligning models with downstream applications while preserving diversity and
19 prior behavior. To complement supervised fine-tuning [Ouyang et al., 2022], Reinforcement Learning
20 (RL) with KL regularization has proven effective and is widely used for this task [Jaques et al., 2017,
21 Uehara et al., 2024], where reward maximization is traded off against maintaining consistency with a
22 pretrained model. In parallel, Generative Flow Networks [GFlowNets; Bengio et al., 2021, 2023]
23 have been proposed as an alternative paradigm to train generative models with techniques inspired
24 by RL. They have since gained popularity in areas such as combinatorial optimization [Zhang et al.,
25 2023], causal discovery [Deleu et al., 2022], and scientific discovery [Jain et al., 2023].

26 Extending ideas from GFlowNets, the *Relative Trajectory Balance* objective [RTB; Venkatraman
27 et al., 2024] was recently proposed as an alternative to KL-regularized RL for fine-tuning generative
28 models, framing the problem in terms of Trajectory Balance [Malkin et al., 2022]. In this paper,
29 we show that RTB is exactly equivalent to an existing RL algorithm called *Trust-PCL* [Nachum
30 et al., 2018], reinforcing the strong ties existing between GFlowNets and (entropy-regularized) RL
31 [Tiapkin et al., 2024, Mohammadpour et al., 2024, Deleu et al., 2024, Jiralerspong et al., 2025]. In
32 particular, while Venkatraman et al. [2024] suggested that KL-regularized RL methods may perform
33 poorly even on simple tasks, we argue and show on their illustrative example that these issues stem
34 from algorithmic and reward choices rather than from any fundamental limitation of RL with KL
35 regularization itself. This new perspective places RTB within the well-established framework of
36 KL-regularized RL and underscores its actual source of effectiveness.

37 2 Background

38 We study the problem of training a *sequential generative model* capable of sampling objects as
 39 a sequence of multiple steps. This covers, for example, the case of diffusion samplers [Sendera
 40 et al., 2024] that generate images \mathbf{x}_T as a sequence of denoising steps $\mathbf{x}_1 \rightarrow \mathbf{x}_2 \rightarrow \dots \rightarrow \mathbf{x}_T$,
 41 starting from some Gaussian noise $\mathbf{x}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and sequentially sampling from a denoising model
 42 $\mathbf{x}_{t+1} \sim P_\phi(\mathbf{x}_{t+1} | \mathbf{x}_t)$. By following this iterative process, we eventually obtain samples from the
 43 *marginal distribution* defined by

$$P_\phi^\top(\mathbf{x}_T) = \int P_\phi(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T) d\mathbf{x}_{1:T-1} = \int \left[P(\mathbf{x}_1) \prod_{t=1}^{T-1} P_\phi(\mathbf{x}_{t+1} | \mathbf{x}_t) \right] d\mathbf{x}_{1:T-1}. \quad (1)$$

44 We assume that we have access to a base sequential generative model, with transition probabilities
 45 $\pi_{\text{prior}}(\mathbf{x}_{t+1} | \mathbf{x}_t)$, such as a model pre-trained with a large dataset of observations, and that we are
 46 also given a fixed *energy function* $\mathcal{E}(\mathbf{x}_T)$ that we can query for any object generated by the model.
 47 Our objective in this paper is to learn a sequential generative model $P_\phi(\mathbf{x}_{t+1} | \mathbf{x}_t)$ such that its
 48 marginal distribution is the *tilted distribution*, modulated by $\mathcal{E}(\mathbf{x}_T)$:

$$P_\phi^\top(\mathbf{x}_T) \propto \pi_{\text{prior}}^\top(\mathbf{x}_T) \exp(-\mathcal{E}(\mathbf{x}_T)/\alpha), \quad (2)$$

49 for some temperature parameter $\alpha > 0$. The role of the energy function is to steer generation towards
 50 samples with desired properties (e.g., generating images based on a description).

