

GFlowNet Training by Policy Gradients

Anonymous Authors¹

Abstract

Generative Flow Networks (GFlowNets) have been shown effective to generate combinatorial objects with desired properties. We here propose a new GFlowNet training framework, with policy-dependent rewards, that bridges keeping flow balance of GFlowNets to optimizing the expected accumulated reward in traditional Reinforcement-Learning (RL). This enables the derivation of new policy-based GFlowNet training methods, in contrast to existing ones resembling value-based RL. It is known that the design of backward policies in GFlowNet training affects efficiency. We further develop a coupled training strategy that jointly solves GFlowNet forward policy training and backward policy design. Performance analysis is provided with a theoretical guarantee of our policy-based GFlowNet training. Experiments on both simulated and real-world datasets verify that our policy-based strategies provide advanced RL perspectives for robust gradient estimation to improve GFlowNet performance.

1. Introduction

Generative Flow Networks (GFlowNets) are a family of generative models on the space of combinatorial objects \mathcal{X} , e.g. graphs composed by organizing nodes and edges in a particular manner, or strings composed of characters in a particular ordering. GFlowNets aim to solve a challenging task, sampling $x \in \mathcal{X}$ with a probability proportional to some non-negative reward function $R(x)$ that defines an unnormalized distribution, where $|\mathcal{X}|$ can be enormous and the distribution modes are highly isolated by its combinatorial nature. GFlowNets (Bengio et al., 2021; 2023) decompose the process of generating or sampling $x \in \mathcal{X}$ by generating incremental trajectories that start from a null state, pass through intermediate states, and end at x as the desired ter-

minating state. These trajectory instances are interpreted as the paths along a Directed Acyclic Graph (DAG). Probability measures of trajectories are viewed as the amount of ‘water’ flows along the DAG, with $R(x)$ being the total flow of trajectories that end at x , so that following the forward generating policy defined by the measure, sampled trajectories will end at x with the probability proportional to $R(x)$.

GFlowNets bear a similar form of reinforcement learning (RL) in that they both operate over Markovian Decision Processes (MDP) with a reward function $R(x)$, where nodes, edges, and node transition distributions defined by Markovian flows are considered as states, actions, and stochastic policies in MDPs. They, however, differ in the following aspects: the goal of RL problems is to learn optimal policies that maximize the expected cumulative trajectory reward by R . For **value-based** RL methods, the key to achieve this is by reducing the Temporal Difference (TD) error of Bellman equations for the estimated state value function V and state-action value function Q (Sutton & Barto, 2018; Mnih et al., 2013). GFlowNets amortize the sampling problem into finding some Markovian flow that assigns the proper probability flow to edges (actions) so that the total flow of trajectories ending at x is $R(x)$. When studying these in the lens of RL, the existing GFlowNet training strategies are also value-based in that they achieve the goal by keeping the balance flow equation over states of the DAG, whose difference can be measured in trajectory-wise and edge-wise ways (Bengio et al., 2021; 2023; Malkin et al., 2022a; Madan et al., 2023).

Due to the similarity of GFlowNet training and RL, investigating the relationships between them can not only deepen understanding of GFlowNets but also help derive better training methods from RL. In this work, we propose policy-dependent rewards for GFlowNet training. This bridges GFlowNets to RL in that keeping the flow balance over DAGs can be reformulated as optimizing the expected accumulated rewards in RL problems. We then derive **policy-based** training strategies for GFlowNets, which optimize the accumulated reward by its gradients with respect to (w.r.t.) the forward policy directly (Sutton et al., 1999; Sutton & Barto, 2018).

In terms of RL, we acknowledge that the existing GFlowNet training methods can be considered value-based and have

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

the advantage of allowing off-policy training over policy-based methods (Malkin et al., 2022b). Value-based methods, however, face the difficulty in designing a powerful sampler that can balance the exploration and exploitation trade-off, especially when the combinatorial space is enormous with well-isolated modes. Besides, employing typical annealing or random-mixing solutions may lead to the learned policy trapped in local optima. Finally, designing strategies for powerful samplers vary according to the structure and setting of modeling environments. Policy-based methods, especially the on-policy ones, transform the design of a powerful sampler into robust estimation of policy gradients, which can be achieved by variance reduction techniques (Schulman et al., 2016) and improvement of gradient descent directions, such as natural policy gradients (Kakade, 2001) and mirror policy descent (Zhan et al., 2023). Conservation policy updates such as Trust-Region Policy Optimization (TRPO) (Schulman et al., 2015; Achiam et al., 2017) and its first-order approximation, Proximal Policy Optimization (PPO), have also been developed, for example as the backbone of ChatGPT (Ouyang et al., 2022). Moreover, policy-based methods can be made off-policy, for example, by importance sampling (Degris et al., 2012). Our work provides alternative ways to improve GFlowNet performance via policy-based training. Our contributions can be summarized as follows:

- We reformulate the GFlowNet training problem as RL over a special MDP where the reward is policy-dependent, and the underlying Markovian Chain is absorbing. We further derive policy gradients for this special MDP, and propose policy-based training strategies for GFlowNets, inspired by policy gradient and TRPO methods for stationary reward.
- We further formulate the design of backward policies in GFlowNets as an RL problem and propose a coupled training strategy. While finding a desired forward policy is the goal of GFlowNet training, well-designed backward policies, as the components of training objectives, are expected to improve training efficiency (Shen et al., 2023).
- We provide performance analyses for theoretical guarantees of our method for GFlowNet training. Our theoretical results are also accompanied by performing experiments in three application domains, hyper-grid modeling, sequence design, and Bayesian Network (BN) structure learning. The obtained experimental results serve as empirical evidence for the validity of our work and also help empirically understand the relationship between GFlowNet training and RL.

2. Preliminaries

In a DAG, $\mathcal{G} := (\mathcal{S}, \mathcal{A})$, modeling a MDP: $s \in \mathcal{S}$ denotes a state and $a \in \mathcal{A}$ denotes a directed edge/action ($s \rightarrow s'$) and $\mathcal{A} \subseteq \mathcal{S} \times \mathcal{S}$. Assuming that there is a topological ordering $\mathcal{S}_0, \dots, \mathcal{S}_T$ for $T + 1$ disjoint subsets of \mathcal{S} , then $\mathcal{S} = \bigcup_{t=0}^T \mathcal{S}_t$ and an element of \mathcal{S}_t is denoted as s_t . We use $\{<, >, \preceq, \succeq\}$ to define the partial orders between states; for example, $\forall t < t' : s_t < s_{t'}$. Furthermore, being *acyclic* means $\forall (s \rightarrow s') \in \mathcal{A} : s < s'$. For any $s \in \mathcal{S}$, we denote its parent set by $Pa_{\mathcal{G}}(s) = \{s' | (s' \rightarrow s) \in \mathcal{A}\}$ and its child set $Ch_{\mathcal{G}}(s) = \{s' | (s \rightarrow s') \in \mathcal{A}\}$. Correspondingly, We denote the edge sets that start and end at s as $\mathcal{A}(s) = \{(s \rightarrow s') | s' \in Ch_{\mathcal{G}}(s)\}$ and $\dot{\mathcal{A}}(s) = \{(s' \rightarrow s) | s' \in Pa_{\mathcal{G}}(s)\}$ respectively. The complete trajectory set is defined as $\mathcal{T} = \{\tau = (s_0 \rightarrow \dots \rightarrow s_T) | \forall (s \rightarrow s') \in \tau : (s \rightarrow s') \in \mathcal{A}\}$. We use $\tau_{\succeq s}$ to denote the sub-trajectory that starts at s , and $\tau_{\geq t}$ the sub-trajectory that starts at s_t .

For the DAG \mathcal{G} in GFlowNets, we have two special states: the initial state s^0 with $Pa(s^0) = \emptyset$ and $\mathcal{S}_0 = \{s^0\}$, and the final state s^f with $Ch(s^f) = \emptyset$ and $\mathcal{S}_T = \{s^f\}$. Furthermore, the terminal state set, \mathcal{S}_{T-1} , covering the object set \mathcal{X} with a reward function $R : \mathcal{X} \rightarrow \mathbb{R}^+$. For notation compactness, we restrict \mathcal{G} to be graded¹ so that $\forall (s \rightarrow s') \in \mathcal{A}, s \in \mathcal{S}_t : s' \in \mathcal{S}_{t+1}$. Accordingly, \mathcal{A} can be decomposed into $\bigcup_{t=0}^{T-1} \mathcal{A}_t$ where $\mathcal{A}_t \cap \mathcal{A}_{t' \neq t} = \emptyset$ and $a_t \in \mathcal{A}_t$ represents an edge $(s_t \rightarrow s_{t+1})$ connecting \mathcal{S} and \mathcal{S}_{t+1} . The complete trajectory set is $\mathcal{T} = \{\tau = (s_0 \rightarrow \dots \rightarrow s_T) | s_0 = s^0, s_T = s^f; \forall t \in \{1, \dots, T\} : (s_{t-1} \rightarrow s_t) \in \mathcal{A}_{t-1}\}$.

2.1. GFlowNets

GFlowNets aim at efficient sampling from $P^*(x) := \frac{R(x)}{Z^*}$, where $Z^* = \sum_{x \in \mathcal{X}} R(x)$ and directly computing Z^* is often challenging with typically large $|\mathcal{X}|$. To achieve this, GFlowNets define a probability measure $F(\tau) : \mathcal{T} \rightarrow \mathbb{R}^+$, termed as ‘flow’ (Bengio et al., 2023), so that for any event E , $F(E) = \sum_{\tau \in E} F(\tau)$ and the total flow $Z = F(s^0) = F(s^f)$. For any event E and E' , $P(E) := F(E)/Z$ and $P(E|E') := \frac{F(E \cup E')}{F(E')}$. Furthermore, F is restricted to be Markovian, which means $\forall \tau \in \mathcal{T}$:

$$P(\tau) = \prod_{t=1}^T P(s_{t-1} \rightarrow s_t | s_{t-1}) := \prod_{t=1}^T \frac{F(s_{t-1} \rightarrow s_t)}{F(s_{t-1})}, \quad (1)$$

where $F(s \rightarrow s') = \sum_{\tau \in \{\tau | (s \rightarrow s') \in \tau\}} F(\tau)$, $F(s) = \sum_{\tau \in \{\tau | s \in \tau\}} F(\tau)$ and $P_F(s_t | s_{t-1}) := P(s_{t-1} \rightarrow s_t | s_{t-1})$. Similarly, $P_B(s_{t-1} | s_t) := P(s_{t-1} \rightarrow s_t | s_t) = \frac{F(s_{t-1} \rightarrow s_t)}{F(s_t)}$. A desired generative flow F is set to have the terminal transition probability, $P^T(x) := P(x \rightarrow s_f)$, to

¹Any DAG can be equivalently converted to be graded by adding dummy non-terminating states. Please refer to Appendix A of Malkin et al. (2022b) for more details.

be $P^*(x)$. As shown in Bengio et al. (2023), the necessary and sufficient condition is that $\forall s' \in \mathcal{S} \setminus \{s^0, s^f\}$:

$$\sum_{s \in Pa(s')} F(s \rightarrow s') = \sum_{s'' \in Ch(s)} F(s' \rightarrow s''). \quad (2)$$

where we clamp $F(x \rightarrow s_f) = R(x)$ for any $x \in \mathcal{X}$.

2.2. GFlowNet training

Directly estimating the transition flow $F(s \rightarrow s')$ via the flow matching objective (Bengio et al., 2021) can suffer from the explosion of F values, of which the numerical issues may lead to the failure of model training. In practice, the Trajectory Balance (TB) objective has been shown to achieve the state-of-the-art training performance (Malkin et al., 2022a). With the TB objective, the desired flow is estimated by the total flow Z and a pair of forward/backward policies, $P_F(s'|s)$ and $P_B(s|s')$. The TB objective $\mathcal{L}_{TB}(P_D)$ of a data sampling policy $P_D(s'|s)$ is defined as:

$$\begin{aligned} \mathcal{L}_{TB}(P_D) &:= \mathbb{E}_{P_D(\tau)}[L_{TB}(\tau)], \\ L_{TB}(\tau) &:= \left(\log \frac{P_F(\tau|s_0)Z}{P_B(\tau|x)R(x)} \right)^2. \end{aligned} \quad (3)$$

In the equation above, $P_F(\tau|s_0) = \prod_{t=1}^T P_F(s_t|s_{t-1})$ with $P_F(\tau) = P_F(\tau|s_0)$, $P_F^\top(x) := P_F(x \rightarrow s^f)$, and $P_F(\tau|x) = P_F(\tau)/P_F^\top(x)$. Correspondingly, $P_B(\tau|x) = \prod_{t=1}^{T-1} P_B(s_{t-1}|s_t)$, $P_B^\top(x) := P^*(x)$, $P_B(\tau) := P_B^\top(x)P_B(\tau|x)$, and $P_B(\tau|s_0) = P_B(\tau)$. Furthermore, we define $\mu(s_0 = s^0) := Z/\hat{Z}$ as the distribution over \mathcal{S}_0 so that $P_{F,\mu}(\tau) := P_F(\tau|s_0)\mu(s_0) = P_F(\tau)$, where \hat{Z} is the normalizing constant whose value is clamped to Z . We define $P_{B,\rho}(\tau) := P_B(\tau|x)\rho(x)$ with an arbitrary distribution ρ over \mathcal{X} .

3. Policy gradients for GFlowNet training

Following Malkin et al. (2022b), we first extend the relationship between the GFlowNet training methods based on the TB objective and KL divergence. With the extended equivalence, we then introduce our policy-based and coupled training strategies for GFlowNets. Finally, we present theoretical analyses on our proposed strategies.