51 2.1 Entropy-regularized Reinforcement Learning

52 One way to approach this problem is to view it from the perspective of Reinforcement Learning
 53 [Korbak et al., 2022]. We consider a finite-horizon Markov Decision Process (MDP) $\mathcal{M} = (\bar{\mathcal{S}}, \mathcal{A}, r)$,
 54 where $\bar{\mathcal{S}} = \mathcal{S} \cup \{\perp\}$ is a state space augmented by a special state $\perp \notin \mathcal{S}$ indicating the end of a
 55 trajectory, \mathcal{A} is the action space, and r is a reward function that we will detail below. We identify a
 56 state $s_0 \in \mathcal{S}$ called the *initial state* from which all trajectories start; these trajectories are guaranteed
 57 to end in \perp since the MDP is finite-horizon. To give a concrete example, we may see the multiple
 58 steps of denoising in a diffusion sampler as an MDP with a specific structure [Fan et al., 2023, Black
 59 et al., 2024]: its states are of the form $s_t = (\mathbf{x}_t, t)$, transitions (actions) correspond to applying one
 60 step of denoising $(\mathbf{x}_t, t) \rightarrow (\mathbf{x}_{t+1}, t+1)$, and trajectories must terminate when a state of the form
 61 (\mathbf{x}_T, T) is reached after T steps (termination being indicated by the transition $(\mathbf{x}_T, T) \rightarrow \perp$).

62 Following Deleu et al. [2024], the reward function $r(s, s')$ obtained when transitioning from $s \rightarrow s'$
 63 is defined such that the sum of rewards along a trajectory only depends on the energy of the final state
 64 it reaches right before terminating; in other words, for a trajectory $\tau = (s_0, s_1, \dots, s_T, \perp)$, we have

$$\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T), \quad (3)$$

65 where we will use the convention $s_{T+1} = \perp$ throughout this paper. This includes in particular the
 66 case where the reward is only obtained at the end of the trajectory (i.e., $r(s_T, \perp) = -\mathcal{E}(s_T)$, and
 67 zero everywhere else, also called *outcome-based reward*; Uesato et al., 2022).

68 Contrary to standard RL, where the objective is to find a policy $\pi(s_{t+1} | s_t)$ that maximizes the
 69 expected sum of rewards, *KL-regularized RL* includes an additional KL regularization term between
 70 the current policy and an anchor policy π_{prior} we don't want to deviate too much from:

$$\pi_{\text{RelEnt}}^* = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^T r(s_t, s_{t+1}) - \alpha \text{KL}(\pi(\cdot | s_t) \| \pi_{\text{prior}}(\cdot | s_t)) \right]. \quad (4)$$

71 It can be shown that with the particular choice of reward function made in (3), the marginal distribution
 72 associated with this optimal policy is the tilted distribution [Nachum et al., 2018, Korbak et al., 2022]

$$\pi_{\text{RelEnt}}^{\star\top}(s_T) \propto \pi_{\text{prior}}^\top(s_T) \exp(-\mathcal{E}(s_T)/\alpha), \quad (5)$$

73 where the temperature parameter α naturally emerges from the regularization constant in (4); a proof
 74 of this is given in Appendix A for completeness. This constitutes the foundation of the recent line
 75 of work using RL for fine-tuning language models [Jaques et al., 2017, Ouyang et al., 2022] and
 76 diffusion models [Fan et al., 2023, Black et al., 2024, Uehara et al., 2024].

77 **2.2 Relative trajectory balance**

78 Taking another perspective, this time rooted in the literature on GFlowNets [Bengio et al., 2023],
 79 Venkatraman et al. [2024] introduced an objective called the *Relative Trajectory Balance* (RTB) loss
 80 to sample from the tilted distribution (2). Given a transition probability $P_\phi(s_{t+1} | s_t)$ and a scalar
 81 $Z_\psi > 0$, the RTB loss is a non-linear least-square objective $\mathcal{L}_{\text{RTB}}(\phi, \psi) = \frac{1}{2} \mathbb{E}_{\pi_b} [\Delta_{\text{RTB}}^2(\tau; \phi, \psi)]$,
 82 where π_b is an arbitrary distribution over trajectories (this is an *off-policy* objective), and the residual
 83 is defined as

$$\Delta_{\text{RTB}}(\tau; \phi, \psi) = \log \frac{\prod_{t=0}^T \pi_{\text{prior}}(s_{t+1} | s_t)}{Z_\psi \prod_{t=0}^T P_\phi(s_{t+1} | s_t)} - \frac{\mathcal{E}(s_T)}{\alpha}. \quad (6)$$

84 Venkatraman et al. [2024] showed that if this RTB loss is minimized, then the marginal distribution
 85 associated with the optimal $P_\phi(s_{t+1} | s_t)$ satisfies (2), and the optimal Z_ψ is the normalization
 86 constant of this tilted distribution.

87 **3 Equivalence between RTB and Trust-PCL**

88 In line with the growing body of work recently connecting GFlowNets with Maximum Entropy
 89 Reinforcement Learning (MaxEnt-RL) [Tiapkin et al., 2024, Mohammadpour et al., 2024, Deleu
 90 et al., 2024], where all the losses introduced in the GFlowNet literature have been shown to have
 91 a counterpart in MaxEnt-RL, one may wonder whether the RTB loss can also be viewed from the
 92 point of view of RL. In fact, we saw that solving (4) already provides a way to sample from the tilted
 93 distribution (5), and there exists an off-policy algorithm called *Trust-PCL* [Nachum et al., 2018] that
 94 does exactly this. For a policy $\pi_\phi(s_{t+1} | s_t)$ and a soft state-value function $V_{\text{soft}}^\psi(s)$, Trust-PCL is
 95 also defined as a non-linear least-square objective $\mathcal{L}_{\text{T-PCL}}(\phi, \psi) = \frac{1}{2} \mathbb{E}_{\pi_b} [\Delta_{\text{T-PCL}}^2(\tau; \phi, \psi)]$, where

$$\Delta_{\text{T-PCL}}(\tau; \phi, \psi) = -V_{\text{soft}}^\psi(s_0) + \sum_{t=0}^T r(s_t, s_{t+1}) + \alpha \log \frac{\pi_{\text{prior}}(s_{t+1} | s_t)}{\pi_\phi(s_{t+1} | s_t)}. \quad (7)$$

96 This form suggests that Relative Trajectory Balance is exactly equivalent to Trust-PCL, in the sense
 97 that both losses are equal up to a constant factor that only depends on the temperature α .

98 **Proposition 1** (Equivalence RTB – Trust-PCL). *The Relative Trajectory Balance loss [RTB; Venka-*
 99 *traman et al., 2024] defined with the residual in (6) is proportional to the Trust-PCL objective*
 100 *[Nachum et al., 2018] defined with the residual in (7) on the MDP in Section 2.1 (in particular, whose*
 101 *reward function satisfies (3)): $\mathcal{L}_{\text{T-PCL}}(\phi, \psi) = \alpha^2 \mathcal{L}_{\text{RTB}}(\phi, \psi)$, with*

$$V_{\text{soft}}^\psi(s_0) = \alpha \log Z_\psi \qquad \pi_\phi(s' | s) = P_\phi(s' | s) \quad (8)$$

102 This result differs from similar ones relating GFlowNet objectives with MaxEnt-RL [Tiapkin et al.,
 103 2024, Deleu et al., 2024] in that unlike those existing connections, the equivalence in Proposition 1
 104 does not require any correction of the reward function (with a backward transition probability). The
 105 proof of this proposition is immediate with the correspondence in (8).

106 **3.1 Reinterpreting the empirical success of RTB**

107 Validating Proposition 1 empirically by comparing the performance of Trust-PCL against RTB would
 108 be of limited interest, since they both optimize the exact same loss (up to a constant factor). Therefore,
 109 we instead choose to revisit the results presented in the RTB paper in light of this equivalence.
 110 While Venkatraman et al. [2024] argued that “*RL methods with KL regularization yield inaccurate*
 111 *inference*”, our proposition suggests that this is not a fundamental issue of KL-regularized RL as a
 112 whole, and the success of RTB should be attributed to the choice of a superior RL algorithm, namely
 113 Trust-PCL.

114 To illustrate the apparent failure of RL with KL regularization for sequential generative modeling,
 115 Venkatraman et al. [2024] considered a 2D generation task where the prior is a uniform mixture of 25
 116 Gaussians, and the target marginal distribution is a weighted mixture (Figures 1a & 1b). While RTB
 117 is capable of perfectly recovering the target distribution (Figure 1c), they observed that RL with KL
 118 regularization seems to not be capturing all the modes of the target distribution (Figure 1d).