3.1. Gradient equivalence

When choosing trajectories sampler $P_D(\tau) = P_F(\tau)$, the gradient equivalence between using the KL divergence and TB objective has been proven (Malkin et al., 2022b). However, this forward gradient equivalence does not take the total flow estimator Z into account. Moreover, the backward gradient equivalence requires computing the expectation over $P^*(x)$, which is not feasible. In this work, we extend the proof of the gradient equivalence to take all gradients

into account and remove the dependency on $P^*(x)$, while keeping feasible computation.

Proposition 3.1. *Given a parametrized forward policy $P_F(\cdot|\cdot; \theta)$, a backward policy $P_B(\cdot|\cdot; \phi)$, and a total flow estimator $Z(\theta)$, the gradient of the TB objective² can be written as:*

$$\begin{aligned} \frac{\nabla_\theta \mathcal{L}_{TB}(P_{F,\mu}; \theta)}{2} &= \nabla_\theta D_{KL}^{\mu(\cdot; \theta)}(P_F(\tau|s_0; \theta), P_B(\tau|s_0)) \\ &\quad + \frac{1}{2} \nabla_\theta (\log Z(\theta) - \log Z^*)^2 \\ &= \nabla_\theta D_{KL}^{\mu(\cdot; \theta)}(P_F(\tau|s_0; \theta), \tilde{P}_B(\tau|s_0)); \\ \frac{\nabla_\phi \mathcal{L}_{TB}(P_{B,\rho}; \phi)}{2} &= \nabla_\phi D_{KL}^\rho(P_B(\tau|x; \phi), P_F(\tau|x)) \\ &= \nabla_\phi D_{KL}^\rho(P_B(\tau|x; \phi), \tilde{P}_F(\tau|x)). \end{aligned} \quad (4)$$

In the equations above, $\tilde{P}_F(\tau|x) = P_F(\tau)$ and $\tilde{P}_B(\tau|s_0) := \frac{R(x)}{Z} P_B(\tau|x)$, denoting two unnormalized distributions of $P_F(\tau|x)$ and $P_B(\tau|s_0)$. For arbitrary distributions p, q , and u , $D_{KL}^u(p(\cdot|s), q(\cdot|s)) := \mathbb{E}_{u(s)}[D_{KL}(p(\cdot|s), q(\cdot|s))]$.

The proof is provided in Appendix A.1. As the TB objective is a special case of the Sub-Trajectory Balance (Sub-TB) objective (Madan et al., 2023), we also provide the proof of the gradient equivalence with respect to the Sub-TB objective in Appendix A.3, where the initial distribution μ becomes more flexible.

3.2. RL formulation of GFlowNet training

Inspired by the equivalence relationship in Proposition 3.1, we propose new reward functions that allow us to formulate GFlowNet training as RL problems with a corresponding policy-based training strategy.

Definition 3.2 (Policy-dependent Rewards). For any action $a = (s \rightarrow s') \in \mathcal{A}(s)$ ($a \in \dot{\mathcal{A}}(s')$), we define two reward functions as:

$$\begin{aligned} R_F(s, a; \theta) &:= \log \frac{\pi_F(s, a; \theta)}{\pi_B(s', a)}, \\ R_B(s', a; \phi) &:= \log \frac{\pi_B(s', a; \phi)}{\pi_F(s, a)}, \end{aligned} \quad (5)$$

where $\pi_F(s, a; \theta) := P_F(s'|s; \theta)$, $\pi_B(s', a; \phi) := P_B(s|s'; \phi)$, $\pi_B(x, a)$ is clamped to $R(x)/Z$ for $a = (x \rightarrow s_f)$. For any $a \notin \mathcal{A}(s)$, $R_F(s, a) := 0$. For any $a \notin \dot{\mathcal{A}}(s')$, $R_B(s', a) := 0$.

Tuples $(\mathcal{S}, \mathcal{A}, \mathcal{G}, R_F)$ and $(\mathcal{S}, \dot{\mathcal{A}}, \mathcal{G}, R_B)$ specify two Markov Decision Processes (MDPs) with policy-dependent

²We note that training via TB was intrinsically done in an off-policy setting, so $\nabla \mathcal{L}_{TB}(P_D) = \mathbb{E}_{P_D(\tau)}[\nabla L_{TB}(\tau)]$ for any choice of P_D .

rewards. In the MDPs, \mathcal{G} specifies a deterministic transition environment such that $P(s'|s, a) = \mathbb{I}[(s \rightarrow s') = a]$ with the indicator function \mathbb{I} . (\mathcal{G}, π_F) and (\mathcal{G}, π_B) corresponds to two absorbing Markovian chains. Accordingly, the nature of DAGs requires that each state has only one order index, and allows us to define time-invariant expected value functions of states and state-action pairs, which are defined as $V_F(s) := \mathbb{E}_{P_F(\tau_{>t}|s_t)}[\sum_{l=t}^{T-1} R_F(s_l, a_l) | s_t = s]$ and $Q_F(s, a) := \mathbb{E}_{P_F(\tau_{>t+1}|s_t, a_t)}[\sum_{l=t}^{T-1} R_F(s_l, a_l) | s_t = s, a_t = a]$. Then we define $J_F := \mathbb{E}_{\mu(s_0)}[V_F(s_0)]$, $A_F(s, a) := Q_F(s, a) - V_F(s)$, and $d_{F,\mu}(s) := \frac{1}{T} \sum_{t=0}^{T-1} P_F(s_t = s)$. We likewise denote the functions for the backward policy as $\{V_B, Q_B, J_B, A_B, d_{B,\rho}\}$. More details are provided in Appendix B.1. By definition, $V_F(s_0) = \mathbb{E}_{P(\tau|s_0)}[\sum_{t=0}^{T-1} R_F(s_t, a_t)] = D_{KL}(P_F(\tau|s_0), \tilde{P}_B(\tau|s_0))$, so $J_F = D_{KL}^\mu(P_F(\tau|s_0), \tilde{P}_B(\tau|s_0))$. Likewise, we can obtain $J_B = D_{KL}^\rho(P_B(\tau|x), \tilde{P}_F(\tau|x))$. We can conclude that GFlowNet training can be converted into minimizing the expected value function J_F and J_B by Proposition 3.1. With the derived ∇J_F and ∇J_B provided in Appendix B.3, we update π_F , π_B , and μ to minimize J_F and J_B based on the correspondingly computed gradients of the following two objectives:

$$\begin{aligned} T \cdot \mathbb{E}_{d_{F,\mu}(s), \pi_F(s, a; \theta)} [A_F(s, a)] + \mathbb{E}_{\mu(s_0; \theta)} [V_F(s_0)], \\ T \cdot \mathbb{E}_{d_{B,\rho}(s), \pi_B(s, a; \phi)} [A_B(s, a)]. \end{aligned} \quad (6)$$

Our policy-based method generalizes the TB-based training with $\pi_D = \pi_F$ as follows: TB-based training corresponds to approximating $A(s, a)$ empirically by $\hat{Q}_F(s, a) - C$, where $\hat{Q}_F(s, a) = \sum_{l=t}^{T-1} R_F(s_l, a_l)$ for $(s_t = s, a_t = a)$, and C is constant baseline for variance reduction. For comparison, our policy-based method can be considered approximating $A_F(s, a)$ functionally by $\hat{A}_F^\lambda(s, a) = \sum_{l=t}^{T-1} \lambda^{l-t} (\hat{Q}_F(s_l, a_l) - \tilde{V}_F(s_l))$, where $\lambda \in [0, 1]$ controls the **bias-variance trade-off** for gradient estimation (Schulman et al., 2016), $\hat{Q}_F(s_t, a_t) = R_F(s_t, a_t) + \tilde{V}_F(s_{t+1})$, and $\tilde{V}_F(s_t)$ is a functional approximation of exact $V_F(s_t)$ serving as a functional baseline. Specifically, our policy-based method with $\lambda = 1$ can provide unbiased gradient estimation of ∇J_F as the TB-based method. This supports the stability of our policy-based method with the theoretical convergence guarantee by Theorem 3.7 in Section 3.4. A formal discussion of their connection is provided in Appendix B.4 and B.6. Additionally, we discuss the relationship between our method and traditional Maximum Entropy (MaxEnt) RL in Appendix B.5

To further exemplify that the proposed rewards bridge policy-based RL techniques to GFlowNet training, we specifically focus on the TRPO method, whose performance is usually more stable than vanilla policy-based methods due

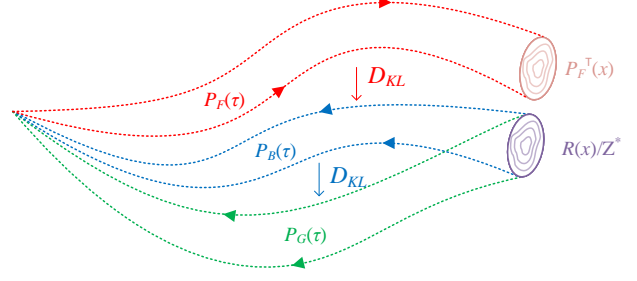


Figure 1. Dotted lines illustrate the spanning range of trajectories for one-dimensional states. P_B and P_G share the ground-truth terminating distribution $R(x)/Z^*$. When pushing P_F to match P_B trajectory-wise, $P_F^\top(x)$ will also be pushed to match $R(x)/Z^*$.

to conservative model updating rules (Schulman et al., 2015; Achiam et al., 2017). Likewise, we propose a TRPO-based objective for updating π_F :

$$\begin{aligned} \min_{\theta'} T \cdot \mathbb{E}_{d_{F,\mu}(s; \theta), \pi_F(s, a; \theta')} [A_F(s, a; \theta)] \\ \text{s.t. } D_{KL}^{d_{F,\mu}(\cdot; \theta)}(\pi_F(s, a; \theta), \pi_F(s, a; \theta')) \leq \zeta_F. \end{aligned} \quad (7)$$

The objective for π_B can be defined similarly and is omitted here. This objective is motivated as the approximation of the upper bound in Theorem 3.6, which generalizes the original results for static rewards and absorbing MDPs. We defer the discussion of their relationship in Section 3.4.

Moreover, the model parameter updating rule based on ∇J_F can be written as $\theta' \leftarrow \theta - \alpha \nabla_{\theta} J_F(\theta)$ or equivalently $\theta' = \arg\min_{\theta'} (\nabla J_F)^T \theta' + \frac{1}{2\alpha} \|\theta' - \theta\|^2$, where $\frac{1}{2\alpha} \|\theta - \theta'\|$ can be generalized to KL divergence as in TRPO, or Bregman divergence corresponding to natural or mirror policy gradients, which we leave for future work.

Details of the model parameter updating rules for proposed methods are provided in Appendix B.6.

3.3. RL formulation of guided backward policy design

During GFlowNet training, (P_B, R) specifies the amount of desired flow that (P_F, Z) is optimized to match. While $P_B(\cdot|\cdot)$ can be chosen free in principle (Bengio et al., 2023), a well-designed P_B that assigns high probabilities over sub-trajectories preceding the terminating state x with a high reward value $R(x)$, will improve training efficiency. Following Shen et al. (2023), we formulate the design problem as minimizing the following objective:

$$\begin{aligned} \mathcal{L}_{TB}^G(P_B^\rho) &:= \mathbb{E}_{P_B^\rho(\tau)} [L_{TB}^G(\tau)], \\ L_{TB}^G(\tau; \phi) &:= \left(\log \frac{P_B(\tau|x; \phi)}{P_G(\tau|x)} \right)^2, \end{aligned} \quad (8)$$

where $P_G(\tau|x) = \prod_{t=1}^{T-1} P_G(s_{t-1}|\tau_{\geq t})$ is called the guided backward trajectory distribution, which is usually non-Markovian, and $P_G(\tau) = P_G(\tau|x)R(x)/Z^*$. As required

by the training w.r.t. P_F , the objective \mathcal{L}_{TB-G} aims at finding the backward policy whose Markovian flow best matches the non-Markovian flow induced by P_G^3 .

Proposition 3.3. *Given P_G and $P_B(\cdot|\cdot; \phi)$, the gradients of \mathcal{L}_{TB-G} can be written as:*

$$\frac{\nabla_{\phi} \mathcal{L}_{TB}^G(P_B^{\rho}; \phi)}{2} = \nabla_{\phi} D_{KL}^{\rho}(P_B(\tau|x; \phi), P_G(\tau|x)). \quad (9)$$

The proof can be found in Appendix A.2. Based on the proposition, we propose a new reward that allows us to formulate the backward policy design problem as an RL problem.

Definition 3.4. Given a guided backward trajectory distribution $P_G(\tau|x)$, we define a reward function for any action $a := (s \rightarrow s') \in \mathcal{A}(s')$ as:

$$R_B^G(s', a; \phi) := \log \frac{\pi_B(s', a; \phi)}{\pi_G(s', a)}, \quad (10)$$

where $\pi_G(s', a) := P_G(s|\tau_{\leq s'})$. For any $a \notin \mathcal{A}(s')$, $R_B^G(s', a) := 0$.

Accordingly, we denote the associated function set as $\{V_B^G, Q_B^G, J_B^G, A_B^G, \delta_B^G, d_{B,\rho}^G\}$, which are defined in a similar way as R_B but replacing P_F by P_G .