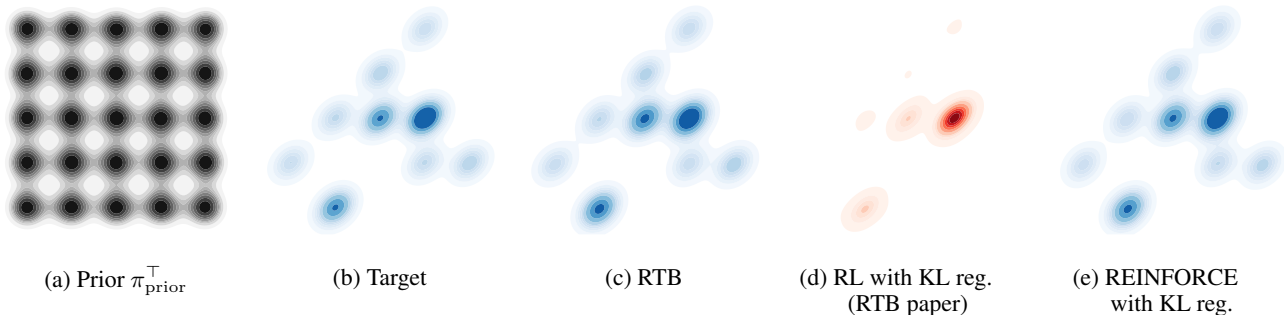


Figure 1: Comparison of RTB and RL with KL regularization on the illustrative example of [Venkatraman et al., 2024]. (a) Prior distribution π_{prior}^T . (b) Target tilted distribution $\propto \pi_{\text{prior}}^T(x) \exp(-\mathcal{E}(x))$. (c) Model trained with RTB. (d) Model trained with REINFORCE with KL regularization, as reported by Venkatraman et al. [2024]. (e) Corrected model trained with off-policy REINFORCE with KL regularization.

119 At first glance, this could be explained by the difference in performance between their RL algorithm
 120 of choice REINFORCE [Williams, 1992] and Trust-PCL/RTB. However, we show in Appendix B
 121 that this happens to be caused by their reward function not satisfying (3). This, in turns, changes the
 122 target distribution (5), which does not match the one in Figure 1b anymore. We show in Figure 1e
 123 that if we apply the same method REINFORCE on the MDP of Section 2.1 with minimal changes
 124 to ensure a fair comparison (e.g., using off-policy data with self-normalized importance sampling
 125 correction [Precup et al., 2000]; see Algorithm 2 for details), then the expected target distribution can
 126 also be recovered as accurately as with RTB.

127 Going beyond this illustrative example, Venkatraman et al. [2024] also compared RTB *on-policy*
 128 (i.e., using $\pi_b \equiv P_\phi$ in Section 2.2) to existing methods based on RL for conditional image generation
 129 [Black et al., 2024, Fan et al., 2023]. Correcting for the discrepancy in the reward function mentioned
 130 above, we show in Appendix C that in this setting, RTB (on-policy) is equivalent to REINFORCE
 131 with KL regularization, up to a different control variate. This mirrors existing results in the GFlowNet
 132 literature [Malkin et al., 2023].

133 4 Conclusion

134 Reinforcement Learning with KL regularization has proven to be a versatile tool for fine-tuning
 135 sequential generative models. Its flexibility is such that even seemingly new methods for addressing
 136 the same problem can often be framed within this paradigm. In this work, we investigated the
 137 recently proposed Relative Trajectory Balance objective [Venkatraman et al., 2024] through the lens
 138 of entropy-regularized Reinforcement Learning. We demonstrated that it is exactly equivalent to the
 139 Trust-PCL algorithm [Nachum et al., 2018], indicating that RTB effectively rephrases an existing
 140 KL-regularized RL approach within the GFlowNet framework. This is further evidence of the strong
 141 theoretical connections between these fields.

142 Leveraging this theoretical correspondence, we revisited the empirical claims by Venkatraman et al.
 143 [2024] regarding the shortcomings of KL-regularized RL. By carefully examining the illustrative
 144 experiment from the RTB paper, we found that simple methods such as REINFORCE with KL
 145 regularization can match RTB in modeling the target tilted distribution on that task. The reported
 146 failures stemmed from algorithmic and reward design choices rather than any fundamental limitation
 147 of Reinforcement Learning. However, on more challenging tasks, advanced methods like RTB/Trust-
 148 PCL are likely to offer clearer benefits, as suggested by recent work incorporating techniques
 149 such as update clipping [Fan et al., 2023, Black et al., 2024, Shao et al., 2024] (inspired by PPO
 150 [Schulman et al., 2017]). These observations emphasize the need for fair comparisons and well-
 151 posed problem formulations, while offering promising directions for capitalizing on the connections
 152 between GFlowNets and RL to explore new algorithms for generative model fine-tuning.

153 **References**

- 154 J. Abramson, J. Adler, J. Dunger, R. Evans, T. Green, A. Pritzel, O. Ronneberger, L. Willmore,
155 A. J. Ballard, J. Bambrick, et al. Accurate structure prediction of biomolecular interactions with
156 AlphaFold 3. *Nature*, 2024.
- 157 E. Bengio, M. Jain, M. Korablyov, D. Precup, and Y. Bengio. Flow Network based Generative Models
158 for Non-Iterative Diverse Candidate Generation. *Advances in Neural Information Processing*
159 *Systems (NeurIPS)*, 2021.
- 160 Y. Bengio, S. Lahlou, T. Deleu, E. J. Hu, M. Tiwari, and E. Bengio. GFlowNet Foundations. *Journal*
161 *of Machine Learning Research (JMLR)*, 2023.
- 162 K. Black, M. Janner, Y. Du, I. Kostrikov, and S. Levine. Training Diffusion Models with Reinforce-
163 ment Learning. *International Conference on Learning Representations (ICLR)*, 2024.
- 164 T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,
165 G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in Neural Information*
166 *Processing Systems (NeurIPS)*, 2020.
- 167 T. Deleu, A. Góis, C. Emezue, M. Rankawat, S. Lacoste-Julien, S. Bauer, and Y. Bengio. Bayesian
168 Structure Learning with Generative Flow Networks. *Conference on Uncertainty in Artificial*
169 *Intelligence (UAI)*, 2022.
- 170 T. Deleu, P. Nouri, N. Malkin, D. Precup, and Y. Bengio. Discrete Probabilistic Inference as Control
171 in Multi-path Environments. *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2024.
- 172 Y. Fan, O. Watkins, Y. Du, H. Liu, M. Ryu, C. Boutilier, P. Abbeel, M. Ghavamzadeh, K. Lee,
173 and K. Lee. DPOK: Reinforcement Learning for Fine-tuning Text-to-Image Diffusion Models.
174 *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- 175 J. Ho, A. Jain, and P. Abbeel. Denoising Diffusion Probabilistic Models. *Advances in Neural*
176 *Information Processing Systems (NeurIPS)*, 2020.
- 177 M. Jain, T. Deleu, J. Hartford, C.-H. Liu, A. Hernandez-Garcia, and Y. Bengio. GFlowNets for
178 AI-Driven Scientific Discovery. *Digital Discovery*, 2023.
- 179 N. Jaques, S. Gu, D. Bahdanau, J. M. Hernández-Lobato, R. E. Turner, and D. Eck. Sequence
180 Tutor: Conservative Fine-Tuning of Sequence Generation Models with KL-control. *International*
181 *Conference on Machine Learning (ICML)*, 2017.
- 182 M. Jiralerspong, E. Derman, D. Vucetic, N. Malkin, B. Sun, T. Zhang, P.-L. Bacon, and G. Gidel.
183 Robust Reinforcement Learning for Discrete Compositional Generation via General Soft Operators.
184 *arXiv Preprint*, 2025.
- 185 T. Korbak, E. Perez, and C. L. Buckley. RL with KL penalties is better viewed as Bayesian inference.
186 *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2022.
- 187 N. Malkin, M. Jain, E. Bengio, C. Sun, and Y. Bengio. Trajectory Balance: Improved Credit
188 Assignment in GFlowNets. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 189 N. Malkin, S. Lahlou, T. Deleu, X. Ji, E. Hu, K. Everett, D. Zhang, and Y. Bengio. GFlowNets and
190 variational inference. *International Conference on Learning Representations (ICLR)*, 2023.
- 191 S. Mohammadpour, E. Bengio, E. Frejinger, and P.-L. Bacon. Maximum entropy GFlowNets with
192 soft Q-learning. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2024.
- 193 O. Nachum, M. Norouzi, K. Xu, and D. Schuurmans. Trust-PCL: An Off-Policy Trust Region Method
194 for Continuous Control. *International Conference on Learning Representations (ICLR)*, 2018.
- 195 L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama,
196 A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in*
197 *Neural Information Processing Systems (NeurIPS)*, 2022.

- 198 D. Precup, R. S. Sutton, and S. Singh. Eligibility traces for off-policy policy evaluation. *International*
199 *Conference on Machine Learning (ICML)*, 2000.
- 200 R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-Resolution Image Synthesis
201 with Latent Diffusion Models. *IEEE Conference on Computer Vision and Pattern Recognition*
202 *(CVPR)*, 2022.
- 203 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal Policy Optimization
204 Algorithms. *arXiv Preprint*, 2017.
- 205 M. Sendera, M. Kim, S. Mittal, P. Lemos, L. Scimeca, J. Rector-Brooks, A. Adam, Y. Bengio, and
206 N. Malkin. On diffusion models for amortized inference: Benchmarking and improving stochastic
207 control and sampling. *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- 208 Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu, et al.
209 DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *arXiv*
210 *Preprint*, 2024.
- 211 D. Tiapkin, N. Morozov, A. Naumov, and D. Vetrov. Generative Flow Networks as Entropy-
212 Regularized RL. *International Conference on Artificial Intelligence and Statistics (AISTATS)*,
213 2024.
- 214 M. Uehara, Y. Zhao, T. Biancalani, and S. Levine. Understanding reinforcement learning-based
215 fine-tuning of diffusion models: A tutorial and review. *arXiv Preprint*, 2024.
- 216 J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins.
217 Solving Math Word Problems with Process-based and Outcome-based Feedback. *arXiv Preprint*,
218 2022.
- 219 S. Venkatraman, M. Jain, L. Scimeca, M. Kim, M. Sendera, M. Hasan, L. Rowe, S. Mittal, P. Lemos,
220 E. Bengio, A. Adam, J. Rector-Brooks, Y. Bengio, G. Berseth, and N. Malkin. Amortizing
221 intractable inference in diffusion models for vision, language, and control. *Advances in Neural*
222 *Information Processing Systems (NeurIPS)*, 2024.
- 223 J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-
224 Thought Prompting Elicits Reasoning in Large Language Models. *Advances in Neural Information*
225 *Processing Systems (NeurIPS)*, 2022.
- 226 R. J. Williams. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement
227 Learning. *Machine learning*, 1992.
- 228 D. Zhang, H. Dai, N. Malkin, A. Courville, Y. Bengio, and L. Pan. Let the Flows Tell: Solving
229 Graph Combinatorial Optimization Problems with GFlowNets. *Advances in Neural Information*
230 *Processing Systems (NeurIPS)*, 2023.