By the definition of J_B^G and Proposition 3.3, we can conclude that $\nabla_{\phi} J_B^G(\phi) = \frac{1}{2} \nabla_{\phi} \mathcal{L}_{TB-G}(P_B^{\rho}; \phi)$ and the design of backward policy can be solved by minimizing J_B^G . The form of P_G will be detailed in the experiment section for the corresponding tasks.

In principle, following the pipeline by Shen et al. (2023), we need to solve the optimization of \mathcal{L}_{TB-G} to find the desired P_B at first. Then, freezing P_B , we can optimize \mathcal{L}_{TB} to find the desired P_F . This gives rise to training inconvenience in practice. To avoid doing two-phase training, the authors mixed P_B and P_G by $\alpha P_B + (1 - \alpha) P_G$ within the training objective w.r.t P_F . This operation, however, lacks theoretical guarantees as the mixed distribution is still non-Markovian. By comparison, the RL formulation allows us to optimize J_F and J_B^G jointly with a theoretical performance guarantee, which we defer to the next section.

The workflow of our coupled training strategy is summarized in Algorithm 1 and depicted by Fig. 1.

3.4. Performance Analysis

In the previous sections, we formulate two RL problems with respect to R_F and R_B^G . Now, we show below that the two problems can be solved jointly.

³By non-Markovian assumption, $P_G(\tau|x)$ can factorize in arbitrary ways conditioning on x . Here it is assumed to factorize in the backward direction for notation compactness.

Algorithm 1 GFlowNet Training Workflow

Require: $P_F(\cdot|\cdot; \theta)$, $Z(\theta)$, $P_B(\cdot|\cdot; \phi)$, $P_G(\cdot|\cdot)$
for $n = \{1, \dots, N\}$ **do**
 $\mathcal{D} \leftarrow \{\hat{\tau} | \hat{\tau} \sim P_F(\tau; \theta)\}$
 Update θ w.r.t. R_F and \mathcal{D}
if $\phi \neq \emptyset$ **then**
 $\dot{\mathcal{D}} \leftarrow \{\hat{\tau} | \forall x \in \mathcal{D} : \hat{\tau}|x \sim P_B(\tau|x)\}$
if $P_G(\hat{\tau}|x) \neq P_B(\hat{\tau}|x)$ **then**
 Update ϕ w.r.t. R_B^G and $\dot{\mathcal{D}}$
else
 Update ϕ w.r.t. R_B and $\dot{\mathcal{D}}$
end if
end if
end for

Theorem 3.5. Denoting J_F^G the corresponding function of R_F when $P_B = P_G$ and choosing $\rho(x) = P_F^{\top}(x)$, J_F^G , J_F and J_B^G satisfy the following inequality:

$$J_F^G \leq J_F + J_B^G + R_B^{G, \max} \sqrt{\frac{(J_F + \log Z^* - \log Z)}{2}}, \quad (11)$$

where $R_B^{G, \max} = \max_{s,a} R_B^G(s, a)$.

The proof is given in Appendix C.1. As shown in Proposition 3.1, minimization of J_F will incur the decrease of $D_{KL}^{\mu}(P_F(\tau|s_0), P_B(\tau|s_0)) = J_F + \log Z^* - \log Z$. Thus, by minimizing J_F and J_B^G jointly, the upper bound of J_F^G decreases.

Moreover, the TRPO-based objective introduced in the previous section is motivated by the following upper bounds.

Theorem 3.6. For two forward policies (π_F, π'_F) with $\max_s D_{KL}(\pi'_F(\cdot, s), \pi_F(\cdot, s)) < \zeta_F$, and two backward policies (π_B, π'_B) with $\max_s D_{KL}(\pi'_B(\cdot, s), \pi_B(\cdot, s)) < \zeta_B$,

$$\begin{aligned} \frac{(J'_F - J_F)}{T} &\leq \mathbb{E}_{s,a \sim d_{F,\mu}, \pi'_F} [A_F(s, a)] + \epsilon_F \sqrt{2\zeta_F} + \zeta_F, \\ \frac{(J'_B - J_B)}{T} &\leq \mathbb{E}_{s,a \sim d_{B,\rho}, \pi'_B} [A_B(s, a)] + \epsilon_B \sqrt{2\zeta_B} + \zeta_B, \end{aligned} \quad (12)$$

where $\epsilon_F = \max_s |\mathbb{E}_{\pi'_F(s,a)} [A_F(s, a)]|$ and $\epsilon_B = \max_s |\mathbb{E}_{\pi'_B(s,a)} [A_B(s, a)]|$. Similar results also apply to J_B^G and A_B^G for the backward policy π_B .

The proof is given in Appendix C.2. The TRPO-based objective can be derived following the similar logic in Schulman et al. (2015) and Achiam et al. (2017): in the worst case, we choose $\pi'_F = \pi_F$ and $J'_F - J_F = 0$; When we make a conservative decision, this is, $\pi'_F \neq \pi_F$ and ζ_F is negligibly small, it can be expected that $\mathbb{E}_{s,a \sim d_{F,\mu}, \pi'_F} [A_F(s, a)] < 0$, thereby resulting in $J'_F - J_F < 0$; therefore, TRPO method

usually provides more stable performance gain than the vanilla policy-based method.

Lastly, we provide a theoretical guarantee that policy-based methods with policy-dependent rewards can asymptotically converge to stationary points, which draws inspiration from the results for static rewards by Agarwal et al. (2019).

Theorem 3.7. *Suppose that: $J_F(\theta)$ is β -smooth; $\mathbb{E}_{P(\cdot|\theta)}[\widehat{\nabla}_\theta J_F(\theta)] = \nabla_\theta J_F(\theta)$; the estimation variance, $\mathbb{E}_{P(\cdot|\theta)}[\|\widehat{\nabla}_\theta J_F(\theta) - \nabla_\theta J_F(\theta)\|^2] \leq \sigma_F$; $|\log Z(\theta) - \log Z^*| \leq \sigma_Z$; we update θ for N ($> \beta$) iterations by $\theta_{n+1} \leftarrow \theta_n - \alpha \widehat{\nabla}_\theta J_F(\theta_n)$ with $n \in \{0, \dots, N-1\}$, $\alpha = \sqrt{1/(\beta N)}$ and initial parameter θ_0 . Then we have:*

$$\begin{aligned} \min_{n \in \{0, \dots, N-1\}} \mathbb{E}_{P(\theta_n)} [\|\nabla_{\theta_n} J_F(\theta_n)\|^2] \\ \leq \frac{\sigma_F + \sigma_Z + \mathbb{E}_{P(\theta_0)}[J_F(\theta_0)]}{(\sqrt{(2N)/\beta} - 1)}. \end{aligned} \quad (13)$$

Similar results also apply to J_B and J_B^G .

The proof is provided in Appendix C.3. The assumption $\mathbb{E}_{P(\cdot|\theta)}[\widehat{\nabla}_\theta J_F(\theta)] = \nabla_\theta J_F(\theta)$ means gradient estimation is unbiased as explained in Appendix B.6.

3.5. Related Work

GFlowNet training GFlowNets were first proposed by Bengio et al. (2021) and trained by an Flow Matching (FM) objective, which aims at minimizing the mismatch of (2) w.r.t. a parameterized state flow estimator $F(s)$ and a parameterized edge flow estimator $F(s \rightarrow s')$ directly. Bengio et al. (2023) reformulated (2) and proposed a Detailed Balance (DB) objective, where edge flows $F(s \rightarrow s')$ are represented by $F(s)P_F(s'|s)$ or $F(s')P_B(s|s')$. Malkin et al. (2022a) claimed that the FM and DB objectives are prone to inefficient credit propagation across long trajectories and showed that the TB objective is the more efficient alternative. Madan et al. (2023) proposed a Sub-TB objective that unified the TB and DB objectives as special cases. They can be considered as Sub-TB objectives with sub-trajectories, which are complete or of length 1 respectively. Zimmermann et al. (2022) proposed KL-based training objectives and Malkin et al. (2022b) first established the equivalence between the KL and TB objectives. Shen et al. (2023) analyzed how the TB objective helps to learn the desired flow under the sequence prepend/append MDP setting, and proposed a guided TB objective.

Hierarchical Variational inference Hierarchical Variational Inference (HVI) (Vahdat & Kautz, 2020; Zimmermann et al., 2021) generalizes amortized VI (Zhang et al., 2018) to better explore specific statistical dependency structures between observed variables and latent variables by

introducing the hierarchy of latent variables. Training HVI models typically involves minimizing the selected divergence measures between the target distribution and the variational distribution parametrized by DNNs (Kingma & Welling, 2014; Burda et al., 2015). GFlowNets can be considered as a special HVI model, where non-terminating states are latent variables, the hierarchy corresponds to a DAG, and the task of minimizing divergences is achieved by keeping flow balance (Malkin et al., 2022b). Our work provides another view of divergence minimization by interpreting the divergence as the expected accumulated reward.

Policy-based RL Policy-based RL optimizes the expected value function J directly based on policy gradients (Sutton et al., 1999). The most relevant policy-based methods are the Actor-Critic method (Sutton & Barto, 2018) and Trust Region Policy Optimization (TRPO) (Schulman et al., 2015) along with its extension – Constrained Policy Optimization (CPO) (Achiam et al., 2017). Compared to our methods, they work under the assumption that the reward functions must be fixed w.r.t. policies. The underlying Markov chains are further assumed to be ergodic by CPO and TRPO.

MaxEnt RL Bengio et al. (2021) have shown that directly applying MaxEnt RL with $R(s, a)$ defined based on the terminating state is problematic as it corresponds to modeling $p(x) \propto n(x)R(x)$, where $n(x)$ is the number of trajectories that can pass through x . As discussed in Appendix B.5, our policy-based methods, when fixing $\log \pi_B(s', a)$ and $\log Z$, can be connected to Soft-Q learning, a typical MaxEnt RL method.

Imitation learning Imitation learning in RL is to learn a policy that mimics the expert demonstrations with limited expert data, by minimizing the gap between the learned policy and expert policy measured by the 0 – 1 loss or Jensen–Shannon divergence empirically (Rajaraman et al., 2020; Ho & Ermon, 2016). For GFlowNet training in this work, we reduce the gap between the forward policy and the expert forward policy at the trajectory level, as the expert trajectory distribution is equal to $P_B(\tau)$, implicitly encouraging the learned policy to match the desired expert policy.

Bi-level optimization Our proposed training strategy can also be seen as a Stochastic Bi-level Optimization method for GFlowNet training (Ji et al., 2021; Hong et al., 2023; Ghadimi & Wang, 2018). The inner problem is the RL problem w.r.t. R_B or R_B^G for designing backward policies. The outer problem is the RL problem w.r.t. R_F for forward policies. For gradient-based solutions to Bi-level optimization in general, the learning rate of inner problems is carefully selected to guarantee the overall convergence, which is not required in our methods designed for GFlowNet training.

4. Experiments

To compare our policy-based training strategies for GFlowNets with the existing value-based methods based on the TB and DB objectives, we have conducted three simulated experiments for hyper-grid modeling, two real-world experiments for biological and molecular sequence design, one on Bayesian Network structure learning, and ablation study of λ . The detailed descriptions of experimental settings can be found in Appendix D. We compare the performance of GFlowNets by the following training strategies: (1) DB-U: DB-based training strategy with a fixed uniform π_B ; (2) DB-B: DB-based training strategy with a parameterized π_B ; (3) TB-U: TB-based training strategy with a fixed uniform π_B ; (4) TB-B: TB-based training strategy with a parameterized π_B ; (5) RL-U: our policy-based training strategy with a fixed uniform π_B ; (6) RL-B: our policy-based training strategy with a parameterized π_B ; (7) RL-T: our TRPO-based training strategy with a fixed uniform π_B ; and (8) RL-G: our joint training strategy based on policy gradients with guided policy π_G . Total variation D_{TV} , Jensen–Shannon divergence D_{JSD} , and mode accuracy Acc are used to measure the gap between $P_F^\top(x)$ and $P^*(x)$. Their definitions are provided in Appendix D.1.

4.1. Hyper-grid Modeling

In this set of experiments, we use the hyper-grid environment following Malkin et al. (2022b). In terms of GFlowNets, states are the coordinate tuples of an D -dimensional hyper-cubic grid with heights equal to N . The initial state s^0 is $\mathbf{0}$. Starting from s^0 , actions correspond to increasing one of D coordinates by 1 for the current state or stopping the process at the current state and outputting it as the terminating state x . A manually designed reward function $R(\cdot)$ assigns high reward values to some grid points while assigning low values to others. We conduct experiments on 256×256 , 128×128 , $64 \times 64 \times 64$, and $32 \times 32 \times 32 \times 32$ grids. For performance evaluation, $P_F^\top(x)$ is computed exactly by dynamic programming.

The training curves by D_{KL} across five runs for 256×256 and 128×128 grids are plotted in Fig. 2 with Table 1 in Appendix D.2 reporting the mean and standard deviation of metric values at the last iteration. The graphical illustrations of $P_F^\top(x)$ are shown in Figs. 10 and 11 in Appendix D.6. **In the first setting**, it can be observed that our policy-based RL-U or RL-B, in terms of convergence rate and converged D_{KL} is much better than the existing GFlowNet training methods by DB-U, TB-U, DB-B, or TB-B. This shows that our policy-based training strategies give a more robust gradient estimation. Besides, RL-G achieves the smallest D_{KL} and converges much faster than all the other competing methods. In RL-G, the guided distribution assigns small values to the probability of terminating at coordinates with

low rewards. This prevents the forward policy from falling into the reward ‘desert’ between the highly isolated modes. Finally, RL-T outperforms DB-U, TB-U, DB-B, TB-B and RL-U and it behaves more stably than RL-U during training. This confirms that with the help of trust regions, the gradient estimator becomes less sensitive to environment noises. Here we use a fixed constant ζ_F for trust region control. It is expected that using a proper scheduler of ζ_F during training may further improve the performance of RL-T. **In the second setting** The converged D_{KL} of RL-U, RL-B, TB-U and TB-B are similar and significantly better than those of DB-U and DB-B. As expected, RL-U and RL-B converge much faster than DB-U, TB-U, DB-B, and TB-B. Thus, the results further support the effectiveness of our policy-based methods. Besides, RL-G and RL-T achieve the second-best and the best performance and RL-T shows faster convergence and better stability than RL-G. This again shows the superiority of coupled and TRPO-based strategies, confirming our theoretical analysis conclusions. More results and discussions for $64 \times 64 \times 64$ and $32 \times 32 \times 32 \times 32$ grids can be found in Appendix D.2.

4.2. Biological and Molecular Sequence Design

In this set of experiments, we use GFlowNets to generate nucleotide strings of length D and molecular graphs composed of D blocks according to given rewards. The initial state $s^0 = \emptyset$ is an empty sequence. The generative process runs as follows: starting from s^0 , an action is taken to pick one of the empty slots and fill it with one element until the sequence is completed. Then the sequence is returned as the terminating state x . We use the QM9 dataset composed of molecular graphs with five blocks, the SIX6 dataset composed of strings of length 8, and the PHO4 dataset composed of strings of length 10 as in Shen et al. (2023).

Following (Shen et al., 2023), the training curves by the mode accuracy Acc for SIX6 and QM9 datasets are shown in Fig. 3, where P_F^\top is computed exactly by dynamic programming. For evaluation consistency, we also provide the curves by D_{TV} as well as the performance metric values at the last iteration summarized in Fig. 6 and Table 3 in Appendix D.3. In both experiments, TB-based and policy-based methods achieve better performance than DB-based methods. While the converged Acc values of TB-based methods and our policy-based methods are similar, the latter converges much faster than TB-based methods with only TB-U achieving a comparable convergence rate. Besides, RL-T has the fastest convergence rates in both experiments. The performances of RL-G are similar to those of RL-B, which has a parameterized π_B , but slightly better than RL-U with a uniform π_B . In summary, experimental results for QM9 and SIX6 datasets align with those of hyper-grid tasks, confirming again the advantage offered by our policy-based methods for robust gradient estimation. More results and

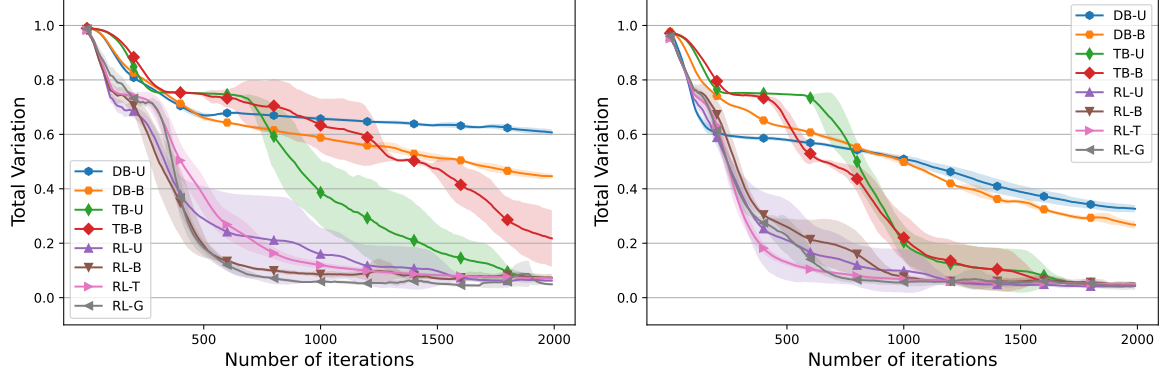


Figure 2. Training curves by D_{TV} between P_F^T and P^* for 256×256 (left) and 128×128 hyper-grids (right). The curves are plotted based on their mean and standard deviation values across five runs.

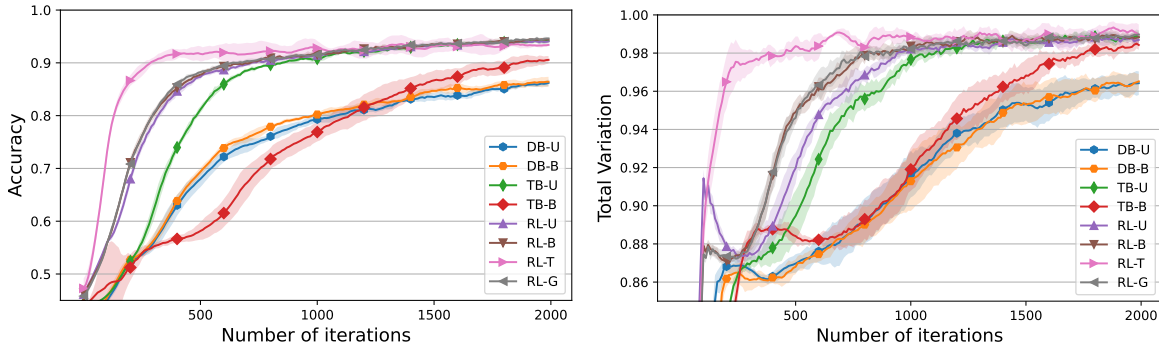


Figure 3. Training curves by Acc of P_F^T w.r.t. P^* for SIX6 (left) and QM9 (right) datasets. The curves are plotted by their mean and standard deviation values across five runs.

discussions for PHO4 dataset can be found in Appendix D.3.

4.3. Ablation Study of λ

To investigate how the setting of λ , which controls the bias-variance-trade-off, may help robust estimation of gradients, we conduct experiments in the 256×256 hyper-grid environment. We compare the performance of RL-U methods with different λ values $\{0.95, 0.96, 0.97, 0.98, 0.99, 1.0\}$ and TB-U methods with the decay rate of the random-mix-in factor taking values $\{0.95, 0.96, 0.97, 0.98, 0.99\}$. The obtained training curves by D_{KL} across five runs are shown in Fig. 5 in Appendix D.2. Among the decay rate choices for TB-U, the values 0.99 and 0.95 yield the best and the worst performances. In contrast, RL-U under all setups except $\lambda = 1$, demonstrates significantly faster convergence than TB-U. It should be pointed out that when $\lambda = 1$, Q_F is approximated empirically as TB-based methods, but V_F is approximated functionally. Additionally, the final performances in all setups achieved by RL-U are better than that of TB-U. These results verify that by controlling λ , our policy-based methods can provide more robust gradient estimation than TB-based methods.

We have also conducted performance comparison using different GFlowNet training methods for Bayesian network

structure learning. The results and discussions can be found in Appendix D.4.

5. Conclusion, Limitations and Future Work

This work bridges the flow-balance-based GFlowNet training to RL problems. We have developed policy-based training strategies, which provide alternative ways to improve training performance compared to the existing value-based strategies. The experimental results support our claims. Our policy-based methods are not limited to the cases where \mathcal{G} must be a DAG as it intrinsically corresponds to minimizing the KL divergence between two distributions. This does not require \mathcal{G} to be a DAG. The consequent work will focus on extending the proposed methods to general \mathcal{G} with the existence of cycles for more flexible modeling of generative processes of object $x \in \mathcal{X}$. While our policy-based training strategies do not require an explicit design of a data sampler and are shown to achieve better GFlowNet training, they may still get trapped into local optima due to the variance of gradient estimation when the state space is very large. Thus, future research will also focus on further improving policy-based methods, with more robust gradient estimation according to the explained gradient equivalence relationship.

6. Impact Statements

The presented research aims at improving GFlowNet training methods to address the convergence efficiency challenge. The implications of our work extend to various societal realms, ranging from medicine to materials design.

References

- Achiam, J., Held, D., Tamar, A., and Abbeel, P. Constrained policy optimization. In *International conference on machine learning*, pp. 22–31. PMLR, 2017.
- Agarwal, A., Jiang, N., Kakade, S. M., and Sun, W. Reinforcement learning: Theory and algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 32, 2019.
- Beck, A. *First-order methods in optimization*. SIAM, 2017.
- Bengio, E., Jain, M., Korablyov, M., Precup, D., and Bengio, Y. Flow network based generative models for non-iterative diverse candidate generation. *Advances in Neural Information Processing Systems*, 34:27381–27394, 2021.
- Bengio, Y., Lahlou, S., Deleu, T., Hu, E. J., Tiwari, M., and Bengio, E. Gflownet foundations. *Journal of Machine Learning Research*, 24(210):1–55, 2023.
- Burda, Y., Grosse, R., and Salakhutdinov, R. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.
- Degrís, T., White, M., and Sutton, R. S. Off-policy actor-critic. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 179–186, 2012.
- Deleu, T., Góis, A., Emezue, C., Rankawat, M., Lacoste-Julien, S., Bauer, S., and Bengio, Y. Bayesian structure learning with generative flow networks. In *Uncertainty in Artificial Intelligence*, pp. 518–528. PMLR, 2022.
- Ghadimi, S. and Wang, M. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- Golpar Raboky, E. and Eftekhari, T. On nilpotent interval matrices. *Journal of Mathematical Modeling*, 7(2):251–261, 2019.
- Grinstead, C. and Snell, L. J. *Introduction to probability*. 2006.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.
- Hestenes, M. R., Stiefel, E., et al. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49(6):409–436, 1952.
- Ho, J. and Ermon, S. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29, 2016.
- Hong, M., Wai, H.-T., Wang, Z., and Yang, Z. A two-timescale stochastic algorithm framework for bilevel optimization: Complexity analysis and application to actor-critic. *SIAM Journal on Optimization*, 33(1):147–180, 2023.
- Ji, K., Yang, J., and Liang, Y. Bilevel optimization: Convergence analysis and enhanced design. In *International conference on machine learning*, pp. 4882–4892. PMLR, 2021.
- Kakade, S. M. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. In Bengio, Y. and LeCun, Y. (eds.), *ICLR*, 2014. URL <http://dblp.uni-trier.de/db/conf/iclr/iclr2014.html#KingmaW13>.
- Kuipers, J., Moffa, G., and Heckerman, D. Addendum on the scoring of gaussian directed acyclic graphical models. 2014.
- Madan, K., Rector-Brooks, J., Korablyov, M., Bengio, E., Jain, M., Nica, A. C., Bosc, T., Bengio, Y., and Malkin, N. Learning GFlowNets from partial episodes for improved convergence and stability. In *International Conference on Machine Learning*, pp. 23467–23483. PMLR, 2023.
- Malkin, N., Jain, M., Bengio, E., Sun, C., and Bengio, Y. Trajectory balance: Improved credit assignment in gflownets. *Advances in Neural Information Processing Systems*, 35:5955–5967, 2022a.
- Malkin, N., Lahlou, S., Deleu, T., Ji, X., Hu, E. J., Everett, K. E., Zhang, D., and Bengio, Y. Gflownets and variational inference. In *The Eleventh International Conference on Learning Representations*, 2022b.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing Atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

- Rajaraman, N., Yang, L., Jiao, J., and Ramchandran, K. Toward the fundamental limits of imitation learning. *Advances in Neural Information Processing Systems*, 33: 2914–2924, 2020.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. High-dimensional continuous control using generalized advantage estimation. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.
- Schulman, J., Chen, X., and Abbeel, P. Equivalence between policy gradients and soft q-learning. *arXiv preprint arXiv:1704.06440*, 2017.
- Shen, M. W., Bengio, E., Hajiramezanali, E., Loukas, A., Cho, K., and Biancalani, T. Towards understanding and improving GFlowNet training. *arXiv preprint arXiv:2305.07170*, 2023.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- Tsitsiklis, J. and Van Roy, B. Analysis of temporal-difference learning with function approximation. *Advances in neural information processing systems*, 9, 1996.
- Vahdat, A. and Kautz, J. Nvae: A deep hierarchical variational autoencoder. *Advances in neural information processing systems*, 33:19667–19679, 2020.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Zhan, W., Cen, S., Huang, B., Chen, Y., Lee, J. D., and Chi, Y. Policy mirror descent for regularized reinforcement learning: A generalized framework with linear convergence. *SIAM Journal on Optimization*, 33(2):1061–1091, 2023.
- Zhang, C., B  tepage, J., Kjellstr  m, H., and Mandt, S. Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):2008–2026, 2018.
- Zimmermann, H., Wu, H., Esmaili, B., and van de Meent, J.-W. Nested variational inference. *Advances in Neural Information Processing Systems*, 34:20423–20435, 2021.
- Zimmermann, H., Lindsten, F., van de Meent, J.-W., and Naesseth, C. A. A variational perspective on generative flow networks. *Transactions on Machine Learning Research*, 2022.

A. Gradient Equivalence

Lemma A.1. (*REINFORCE trick* (Williams, 1992)) Given a random variable u following a distribution $p(\cdot|\psi)$ parameterized by ψ and a arbitrary function f , we have $\nabla_{\psi} \mathbb{E}_{p(u;\psi)}[f(u)] = \mathbb{E}_{p(u;\psi)}[f(u) \nabla_{\psi} \log p(u; \psi)]$.

A.1. Proof of Proposition 3.1

Proof. First of all, we split the parameters of the total flow estimator and forward transition probability and denote them as $Z(\theta_Z)$ and $P_F(\cdot|\cdot; \theta_F)$ respectively. We further define $c(\tau) = \left(\log \frac{P_F(\tau|s_0)}{R(x)P_B(\tau|x)} \right)$.