231 Appendix

232 A Marginal distribution of the optimal KL-regularized policy

233 In this section, we prove that the marginal distribution associated with the optimal policy maximizing
 234 (4) is the tilted distribution (5) for completeness, although this result can be found in various contexts
 235 in the literature [Nachum et al., 2018, Korbak et al., 2022]. First, recall that the KL-regularized RL
 236 objective can be written as

$$\pi_{\text{RelEnt}}^* = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^T r(s_t, s_{t+1}) - \alpha \text{KL}(\pi(\cdot | s_t) \| \pi_{\text{prior}}(\cdot | s_t)) \right]. \quad (9)$$

237 It can be shown that the optimal policy π_{RelEnt}^* can be written in terms of a soft state-value function
 238 $V_{\text{soft}}^*(s)$ and a soft state-action value function $Q_{\text{soft}}^*(s, s')$ as

$$\pi_{\text{RelEnt}}^*(s' | s) = \pi_{\text{prior}}(s' | s) \exp \left(\frac{1}{\alpha} (Q_{\text{soft}}^*(s, s') - V_{\text{soft}}^*(s)) \right), \quad (10)$$

239 where the soft value functions, adapted to our setting where the MDP \mathcal{M} is finite horizon and
 240 deterministic, satisfy

$$Q_{\text{soft}}^*(s, s') = r(s, s') + V_{\text{soft}}^*(s') \quad (11)$$

$$V_{\text{soft}}^*(s) = \alpha \log \sum_{s' \in \text{Ch}(s)} \pi_{\text{prior}}(s' | s) \exp \left(\frac{1}{\alpha} Q_{\text{soft}}^*(s, s') \right), \quad (12)$$

241 where $\text{Ch}(s)$ are the children of the state s in the MDP (e.g., all the states of the form $(x_{t+1}, t+1)$
 242 for a state $s_t = (x_t, t)$ in the example MDP given in Section 2.1 [Fan et al., 2023, Black et al., 2024]).
 243 We defined the marginal distribution (1) in the main text for the particular example of diffusion
 244 samplers. However, this naturally generalizes in the context of KL-regularized RL, by viewing it as
 245 a marginal over trajectories in the MDP that terminate in s_T (i.e., the trajectory terminates with a
 246 transition $s_T \rightarrow \perp$):

$$\pi^\top(s_T) = \int \pi(\tau | s_0) \mathbf{1}(s_T \rightarrow \perp \in \tau) d\tau \quad (13)$$

247 For a trajectory $\tau = (s_0, s_1, \dots, s_T, \perp)$, we have

$$\pi_{\text{RelEnt}}^*(\tau | s_0) = \prod_{t=0}^T \pi_{\text{RelEnt}}^*(s_{t+1} | s_t) \quad (14)$$

$$= \left[\prod_{t=0}^T \pi_{\text{prior}}(s_{t+1} | s_t) \right] \exp \left(\frac{1}{\alpha} \sum_{t=0}^T Q_{\text{soft}}^*(s_t, s_{t+1}) - V_{\text{soft}}^*(s_t) \right) \quad (15)$$

$$= \pi_{\text{prior}}(\tau | s_0) \exp \left(\frac{1}{\alpha} \sum_{t=0}^T r(s_t, s_{t+1}) + V_{\text{soft}}^*(s_{t+1}) - V_{\text{soft}}^*(s_t) \right) \quad (16)$$

$$= \pi_{\text{prior}}(\tau | s_0) \exp \left(\frac{1}{\alpha} (-\mathcal{E}(s_T) + V_{\text{soft}}^*(\perp) - V_{\text{soft}}^*(s_0)) \right) \quad (17)$$

$$\propto \pi_{\text{prior}}(\tau | s_0) \exp(-\mathcal{E}(s_T)/\alpha), \quad (18)$$