For the gradients w.r.t. θ_F :

$$\begin{aligned}
 \frac{1}{2} \mathbb{E}_{P_F(\tau|s_0; \theta_F)} [\nabla_{\theta_F} L_{TB}(\tau; \theta_F)] &= \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\tau; \theta_F)} \left[\nabla_{\theta_F} (c(\tau; \theta_F) + \log Z)^2 \right] \\
 &= \mathbb{E}_{P_{F,\mu}(\tau; \theta_F)} [(c(\tau; \theta_F) + \log Z) \nabla_{\theta_F} \log P_F(\tau|s_0; \theta_F)] \\
 &= \mathbb{E}_{P_{F,\mu}(\tau; \theta_F)} [(c(\tau; \theta_F) + \log Z) \nabla_{\theta_F} \log P_F(\tau|s_0; \theta_F)] + \underbrace{\mathbb{E}_{P_{F,\mu}(\tau; \theta_F)} [\nabla_{\theta_F} (c(\tau; \theta_F) + \log Z)]}_{(a)} \\
 &= \mathbb{E}_{\mu(s_0)} [\nabla_{\theta_F} \mathbb{E}_{P_F(\tau|s_0; \theta_F)} [c(\tau; \theta_F) + \log Z]] \\
 &= \mathbb{E}_{\mu(s_0)} [\nabla_{\theta_F} D_{KL}(P_F(\tau|s_0; \theta_F), \tilde{P}_B(\tau|s_0))] \\
 &= \nabla_{\theta_F} D_{KL}^{\mu}(P_F(\tau|s_0; \theta_F), \tilde{P}_B(\tau|s_0)), \tag{14}
 \end{aligned}$$

where (a) is equal to zero as $\mathbb{E}_{P_F(\tau|s_0; \theta_F)} [\nabla_{\theta_F} c(\tau; \theta_F)] = \mathbb{E}_{P_F(\tau|s_0; \theta_F)} [1 \cdot \nabla_{\theta_F} \log P_F(\tau|s_0; \theta_F)] = \nabla_{\theta} \mathbb{E}_{P_F(\tau|s_0; \theta_F)} [1] = 0$ by Lemma A.1. We also have:

$$\begin{aligned}
 \nabla_{\theta_F} D_{KL}^{\mu}(P_F(\tau|s_0; \theta_F), \tilde{P}_B(\tau|s_0)) &= \nabla_{\theta_F} D_{KL}^{\mu}(P_F(\tau|s_0; \theta_F), \tilde{P}_B(\tau|s_0)) + \underbrace{\nabla_{\theta_F} \mathbb{E}_{P_{F,\mu}(\tau; \theta_F)} [\log Z^* - \log Z]}_{= \nabla_{\theta_F} (\log Z^* - \log Z) = 0} \\
 &= \nabla_{\theta_F} D_{KL}^{\mu}(P_F(\tau|s_0; \theta_F), P_B(\tau|s_0)). \tag{15}
 \end{aligned}$$

It should be emphasized that $P_B(\tau)$ is the ground-truth distribution with $P_B(x) := R(x)/Z^*$, while $\tilde{P}_B(\tau)$ is the approximated one with $\tilde{P}_B(x) := R(x)/Z$. The gradients w.r.t. θ_Z can be written as:

$$\begin{aligned}
 \frac{1}{2} \mathbb{E}_{P_F(\tau)} [\nabla_{\theta_Z} L_{TB}(\tau; \theta_Z)] &= \frac{1}{2} \mathbb{E}_{P_F(\tau)} \left[\nabla_{\theta_Z} (c(\tau) + \log Z(\theta_Z))^2 \right] \\
 &= \mathbb{E}_{P_F(\tau)} [(c(\tau) + \log Z(\theta_Z)) \nabla_{\theta_Z} \log Z(\theta_Z)] \\
 &= [D_{KL}(P_F(\tau|s_0), \tilde{P}_B(\tau|s_0))] [\nabla_{\theta_Z} \log Z(\theta_Z)] \\
 &= \nabla_{\theta_Z} \frac{Z(\theta_Z)}{\hat{Z}} D_{KL}(P_F(\tau|s_0), \tilde{P}_B(\tau|s_0)) \\
 &= \nabla_{\theta_Z} D_{KL}^{\mu(\cdot; \theta_Z)}(P_F(\tau|s_0), \tilde{P}_B(\tau|s_0)). \tag{16}
 \end{aligned}$$

Besides, we have:

$$\begin{aligned}
 \nabla_{\theta_Z} D_{KL}^{\mu(\cdot; \theta_Z)}(P_F(\tau|s_0), \tilde{P}_B(\tau|s_0)) &= [\nabla_{\theta_Z} \log Z(\theta_Z)] [D_{KL}(P_F(\tau|s_0; \theta_F), \tilde{P}_B(\tau|s_0)) + \log \frac{Z^*}{Z^*}] \\
 &= [\nabla_{\theta_Z} \log Z(\theta_Z)] [D_{KL}(P_F(\tau|s_0; \theta_F), P_B(\tau|s_0)) + \log \frac{Z(\theta_Z)}{Z^*}] \\
 &= \nabla_{\theta_Z} \frac{Z(\theta_Z)}{\hat{Z}} [D_{KL}(P_F(\tau|s_0; \theta_F), P_B(\tau|s_0))] + \left[\nabla_{\theta_Z} \log \frac{Z(\theta_Z)}{Z^*} \right] \left[\log \frac{Z(\theta_Z)}{Z^*} \right] \\
 &= \nabla_{\theta_Z} D_{KL}^{\mu(\cdot; \theta_Z)}(P_F(\tau|s_0; \theta_F), P_B(\tau|s_0)) + \frac{1}{2} \nabla_{\theta_Z} \left[\log \frac{Z(\theta_Z)}{Z^*} \right]^2. \tag{17}
 \end{aligned}$$

Combining equations (14) and (16), we obtain:

$$\frac{1}{2} \mathbb{E}_{P_F(\tau; \theta)} [\nabla_{\theta} L_{TB}(\tau; \theta)] = \nabla_{\theta} D_{KL}^{\mu(\cdot; \theta)}(P_F(\tau|s_0; \theta), \tilde{P}_B(\tau|s_0)). \tag{18}$$

Combining equations (15) and (17), we obtain:

$$\frac{1}{2} \mathbb{E}_{P_F(\tau; \theta)} [\nabla_{\theta} L_{TB}(\tau; \theta)] = \nabla_{\theta} \left\{ D_{KL}^{\mu(\cdot; \theta)}(P_F(\tau|s_0; \theta), P_B(\tau|s_0)) + \frac{1}{2} (\log Z(\theta_Z) - \log Z^*)^2 \right\}. \quad (19)$$

Now let's consider the backward gradients and denote $c(\tau) = \left(\log \frac{P_B(\tau|x)}{P_F(\tau)} \right)$,

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_{P_{B,\rho}(\tau; \phi)} [\nabla_{\phi} \mathcal{L}_{TB}(\tau)] \\ &= \mathbb{E}_{P_{B,\rho}(\tau; \phi)} [(c(\tau; \phi) + \log R(x) - \log Z) \nabla_{\phi} \log P_B(\tau|x; \phi)] \\ &= \mathbb{E}_{P_{B,\rho}(\tau; \phi)} [c(\tau; \phi) \nabla_{\phi} \log P_B(\tau|x; \phi)] + \underbrace{\mathbb{E}_{\rho(x)} [(\log R(x) - \log Z) \mathbb{E}_{P_B(\tau|x; \phi)} [\nabla_{\phi} \log P_B(\tau|x; \phi)]]}_{=0 \text{ by Lemma A.1}} \\ &= \mathbb{E}_{P_{B,\rho}(\tau; \phi)} [c(\tau; \phi) \nabla_{\phi} \log P_B(\tau|x; \phi)] + \underbrace{\mathbb{E}_{P_{B,\rho}(\tau; \phi)} [\nabla_{\phi} c(\tau; \phi)]}_{=0 \text{ by Lemma A.1}} \\ &= \mathbb{E}_{\rho(x)} [\nabla_{\phi} D_{KL}(P_B(\tau|x; \phi), \tilde{P}_F(\tau|x))] \\ &= \nabla_{\phi} D_{KL}^{\rho}(P_B(\tau|x; \phi), \tilde{P}_F(\tau|x)). \end{aligned} \quad (20)$$

Besides, we have

$$\begin{aligned} \nabla_{\phi} D_{KL}^{\rho}(P_B(\tau|x; \phi), \tilde{P}_F(\tau|x)) &= \nabla_{\phi} D_{KL}^{\rho}(P_{B,\rho}(\tau|x; \phi), \tilde{P}_F(\tau|x)) + \underbrace{\mathbb{E}_{\rho(x)} [\nabla_{\phi} \mathbb{E}_{P_B(\tau|x; \phi)} [\log P_F^{\top}(x)]]}_{=\log P_F^{\top}(x) \nabla_{\phi} 1=0} \\ &= \nabla_{\phi} D_{KL}^{\rho}(P_B(\tau|x; \phi), \tilde{P}_F(\tau|x)) + \nabla_{\phi} \mathbb{E}_{P_{B,\rho}(\tau; \phi)} [\log P_F^{\top}(x)] \\ &= \nabla_{\phi} D_{KL}^{\rho}(P_B(\tau|x; \phi), P_F(\tau|x)). \end{aligned} \quad (21)$$

$$= \nabla_{\phi} D_{KL}^{\rho}(P_B(\tau|x; \phi), P_F(\tau|x)). \quad (22)$$

Equations (20) and (22) are the expected results. \square

A.2. Proof of Proposition 3.3

Proof. The proof can be done by a procedure similar to that of P_B in Proposition 3.1 by replacing $\tilde{P}_F(\tau|x)$ with $P_G(\tau|x)$. \square

A.3. Sub-trajectory equivalence

Proposition 2 in the paper by Malkin et al. (2022b) only considered the gradients of the Sub-TB objective (Madan et al., 2023) w.r.t. $P_F(\cdot)$ and $P_B(\cdot)$. We provide an extended proposition below that also takes the gradients w.r.t. state flow estimator $F(\cdot)$ into consideration. For any $m < n$ and $n, m \in \{1, T-1\}$, we denote the set of sub-trajectories that start at some state in \mathcal{S}_m and end in some state in \mathcal{S}_n as $\tilde{\mathcal{T}} = \{\tilde{\tau} = (s_m \rightarrow \dots \rightarrow s_n) | \forall i \in \{m, \dots, n-1\} : (s_i \rightarrow s_{i+1}) \in \mathcal{A}_i\}$. The sub-trajectory objective $\mathcal{L}_{Sub-TB}(P_D) = \mathbb{E}_{P_D(\tilde{\tau})} [L_{Sub-TB}(\tilde{\tau})]$ is defined by:

$$L_{Sub-TB}(\tilde{\tau}) = \log \left(\frac{P_F(\tilde{\tau}|s_m) F(s_m)}{P_B(\tilde{\tau}|s_n) F(s_n)} \right)^2. \quad (23)$$

In the equations above, $P_F(\tilde{\tau}|s_m) = \prod_{t=m+1}^n P_F(s_t|s_{t-1})$, $P_B(\tilde{\tau}|s_n) = \prod_{t=m+1}^n P_B(s_{t-1}|s_t)$ and $F(s_n = x)$ is clamped to $R(x)$. Besides, we define $\mu(s_m) := F(s_m)/\hat{Z}_m$ and $\rho(s_n) := F(s_n)/\hat{Z}_n$ where \hat{Z}_m and \hat{Z}_n are the two normalizing constants whose values are clamped to $\sum_{s_m} F(s_m)$ and $\sum_{s_n} F(s_n)$. Furthermore, we define $P_{F,\mu}(\tilde{\tau}) := \mu_F(s_m) P_F(\tilde{\tau}|s_m)$ and $P_B(\tilde{\tau}) := \rho_B(s_n) P_B(\tilde{\tau}|s_n)$ so that $P_{F,\mu}(\tilde{\tau})/\rho^*(s_n) = P_{F,\mu}(\tilde{\tau}|s_n)$ and $P_{B,\rho}(\tilde{\tau}|s_m) = P_{B,\rho}(\tilde{\tau})/\mu^*(s_m)$, where $\rho^*(s_n) = F^*(s_n)/\hat{Z}_n^*$, $F^*(s_n) = \sum_{\tilde{\tau}: s_n \in \tilde{\tau}} F(s_m) P_F(\tilde{\tau}|s_m)$ is the ground-truth state flow over \mathcal{S}_n implied by P_F , $\hat{Z}_n^* = \sum_{s_n} F^*(s_n)$, $\mu^*(s_m) = F^*(s_m)/\hat{Z}_m^*$, $F^*(s_m) = \sum_{\tilde{\tau}: s_m \in \tilde{\tau}} F(s_n) P_B(\tilde{\tau}|s_n)$ is the ground-truth state flow over \mathcal{S}_m implied by P_B , and $\hat{Z}_m^* = \sum_{s_m} F^*(s_m)$.

Proposition A.2. For a forward policy $P_F(\cdot|\cdot; \theta)$, a backward distribution $P_B(\cdot|\cdot; \phi)$, state flow $F(\cdot; \theta)$, and state flow

$F(\cdot; \phi)^4$, the gradients of Sub-TB can be written as:

$$\begin{aligned} \frac{1}{2} \nabla_{\theta} \mathcal{L}_{Sub-TB}(P_{F,\mu}; \theta) &= \nabla_{\theta} D_{KL}^{\mu(\cdot; \theta)}(P_F(\bar{\tau}|s_m; \theta), P_{B,\rho}(\bar{\tau}|s_m)) + \nabla_{\theta} D_{KL}(\mu(s_m), \mu^*(s_m)) \\ &= \nabla_{\theta} D_{KL}^{\mu(\cdot; \theta)}(P_F(\bar{\tau}|s_m; \theta), \tilde{P}_{B,\rho}(\bar{\tau}|s_m)), \\ \frac{1}{2} \nabla_{\phi} \mathcal{L}_{Sub-TB}(P_{B,\rho}; \phi) &= \nabla_{\phi} D_{KL}^{\rho(\cdot; \phi)}(P_B(\bar{\tau}|s_n; \phi), P_{F,\mu}(\bar{\tau}|s_n)) + \nabla_{\phi} D_{KL}(\rho(s_m), \rho^*(s_m)) \\ &= \nabla_{\phi} D_{KL}^{\rho(\cdot; \phi)}(P_B(\bar{\tau}|s_n; \phi), \tilde{P}_{F,\mu}(\bar{\tau}|s_n)), \end{aligned} \quad (24)$$

where $\tilde{P}_{F,\mu}(\bar{\tau}|s_n) := P_{F,\mu}(\bar{\tau})/\rho(s_n)$ and $\tilde{P}_{B,\rho}(\bar{\tau}|s_m) := P_{B,\rho}(\bar{\tau})/\mu(s_m)$ are approximation to $P_{F,\mu}(\bar{\tau}|s_n)$ and $P_{B,\rho}(\bar{\tau}|s_m)$.

Proof. First of all, we split the parameters of the state flow estimator and forward transition probability and denote them as $F(\cdot; \theta_M)$ and $P_F(\cdot; \theta_F)$ respectively. We further define $c(\bar{\tau}) = \left(\log \frac{P_F(\bar{\tau}|s_m)}{F(s_n)P_B(\bar{\tau}|s_n)} \right)$.

For the gradients w.r.t. θ_F :

$$\begin{aligned} \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [\nabla_{\theta_F} L_{Sub-TB}(\bar{\tau}; \theta_F)] &= \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} \left[\nabla_{\theta_F} (c(\bar{\tau}; \theta_F) + \log F(s_m))^2 \right] \\ &= \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} \left[(c(\bar{\tau}; \theta_F) + \log F(s_m)) \nabla_{\theta_F} \log P_F(\bar{\tau}|s_m; \theta_F) \right] + \underbrace{\mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [\nabla_{\theta_F} (c(\bar{\tau}; \theta_F) + \log F(s_m))]}_{=0 \text{ by Lemma A.1}} \\ &= \mathbb{E}_{\mu(s_m)} [\nabla_{\theta_F} D_{KL}(P_F(\bar{\tau}|s_m; \theta_F), \tilde{P}_{B,\rho}(\bar{\tau}|s_m))] \\ &= \nabla_{\theta_F} D_{KL}^{\mu}(P_F(\bar{\tau}|s_m; \theta_F), \tilde{P}_{B,\rho}(\bar{\tau}|s_m)). \end{aligned} \quad (25)$$

Besides,

$$\begin{aligned} \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [\nabla_{\theta_F} L_{Sub-TB}(\bar{\tau}; \theta_F)] &= \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} \left[\nabla_{\theta_F} (c(\bar{\tau}; \theta_F) + \log F(s_m))^2 \right] \\ &= \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [c(\bar{\tau}; \theta_F) \nabla_{\theta_F} \log P_F(\bar{\tau}|s_m; \theta_F)] + \underbrace{\mathbb{E}_{\mu(s_m)} [\log F(s_m) \mathbb{E}_{P_F(\bar{\tau}|s_m; \theta_F)} [\nabla_{\theta_F} \log P_F(\bar{\tau}|s_m; \theta_F)]]}_{=0 \text{ by Lemma A.1}} \\ &= \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [c(\bar{\tau}; \theta_F) \nabla_{\theta_F} \log P_F(\bar{\tau}|s_m; \theta_F)] + \underbrace{\mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [\nabla_{\theta_F} c(\bar{\tau}; \theta_F)]}_{=0} \\ &= \nabla_{\theta_F} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [c(\bar{\tau}; \theta_F)] + \underbrace{\mathbb{E}_{\mu(s_m)} [\nabla_{\theta_F} \mathbb{E}_{P_F(\bar{\tau}|s_m; \theta_F)} [\log \mu^*(s_m)]]}_{=\log \mu^*(s_m) \nabla_{\theta_F} 1=0} + \nabla_{\theta_F} \log \hat{Z}_n \\ &= \nabla_{\theta_F} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [c(\bar{\tau}; \theta_F)] + \nabla_{\theta_F} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_F)} [\log \mu^*(s_m) + \log \hat{Z}_n] \\ &= \nabla_{\theta_F} D_{KL}^{\mu}(P_F(\bar{\tau}|s_m; \theta_F), P_{B,\rho}(\bar{\tau}|s_m)). \end{aligned} \quad (26)$$

For the gradients w.r.t. θ_M , we have:

$$\begin{aligned} \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} [\nabla_{\theta_M} L_{Sub-TB}(\bar{\tau}; \theta_M)] &= \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} \left[\nabla_{\theta_M} (c(\bar{\tau}) + \log F(s_m; \theta_M))^2 \right] \\ &= \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} [(c(\bar{\tau}) + \log F(s_m; \theta_M)) \nabla_{\theta_F} \log F(s_m; \theta_M)] \\ &= \mathbb{E}_{\mu(s_m; \theta_M)} \left[\nabla_{\theta_F} \log \frac{F(s_m; \theta_M)}{\hat{Z}_m} D_{KL}(P_F(\bar{\tau}|s_m; \theta_M), \tilde{P}_B(\bar{\tau}|s_m)) \right] \\ &= \nabla_{\theta_M} D_{KL}^{\mu(\cdot; \theta_M)}(P_F(\bar{\tau}|s_m), \tilde{P}_B(\bar{\tau}|s_m)). \end{aligned} \quad (27)$$

⁴Here F_{θ} and F_{ϕ} actually share the same parameters and represent the same flow estimator F . Model parameters are duplicated just for the clarity of gradient equivalences. Therefore the true gradient of the state flow estimator F is $\nabla_{\theta} F_{\theta} + \nabla_{\phi} F_{\phi}$.

Besides,

$$\begin{aligned}
 & \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} [\nabla_{\theta_M} L_{sub-TB}(\bar{\tau}; \theta_M)] = \frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} \left[\nabla_{\theta_M} (c(\bar{\tau}) + \log F(s_m; \theta_M))^2 \right] \\
 & = \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} [c(\bar{\tau}) \nabla_{\theta_F} \log F(s_m; \theta_M)] + \mathbb{E}_{\mu(s_m; \theta_M)} [\log F(s_m; \theta_M) \nabla_{\theta_M} \log F(s_m; \theta_M)] \\
 & = \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} \left[(c(\bar{\tau}) + \log \mu^*(s_m) + \log \hat{Z}_n) \nabla_{\theta_M} \log F(s_m; \theta_M) \right] \\
 & \quad + \mathbb{E}_{\mu(s_m; \theta_M)} \left[(\log F(s_m; \theta_M) - \log \mu^*(s_m) - \log \hat{Z}_n) \nabla_{\theta_M} \log F(s_m; \theta_M) \right], \tag{28}
 \end{aligned}$$

$$\begin{aligned}
 & = \mathbb{E}_{\mu(s_m; \theta_M)} \left[\nabla_{\theta_M} \log \frac{F(s_m; \theta_M)}{\hat{Z}_n} D_{KL}(P_F(\bar{\tau}|s_m), P_{B,\rho}(\bar{\tau}|s_m)) \right] \\
 & \quad + \mathbb{E}_{\mu(s_m; \theta_M)} \left[\left(\log \frac{F(s_m; \theta_M)}{\hat{Z}_n} - \log \mu^*(s_m) \right) \nabla_{\theta_M} \log \frac{F(s_m; \theta_M)}{\hat{Z}_n} \right] \\
 & = \nabla_{\theta_M} D_{KL}^{\mu(\cdot; \theta_M)}(P_F(\bar{\tau}|s_m), P_{B,\rho}(\bar{\tau}|s_m)) + \nabla_{\theta_M} D_{KL}(\mu(s_m; \theta_M), \mu^*(s_m)), \tag{29}
 \end{aligned}$$

where $\mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} [\log \hat{Z}_n \nabla_{\theta_M} \log F(s_m; \theta_M)] = \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta_M)} [\log \hat{Z}_n \nabla_{\theta_M} \log F(s_m; \theta_M)] = 0$ in (28) by Lemma A.1. Combining equations (25) and (27), we obtain

$$\frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta)} [\nabla_{\theta} L_{sub-TB}(\bar{\tau}; \theta)] = \nabla_{\theta} D_{KL}^{\mu(\cdot; \theta)}(P_F(\bar{\tau}|s_m; \theta), \tilde{P}_B(\bar{\tau}|s_m)). \tag{30}$$

Combining equations (26) and (29), we obtain

$$\frac{1}{2} \mathbb{E}_{P_{F,\mu}(\bar{\tau}; \theta)} [\nabla_{\theta} L_{sub-TB}(\bar{\tau}; \theta)] = \nabla_{\theta} \{ D_{KL}^{\mu(\cdot; \theta)}(P_F(\bar{\tau}|s_m; \theta), P_B(\bar{\tau}|s_m)) + D_{KL}(\mu(s_m; \theta), \mu^*(s_m)) \}. \tag{31}$$

Splitting ϕ into ϕ_B and ϕ_M and denoting $c(\bar{\tau}) = \log \frac{P_B(\bar{\tau}|s_m)}{P_F(\bar{\tau}|s_m)}$, the gradient derivation of ϕ follows the similar way as θ , and is omitted here. \square

B. RL framework

B.1. Derivation of RL functions

Let's first consider the case of forward policies. For any $s \in \mathcal{S}_t$ and $a = (s \rightarrow s') \in \mathcal{A}(s)$ with $t \in \{0, \dots, T-1\}$, we define the $V_{F,t}$ and $Q_{F,t}$ as:

$$\begin{aligned}
 V_{F,t}(s) &:= \mathbb{E}_{P_F(\tau_{>t}|s_t)} \left[\sum_{l=t}^{T-1} R_F(s_l, a_l) \middle| s_t = s \right] \\
 &= R_F(s) + \mathbb{E}_{P_F(s_{t+1}|s_t)} \left[\mathbb{E}_{P_F(\tau_{>t+1}|s_{t+1})} \left[\sum_{l=t+1}^T R_F(s_l, a_l) \middle| s_{t+1} = s' \right] \middle| s_t = s \right] \\
 &= R_F(s) + \mathbb{E}_{\pi_F(s,a)} [V_{F,t+1}(s')] \\
 Q_{F,t}(s, a) &:= \mathbb{E}_{P_F(\tau_{>t+1}|s_t, a_t)} \left[\sum_{l=t}^{T-1} R_F(s_l, a_l) \middle| s_t = s, a_t = a \right] \\
 &= R_F(s, a) + \mathbb{E}_{P_F(\tau_{>t+1}|s_{t+1})} \left[\sum_{l=t+1}^T R_F(s_l, a_l) \middle| s_{t+1} = s' \right] \\
 &= R_F(s, a) + V_{F,t+1}(s')
 \end{aligned} \tag{32}$$

where $R_F(s) := \mathbb{E}_{\pi_F(s,a)} [R_F(s, a)]$, $V_{F,T}(\cdot) := 0$, and $Q_{F,T}(\cdot, \cdot) := 0$. Since $S_t \cap S_{t'} = \emptyset$ for any $t \neq t'$, we can read off the time indices (topological orders) from state values. Plus the fact that $R_F(s, a) := 0$ for any $a \notin \mathcal{A}(s)$, we are allowed to define two universal $V_F : \mathcal{S} \rightarrow \mathbb{R}$ and $Q_F : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ such that $V_F(s_t = s) := V_{F,t}(s)$ and $Q_F(s_t = s, a) := Q_{F,t}(s, a)$. *Remark B.1.* While the transition environment \mathcal{G} is exactly known, the state space \mathcal{S} can be exponentially large with the exact values of V and Q to be intractable. This fact, in spirit, corresponds to a regular RL problem where the exact values of V and Q are infeasible due to the unknown and uncertain transition model $P(s'|s, a)$.

For backward policies, rewards are accumulated from time T to 1. Similarly, for $s' \in \mathcal{S}_t$ and $a = (s \rightarrow s') \in \dot{\mathcal{A}}(s')$ we can define:

$$\begin{aligned} V_{B,t}(s') &:= \mathbb{E}_{P_B(\tau < t | s_t)} \left[\sum_{l=1}^t R_B(s_l, a_l) \middle| s_t = s' \right] = R_B(s') + \mathbb{E}_{\pi_B(s', a)} [V_{B,t-1}(s)] \\ Q_{B,t}(s', a) &:= \mathbb{E}_{P_B(\tau < t-1 | s_t, a_t)} \left[\sum_{l=1}^t R_B(s_l, a_l) \middle| s_t = s', a_t = a \right] = R_B(s', a) + V_{B,t-1}(s) \end{aligned} \quad (33)$$

where $R_B(s') = \mathbb{E}_{\pi_B(s', a)} [R_B(s', a)]$, $V_{B,0}(\cdot) := 0$, and $Q_{B,0}(\cdot, \cdot) := 0$. For the same reason as forward policies, we can define universal functions $V_B : \mathcal{S} \rightarrow R$ and $Q_B : \mathcal{S} \times \dot{\mathcal{A}} \rightarrow R$ such that $V_B(s_t = s') := V_{B,t}(s')$ and $Q_B(s_t = s', a) := Q_{B,t}(s', a)$.