248 where we used the definition of π_{RelEnt}^* in (15), the definition of $Q_{\text{soft}}^*(s_t, s_{t+1})$ in (16), the property
 249 that the sum of reward functions only depends on the energy at the end of the trajectory (3) and a
 250 telescoping sum in (17), and the fact that $V_{\text{soft}}^*(s_0)$ is a constant independent of τ (and $V_{\text{soft}}^*(\perp) = 0$)
 251 in (18). Plugging this in the definition of the marginal distribution (13), we get

$$\pi_{\text{RelEnt}}^{\star\top}(s_T) = \int \pi_{\text{RelEnt}}^*(\tau | s_0) \mathbf{1}(s_T \rightarrow \perp \in \tau) d\tau \quad (19)$$

$$\propto \exp(-\mathcal{E}(s_T)/\alpha) \int \pi_{\text{prior}}(\tau | s_0) \mathbf{1}(s_T \rightarrow \perp \in \tau) d\tau \quad (20)$$

$$\propto \pi_{\text{prior}}^\top(s_T) \exp(-\mathcal{E}(s_T)/\alpha). \quad (21)$$

Algorithm 1 RL with KL reg. [Venkatraman et al., 2024]

- 1: Sample a batch of trajectories $\{\tau_1, \dots, \tau_N\}$ using π_ϕ
- 2: Compute the cumulative reward for τ_n :

$$r_n \leftarrow \exp(-\mathcal{E}(s_T^{(n)}))$$

- 3: Compute the advantage: $\bar{r}_n = r_n - \frac{1}{N} \sum_{n=1}^N r_n$
- 5: Compute the REINFORCE loss with KL reg.

$$\mathcal{L}(\phi) = \frac{1}{N} \sum_{n=1}^N -\text{sg}(\bar{r}_n) \log \pi_\phi(\tau_n) + \frac{\lambda}{2} \left(\log \frac{\pi_\phi(\tau_n)}{\pi_{\text{prior}}(\tau_n)} \right)^2$$

Algorithm 2 Off-policy REINFORCE with KL reg.

- 1: Sample a batch of trajectories $\{\tau_1, \dots, \tau_N\}$ using π_b
- 2: Compute the cumulative reward for τ_n :

$$r_n \leftarrow -\mathcal{E}(s_T^{(n)}) - \alpha \log \frac{\pi_\phi(\tau_n)}{\pi_{\text{prior}}(\tau_n)}$$

- 3: Compute the advantage: $\bar{r}_n = r_n - \frac{1}{N} \sum_{n=1}^N r_n$
- 4: Get the SNIS weights: $w_n \propto (\pi_\phi(\tau_n)/\pi_b(\tau_n))^{1/T}$
- 5: Compute the REINFORCE loss

$$\mathcal{L}(\phi) = -\frac{1}{N} \sum_{n=1}^N \text{sg}(w_n \bar{r}_n) \log \pi_\phi(\tau_n)$$

252 B Off-policy REINFORCE with KL regularization

253 We saw in Section 2.1 that the marginal distribution associated with the optimal policy π_{RelEnt}^*
 254 matches the tilted distribution (5) in the case where the reward function of the MDP satisfies

$$\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T). \quad (22)$$

255 In their comparison with KL-regularized RL, Venkatraman et al. [2024] instead used a sparse reward
 256 function that equals $\tilde{r}(s_T, \perp) = \exp(-\mathcal{E}(s_T))$ at the terminating transition, and zero everywhere
 257 else, meaning that their optimum would be

$$\pi^{*\top}(s_T) \propto \pi_{\text{prior}}^\top(s_T) \exp(\exp(-\mathcal{E}(s_T))/\alpha), \quad (23)$$

258 where we emphasize the extra inner “exp”. This error likely came from the naming conflict between
 259 the “reward function” in RL and the “reward” in the GFlowNet literature, where the latter should
 260 indeed equal the exponential of the corresponding reward in RL [Deleu et al., 2024]. This can be
 261 seen in the residual for RTB (6), that can be re-written as

$$\Delta_{\text{RTB}}(\tau; \phi, \psi) = \log \frac{\overbrace{\exp(-\mathcal{E}(s_T)/\alpha)}^{\text{GFlowNet reward } "R(s_T)"}}{\prod_{t=0}^T \pi_{\text{prior}}(s_{t+1} | s_t)} \cdot \frac{1}{Z_\psi \prod_{t=0}^T P_\phi(s_{t+1} | s_t)}. \quad (24)$$