Based on the definitions above, the expected value functions are defined as:

$$J_F := \mathbb{E}_{\mu(s_0)} [V_F(s_0)], \quad J_B := \mathbb{E}_{\rho(x)} [V_B(x)]. \quad (34)$$

By definitions, $V_F(s_0) = \mathbb{E}_{P_F(\tau > t | s_0)} [\sum_{l=0}^{T-1} R_F(s_l, a_l) | s_0] = D_{KL}(P_F(\tau | s_0), \tilde{P}_B(\tau | s_0))$, so $J_F = D_{KL}^\mu(P_F(\tau | s_0), \tilde{P}_B(\tau | s_0))$. Likewise, we can obtain $J_B = D_{KL}^\rho(P_B(\tau | x), \tilde{P}_F(\tau | x))$. The advantages functions are defined as:

$$A_F(s, a) = Q_F(s, a) - V_F(s), \quad A_B(s', a) = Q_B(s', a) - V_B(s'). \quad (35)$$

We define forward accumulated state distribution as $d_{F,\mu}(s) := \frac{1}{T} \sum_{t=0}^{T-1} P_{F,\mu}(s_t = s)$ such that for arbitrary function $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$,

$$\begin{aligned} \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} f(s_t, a_t) \right] &= \sum_{t=0}^T \mathbb{E}_{P_{F,\mu}(s_t, s_{t+1})} [f(s_t, a_t)] = \sum_{t=0}^T \mathbb{E}_{P_{F,\mu}(s_t), \pi(s_t, a_t)} [f(s_t, a_t)] \\ &= \sum_t^T \sum_s P_{F,\mu}(s_t = s) \sum_a^{\mathcal{A}} \pi(s, a) f(s, a) = \sum_s^{\mathcal{S}} \sum_a^{\mathcal{A}} \left(\sum_t^T P_{F,\mu}(s_t = s) \right) \pi(s, a) f(s, a) \\ &= T \cdot \sum_s^{\mathcal{S}} \sum_a^{\mathcal{A}} d_{F,\mu}(s) \pi(s, a) f(s, a) = T \cdot \mathbb{E}_{d_{F,\mu}(s), \pi_F(s, a)} [f(s, a)] \end{aligned} \quad (36)$$

where equation (36) holds in that $\forall s \notin \mathcal{S}_t : P(s_t = s) = 0$ and $\forall a \notin \mathcal{A}(s) : \pi(s, a) = 0$. By the fact that $S_t \cap S_{t'} = \emptyset$ for any $t \neq t'$ and any $\tau \in \mathcal{T}$ must pass some $s_t \in \mathcal{S}_t$ with $t \in \{0, \dots, T-1\}$, $P_{F,\mu}(s_t)$ is a valid distribution over \mathcal{S}_t and $\sum_{s_t} P_{F,\mu}(s_t) = 1$. Accordingly, $d_{F,\mu}(s)$ is a valid distribution over \mathcal{S} and $T \cdot d_{F,\mu}(s) = P_{F,\mu}(s_t)$. Analogically, we can define $d_{B,\rho}(s') := \frac{1}{T} \sum_{t=1}^T P_{B,\rho}(s_t = s')$ such that for arbitrary function $f : \mathcal{S} \times \dot{\mathcal{A}} \rightarrow \mathbb{R}$,

$$\mathbb{E}_{P_{B,\rho}(\tau)} \left[\sum_{t=1}^T f(s_t, a_t) \right] = T \cdot \mathbb{E}_{d_{B,\rho}(s'), \pi_B(s', a)} [f(s', a)] \quad (37)$$

B.2. DAGs as transition environments

Theorem B.2. (Golpar Raboky & Eftekhari, 2019): Let $P \in \mathbb{R}^{N \times N}$ be a non-negative matrix, the following statements are equivalent:

1. P is nilpotent;
2. $P^N = 0$;
3. The directed graph $\mathcal{G}(\mathcal{S}, \mathcal{A})$ associated with P is a DAG graph;
4. There exists a permutation matrix U such that $U^T P U$ is a strictly triangular matrix.

where $\mathcal{S} = \{s^0, \dots, s^{N-1}\}$ and $\mathcal{A} = \{(s^i \rightarrow s^j) | P_{i,j} \neq 0\}$ are node and edge sets.

Lemma B.3. : For any DAG graph $\mathcal{G}(\mathcal{S}, \mathcal{A})$ associated with $P \in \mathbb{R}^{N \times N}$ with $T+1 (\leq N)$ different topological node orders indexed by integers $[0, T]$,

$$\forall t > T, \quad P^t = 0 \quad (38)$$

Proof. We prove the result by contradiction. Assuming $P^t (t > T)$ is not zero, then $\forall i \neq j$:

$$(P^t)_{i,j} = \sum_{k_1:t-1} P_{i,k_1} P_{k_1,k_2} \dots P_{k_{t-1},j} \quad (39)$$

By the nature of DAG graphs, $\forall (s' \rightarrow s) \in \mathcal{A} : s' \prec s$. Then the above expression is equal to:

$$\begin{aligned} &= \sum_{k_1:s^i \prec s^{k_1}} P_{i,k_1} \left(\sum_{k_2:s^{k_1} \prec s^{k_2}} P_{k_1,k_2} \dots \left(\sum_{k_{t-1}:s^{k_{t-2}} \prec s^{k_{t-1}}} P_{k_{t-2},k_{t-1}} P_{k_{t-1},j} \right) \right) \\ &= \sum_{k_1:t-1:(s^i \prec s^{k_1} \prec \dots \prec s^{k_{t-1}} \prec s^j)} P_{i,k_1} P_{k_1,k_2} \dots P_{k_{t-1},j} \\ &> 0 \end{aligned} \quad (40)$$

Then it means that there at least exists a trajectory $(s^i \prec s^{k_1} \prec \dots \prec s^{k_{t-1}} \prec s^j)$ with non-zero probability. However, there are at least $t+1$ distinct topological orders in the path, which contradicts the assumption that there are $T+1$ different node orders. \square

Let's return to the graded DAG, $\mathcal{G}(\mathcal{S}, \mathcal{A})$ in GFlowNets. For the easiness of analysis, we restrict forward and backward policies and initial distribution to be tabular forms, $P_F \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, $\mu \in \mathbb{R}^{|\mathcal{S}|}$, $P_B \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, and $\rho \in \mathbb{R}^{|\mathcal{S}|}$ such that $P_F(s^j | s^i) = (P_F)_{j,i}$ and $P_B(s^j | s^i) = (P_B)_{j,i}$. Besides, we split initial distribution vectors by $\mu = [\bar{\mu}; 0] \in \mathbb{R}^{|\mathcal{S}|}$ and $\rho = [0; \bar{\rho}] \in \mathbb{R}^{|\mathcal{S}|}$, where $\bar{\mu}$ and $\bar{\rho}$ denote the probabilities of states except s^f and s^0 respectively. We denote the graph equipped with self-loop over s^f as $\mathcal{G}_F(\mathcal{S}, \mathcal{A} \cup \{(s^f \rightarrow s^f)\})$, and the reverse graph equipped with self-loop over s^0 as $\mathcal{G}_B(\mathcal{S}, \mathcal{A} \cup \{(s^0 \rightarrow s^0)\})$. Accordingly, we enhanced P_F and P_B by defining $P_F(s^f | s^f) := 1$ and $P_B(s^0 | s^0) := 1$. (\mathcal{G}_F, P_F) specifies an absorbing Markov Chains: s^f is the only absorbing state as only self-loop is allowed once entering s^f ; the sub-graph over $\mathcal{S} \setminus \{s^f\}$ is still a DAG, so any state $s \in \mathcal{S} \setminus \{s^f\}$ is transient as it can be visited at least one time. Similarly, (\mathcal{G}_B, P_B) specifies another absorbing Markov Chains with absorbing state s^0 . For graph \mathcal{G}_F and \mathcal{G}_B , their transition matrices P_F and P_B can be decomposed into:

$$P_F = \begin{pmatrix} \bar{P}_F & 0 \\ r_F & 1 \end{pmatrix}, P_B = \begin{pmatrix} 1 & r_B \\ 0 & \bar{P}_B \end{pmatrix} \quad (41)$$

Where $r_F^\top \in \mathbb{R}^{|\mathcal{S}|-1}$ and $r_B^\top \in \mathbb{R}^{|\mathcal{S}|-1}$ denotes the forward probability of $(s \rightarrow s^f)$ for any $s \in \mathcal{S} \setminus \{s^f\}$ and $(s^0 \leftarrow s)$, for any $s \in \mathcal{S} \setminus \{s^0\}$, $\bar{P}_F \in \mathbb{R}^{(|\mathcal{S}|-1) \times (|\mathcal{S}|-1)}$ and $\bar{P}_B \in \mathbb{R}^{(|\mathcal{S}|-1) \times (|\mathcal{S}|-1)}$ denote probability of $(s \rightarrow s')$ for any $s, s' \in \mathcal{S} \setminus \{s^f\}$ and $(s \leftarrow s')$ for any $s, s' \in \mathcal{S} \setminus \{s^0\}$.

Lemma B.4. : For (\mathcal{G}, P_F, μ) and $(\mathcal{G}_B, P_B, \rho)$, $d_{F,\mu} \in \mathbb{R}^{|\mathcal{S}|}$ and $d_{B,\rho} \in \mathbb{R}^{|\mathcal{S}|}$, can be written in the following form:

$$d_{F,\mu} = (\bar{d}_{F,\mu}, 0), \bar{d}_{F,\mu} = \frac{1}{T} (I - \bar{P}_F)^{-1} \bar{\mu} \quad (42)$$

$$d_{B,\rho} = (0, \bar{d}_{B,\rho}), \bar{d}_{B,\rho} = \frac{1}{T} (I - \bar{P}_B)^{-1} \bar{\rho} \quad (43)$$

Proof. We prove the result for the forward case. A similar proof procedure can easily derived for the backward case, so it is omitted. First of all, by the nature of Markov Chains, $P_{F,\mu}(s_t = s^i) = [(P_F)^t \mu]_i$, and $d_{F,\mu} = \frac{1}{T} \sum_{t=0}^T (P_F)^t \mu$. Then, it can be easily verified (Grinstead & Snell, 2006):

$$(P_F)^t = \begin{pmatrix} (\bar{P}_F)^t & 0 \\ * & 1 \end{pmatrix}, \quad (44)$$

where the explicit expression of the upper right corner is omitted. By theorem B.2, \bar{P}_F is a nilpotent matrix and by lemma B.3, $\frac{1}{T} \sum_{t=0}^T (\bar{P}_F)^t = \frac{1}{T} \sum_{t=0}^{\infty} (\bar{P}_F)^t = \frac{1}{T} (I - \bar{P}_F)^{-1}$ (The second equality follows that fact that: $(I - \bar{P}_F) \sum_{t=0}^{\infty} (\bar{P}_F)^t = \sum_{t=0}^{\infty} (\bar{P}_F)^t - \sum_{t=1}^{\infty} (\bar{P}_F)^t = I$). Therefore,

$$d_{F,\mu} = \frac{1}{T} \begin{pmatrix} \sum_{t=0}^T \bar{P}_F^t & 0 \\ * & 1 \end{pmatrix} \mu = \left(\frac{1}{T} (I - \bar{P}_F)^{-1} \bar{\mu}, * \bar{\mu} \right)$$

By Theorem 11.4 in Grinstead & Snell (2006), $(I - \bar{P}_F)_{i,j}^{-1}$ is the expected number of times the chain is in state s_j starting from s_i before absorbing in s^f . And $[(I - \bar{P}_F)^{-1} \bar{\mu}]_j$ is the expected number of times the chain is in state s_j before absorbing. Since $\forall s \notin \mathcal{S}_0 : \mu(s) = 0$ and the fact \mathcal{G} is graded, any trajectories over \mathcal{G} must start from states in \mathcal{S}_0 to state in \mathcal{S}_T i.e. $\sum_j [(I - \bar{P}_F)^{-1} \bar{\mu}]_j = T$. Thus, $\frac{1}{T} [(I - \bar{P}_F)^{-1} \bar{\mu}]_i$ denote the fraction of staying in transient state s_i before absorbing, this is, the probability observing state s_i within T time steps. By the same reasoning, we can conclude that $*\bar{\mu} = 0$ as s^f can not be reached within $T - 1$ transitions. \square

Lemma B.5. : For two forward policy, π_F and π'_F , and two backward policy, π_B and π'_B , we have:

$$\begin{aligned} D_{TV}(d_{F,\mu}(\cdot), d'_{F,\mu}(\cdot)) &\leq D_{TV}^{d'_{F,\mu}}(\pi_F(s, \cdot), \pi'_F(s, \cdot)), \\ D_{TV}(d_{B,\rho}(\cdot), d'_{B,\rho}(\cdot)) &\leq D_{TV}^{d'_{B,\rho}}(\pi_B(s, \cdot), \pi'_B(s, \cdot)), \end{aligned} \quad (45)$$

where for three arbitrary distributions p, q and u , $D_{TV}(p(\cdot), q(\cdot)) := \frac{1}{2} \|p(\cdot) - q(\cdot)\|_1$ and $D_{TV}^u(p(\cdot|s), q(\cdot|s)) := \frac{1}{2} \mathbb{E}_{u(s)} \|p(\cdot|s) - q(\cdot|s)\|_1$.