262 Besides this mismatch in reward functions, we made some minor updates to the REINFORCE
 263 [Williams, 1992] algorithm from Algorithm 1 considered in the RTB paper to Algorithm 2 (“sg” is
 264 the stop-gradient operator), to ensure a fair comparison with RTB:

- 265 • The trajectories are collected using the same exploratory behavior policy π_b as RTB, to address
 266 the well known limitations of on-policy methods in terms of exploration (this was also noted by
 267 Venkatraman et al. [2024]). To that end, we use importance sampling to address this shift in
 268 distributions (between π_b and the policy π_ϕ being learned), and more particularly self-normalized
 269 importance sampling [SNIS; Precup et al., 2000] to control the variance of the gradient estimate.
- 270 • The log-ratio $\log \pi_\phi/\pi_{\text{prior}}$ is added to the reward itself, as opposed to being a squared reg-
 271 ularization, so that it can be incorporated into the average control variate. Note that while
 272 Venkatraman et al. [2024] treated the regularization constant λ as a separate hyperparameter they
 273 had to tune carefully (with $\lambda = \alpha = 1$ being the best value they found), we saw in Section 2.1
 274 that this exactly corresponds to the temperature parameter α of the tilted distribution.

275 Going beyond the illustrative example of Figure 1, we note that Venkatraman et al. [2024] only
 276 considered baselines based on vanilla REINFORCE in their comparisons with KL-regularized RL
 277 methods, and they did not use the advanced techniques such as update clipping that were used in
 278 practice by DDPO [Black et al., 2024] & DPOK [Fan et al., 2023].

279 **C On-policy RTB is equivalent to REINFORCE with KL regularization**

280 In this section, we will show that the gradient of on-policy RTB is equal to the gradient of REIN-
 281 FORCE with KL regularization, up to a constant α . First, recall that the gradient of the REINFORCE
 282 loss with KL regularization can be written as

$$\nabla_{\phi} \mathcal{L}_{\text{RL}}(\phi) = -\mathbb{E}_{\pi_{\phi}} \left[\left(-\mathcal{E}(s_T) - \alpha \log \frac{\pi_{\phi}(\tau)}{\pi_{\text{prior}}(\tau)} - b \right) \nabla_{\phi} \log \pi_{\phi}(\tau) \right], \quad (25)$$

283 where b is a control variate (baseline), which is often estimated using an average over the batch of
 284 trajectories (see [Algorithm 2](#)). Similarly, recall that the (on-policy, *i.e.*, where $\pi_b \equiv \pi_{\phi}$) RTB loss in
 285 (6) can be written as

$$\mathcal{L}_{\text{RTB}}(\phi, \psi) = \frac{1}{2} \mathbb{E}_{\pi_{\phi}} \left[\left(\log \frac{\pi_{\text{prior}}(\tau)}{Z_{\psi} \pi_{\phi}(\tau)} - \frac{\mathcal{E}(s_T)}{\alpha} \right)^2 \right] \quad (26)$$

286 where we use the notation π_{ϕ} instead of P_{ϕ} to match the policy in KL-regularized RL. Taking the
 287 gradient of \mathcal{L}_{RTB} wrt. ϕ , and ignoring differentiation through the policy over which we take the
 288 expectation (which is standard in the GFlowNet literature [[Malkin et al., 2023](#)]), we get

$$\nabla_{\phi} \mathcal{L}_{\text{RTB}}(\phi, \psi) = \frac{1}{2} \mathbb{E}_{\pi_{\phi}} [\nabla_{\phi} \Delta_{\text{RTB}}^2(\tau; \phi, \psi)] \quad (27)$$

$$= \mathbb{E}_{\pi_{\phi}} \left[\left(\frac{\mathcal{E}(s_T)}{\alpha} + \log \frac{\pi_{\phi}(\tau)}{\pi_{\text{prior}}(\tau)} + \log Z_{\psi} \right) \nabla_{\phi} \log \pi_{\phi}(\tau) \right] \quad (28)$$

$$= \frac{1}{\alpha} \nabla_{\phi} \mathcal{L}_{\text{RL}}(\phi), \quad (29)$$

289 If the baseline $b = \alpha \log Z_{\psi}$. Therefore both methods are equivalent, only differing their control
 290 variate: REINFORCE with KL regularization often use a local baseline estimated using the current
 291 batch of trajectories, whereas Z_{ψ} in on-policy RTB acts as a *global* baseline. This result is similar to
 292 the one derived by [Malkin et al. \[2023\]](#).