Proof. The proof procedure follows that of Lemma 3 in Achiam et al. (2017). For two forward policy π_F and π'_F , denoting $\bar{N}_F := (I - \bar{P}_F)^{-1}$ and $\bar{N}'_F := (I - \bar{P}'_F)^{-1}$, then

$$\bar{N}_F^{-1} - (\bar{N}'_F)^{-1} = \bar{P}'_F - \bar{P}_F := \Delta \quad (46)$$

and

$$\bar{N}'_F - \bar{N}_F = \bar{N}'_F \Delta \bar{N}_F \quad (47)$$

Then,

$$\begin{aligned} \|d_{F,\mu} - d'_{F,\mu}\|_1 &= \|\bar{d}_{F,\mu} - \bar{d}'_{F,\mu}\|_1 \\ &= \frac{1}{T} \|(\bar{N}_F - \bar{N}'_F) \bar{\mu}\|_1 = \frac{1}{T} \|\bar{N}_F \Delta \bar{d}'_{F,\mu}\|_1 \\ &\leq \frac{1}{T} \|\bar{N}_F\|_1 \|\Delta \bar{d}'_{F,\mu}\|_1 \leq \|\Delta \bar{d}'_{F,\mu}\|_1 \end{aligned} \quad (48)$$

Therefore, we have

$$\begin{aligned} \|\Delta \bar{d}'_{F,\mu}\|_1 &= \|(P'_F - P_F) d'_{F,\mu}\|_1 = \sum_s \left| \sum_{s'} (P'_F(s'|s) - P_F(s'|s)) d'_{F,\mu}(s) \right| \\ &\leq \sum_{s,s'} |P'_F(s'|s) - P_F(s'|s)| d'_{F,\mu}(s) = \sum_{s,a} |\pi'(s, a) - \pi(s, a)| d'_{F,\mu}(s) \\ &= \mathbb{E}_{d'_{F,\mu}(s)} [\|\pi'_F(s, \cdot) - \pi_F(s, \cdot)\|_1] \end{aligned}$$

Where first equality holds by the fact that $P_F(s^f|s^f) = P'_F(s^f|s^f) = 1$ for any two forward policies consistent with \mathcal{G}_F . The result of backward policies can be derived analogically and is omitted here. \square

B.3. Derivation of gradients

Proposition B.6. The gradients of $J_F(\theta)$ and J_B^ϕ w.r.t θ and ϕ can be written as:

$$\begin{aligned} \nabla_\theta J_F(\theta) &= T \mathbb{E}_{d_{F,\mu}(s) \pi_F(s,a)} [Q_F(s, a) \nabla_\theta \log \pi_F(s, a; \theta)] + \mathbb{E}_{\mu(s_0)} [V_F(s_0) \nabla_\theta \log \mu(s_0; \theta)] \\ &= T \mathbb{E}_{d_{F,\mu}(s) \pi_F(s,a)} [A_F(s, a) \nabla_\theta \log \pi_F(s, a; \theta)] + \mathbb{E}_{\mu(s_0)} [V_F(s_0) \nabla_\theta \log \mu(s_0; \theta)] \\ \nabla_\phi J_B(\phi) &= T \mathbb{E}_{d_{B,\rho}(s) \pi_B(s,a)} [Q_B(s, a) \nabla_\phi \log \pi_B(s, a; \phi)] \\ &= T \mathbb{E}_{d_{B,\rho}(s) \pi_B(s,a)} [A_B(s, a) \nabla_\phi \log \pi_B(s, a; \phi)] \end{aligned} \quad (49)$$

Remark B.7. This result implies that an estimated value function, which may differ from the exact one, does not lead to biased gradient estimation.

Proof.

$$\begin{aligned}
 \nabla_{\theta} J_F(\theta) &= \nabla_{\theta} \frac{Z^{\theta}}{\bar{Z}} V_F(s_0; \theta) = V_F(s_0) \nabla_{\theta} \frac{Z(\theta)}{\bar{Z}} + \frac{Z}{\bar{Z}} \nabla_{\theta} V_F(s_0; \theta) \\
 &= \mathbb{E}_{\mu(s_0)} [V_F(s_0) \nabla_{\theta} \log \mu(s_0; \theta)] + \underbrace{\mathbb{E}_{\mu(s_0)} [\nabla_{\theta} V_F(s_0; \theta)]}_{(1)} \\
 &\stackrel{(1)}{=} \mathbb{E}_{P_{F,\mu}(s_0)} [\nabla_{\theta} \mathbb{E}_{\pi_F(s_0, a_0; \theta)} [Q_F(s_0, a_0; \theta)]] \\
 &= \mathbb{E}_{P_{F,\mu}(s_0, a_0)} [\nabla_{\theta} \log \pi_F(s_0, a_0; \theta) Q_F(s_0, a_0; \theta) + \nabla_{\theta} Q_F(s_0, a_0; \theta)] \\
 &= \mathbb{E}_{P_{F,\mu}(s_0, a_0)} [\nabla_{\theta} \log \pi_F(s_0, a_0; \theta) Q_F(s_0, a_0; \theta)] + \underbrace{\mathbb{E}_{P_{F,\mu}(s_0, a_0)} [\nabla_{\theta} R_F(s_0, a_0; \theta) + V_F(s_1; \theta)]}_{(2)} \\
 &\stackrel{(2)}{=} \underbrace{\mathbb{E}_{P_{F,\mu}(s_0, a_0)} \left[\nabla_{\theta} \log \frac{\pi_F(s_0, a_0; \theta)}{\pi_B(s_1, a_0)} \right]}_{(3)} + \mathbb{E}_{P_{F,\mu}(s_1)} [\nabla_{\theta} V_F(s_1; \theta)] \\
 &\stackrel{(3)}{=} \mathbb{E}_{P_{F,\mu}(s_0)} [\mathbb{E}_{\pi_F(s_0, a_0)} [\nabla_{\theta} \log \pi_F(s_0, a_0; \theta)]] \\
 &= \underbrace{\mathbb{E}_{P_{F,\mu}(s_0)} [\nabla_{\theta} 1]}_{\text{By Lemma A.1}} = 0
 \end{aligned} \tag{50}$$

Therefore,

$$\mathbb{E}_{P_{F,\mu}(s_0)} [\nabla_{\theta} V_F(s_0; \theta)] = \mathbb{E}_{P_{F,\mu}(s_0, a_0)} [\nabla_{\theta} \log \pi_F(s_0, a_0; \theta) Q_F(s_0, a_0; \theta)] + \mathbb{E}_{P_{F,\mu}(s_1)} [\nabla_{\theta} V_F(s_1; \theta)] \tag{51}$$

Keep doing the process, we have

$$\mathbb{E}_{P_{F,\mu}(s_t)} [\nabla_{\theta} V_F(s_t; \theta)] = \mathbb{E}_{P_{F,\mu}(s_t, a_t)} [\nabla_{\theta} \log \pi_F(s_t, a_t; \theta) Q_F(s_t, a_t; \theta)] + \mathbb{E}_{P_{F,\mu}(s_{t+1})} [\nabla_{\theta} \underbrace{V_F(s_{t+1}; \theta)}_{V(s_T)=0}] \tag{52}$$

Then,

$$\stackrel{(1)}{=} \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_F(s_t, a_t; \theta) Q_F(s_t, a_t; \theta) \right] = \mathbb{E}_{d_{F,\mu}(s) \pi_F(s, a)} [\nabla_{\theta} \log \pi_F(s, a; \theta) Q_F(s, a; \theta)] \tag{53}$$

Besides,

$$\begin{aligned}
 &\stackrel{(1)}{=} \mathbb{E}_{d_{F,\mu}(s) \pi_F(s, a)} [\nabla_{\theta} \log \pi_F(s, a; \theta) Q_F(s, a; \theta)] - \mathbb{E}_{d_{F,\mu}(s)} [V_F(s) \underbrace{\mathbb{E}_{\pi_F(s, a)} [\nabla_{\theta} \log \pi_F(s, a; \theta)]}_{=0}] \\
 &= \mathbb{E}_{d_{F,\mu}(s) \pi_F(s, a)} [\nabla_{\theta} \log \pi_F(s, a; \theta) A_F(s, a; \theta)]
 \end{aligned} \tag{54}$$

The derivation of $\nabla_{\phi} J_B(\phi)$ follows the similar way as $\nabla_{\theta} J_F(\theta)$, and is omitted here. \square

B.4. Connection of policy-based training to TB-based training

The gradient of TB objective w.r.t. θ_F can be written as:

$$\frac{1}{2} \nabla_{\theta_F} \mathcal{L}_{TB}(P_{F,\mu}; \theta_F) = \sum_{t=1}^T \mathbb{E}_{P_{F,\mu}(\tau)} \left[\nabla_{\theta_F} \log P_F(s_t | s_{t-1}; \theta_F) \left(\sum_{t=1}^T \log \frac{P_F(s_t | s_{t-1})}{P_B(s_{t-1} | s_t)} \right) \right]. \tag{55}$$

In equation (55), each term for $t > 0$ can be expanded as:

$$\begin{aligned}
 & \mathbb{E}_{P_{F,\mu}(\tau \geq t-1)} \left[\nabla_{\theta_F} \log P_F(s_t | s_{t-1}; \theta_F) \left(\sum_{l=t}^T \log \frac{P_F(s_l | s_{l-1})}{P_B(s_{l-1} | s_l)} \right) \right] \\
 & + \mathbb{E}_{P_{F,\mu}(\tau \leq t)} \left[\nabla_{\theta_F} \log P_F(s_t | s_{t-1}; \theta_F) \left(\sum_{l=1}^{t-1} \log \frac{P_F(s_l | s_{l-1})}{P_B(s_{l-1} | s_l)} \right) \right] \\
 & = \mathbb{E}_{P_{F,\mu}(\tau \geq t-1)} \left[\nabla_{\theta_F} \log P_F(s_t | s_{t-1}; \theta_F) \left(\sum_{l=t}^T \log \frac{P_F(s_l | s_{l-1})}{P_B(s_{l-1} | s_l)} \right) \right] \\
 & + \mathbb{E}_{P_{F,\mu}(\tau \leq t-1)} \left[\left(\sum_{l=1}^{t-1} \log \frac{P_F(s_l | s_{l-1})}{P_B(s_{l-1} | s_l)} \right) \underbrace{\mathbb{E}_{P_{F,\mu}(s_t | s_{t-1})} [\nabla_{\theta_F} \log P_F(s_t | s_{t-1}; \theta_F)]}_{=0 \text{ by Lemma A.1}} \right] \tag{56}
 \end{aligned}$$

Thus,

$$\begin{aligned}
 \frac{1}{2} \nabla_{\theta_F} \mathcal{L}_{TB}(P_{F,\mu}; \theta_F) &= \sum_{t=1}^T \mathbb{E}_{P_{F,\mu}(\tau)} \left[\nabla_{\theta_F} \log P_F(s_t | s_{t-1}; \theta_F) \left(\sum_{l=t}^T \log \frac{P_F(s_l | s_{l-1})}{P_B(s_{l-1} | s_l)} \right) \right] \\
 &\quad - C \sum_{t=0}^{T-1} \underbrace{\mathbb{E}_{P_{F,\mu}(\tau)} [\nabla_{\theta} \log P_F(s_t, a_t; \theta)]}_{=0 \text{ by Lemma A.1}} \\
 &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta_F} \log P_F(s_t | s_{t-1}; \theta_F) \left(\sum_{l=t}^T \log \frac{P_F(s_l | s_{l-1})}{P_B(s_{l-1} | s_l)} - C \right) \right], \tag{57}
 \end{aligned}$$

where C is an added baseline and constant w.r.t. θ for variance reduction during gradient estimation. As shown in Appendix B.3, the gradient of J_F w.r.t. θ_F can be written as:

$$\begin{aligned}
 \nabla_{\theta_F} J_F(\theta_F) &= T \mathbb{E}_{d_{F,\mu}(s) \pi_F(s,a)} [A_F(s,a) \nabla_{\theta_F} \log \pi_F(s,a; \theta_F)] \\
 &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_F(s_t, a_t; \theta) Q_F(s_t, a_t) \right] \\
 &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_F(s_t, a_t; \theta) \mathbb{E}_{P_{F,\mu}(\tau > t+1 | s_t, s_{t+1})} \left[\sum_{l=t}^{T-1} R_F(s_l, a_l) \middle| s_t, a_t \right] \right] \\
 &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_F(s_t, a_t; \theta) \left(\sum_{l=t}^{T-1} R_F(s_l, a_l) \right) \right] \\
 &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_F(s_t, a_t; \theta) \left(\sum_{l=t}^{T-1} R_F(s_l, a_l) - C \right) \right] \tag{58}
 \end{aligned}$$

This result implies that: when we update the forward policy by the estimation of $\nabla_{\theta_F} \mathcal{L}_{TB}(\theta_F)$ based on a batch of sampled trajectories, we approximate $Q_F(s_t, a_t)$ empirically by $\hat{Q}_F(s_t, a_t) = \sum_{l=t}^{T-1} R_F(s_l, a_l)$ for each sample, and can further reduce the estimation variance by some unbiased constant baseline C . By comparison, the RL formulation generalizes the constant to an unbiased functional baseline $\tilde{V}_F(s; \eta)$, which is the approximation of exact $V_F(s)$. This enables to approximate $Q_F(s_t, a_t)$ and $A_F(s_t, a_t)$ functionally by $\hat{Q}_F(s_t, a_t) = R(s_t, a_t) + \tilde{V}_F(s_{t+1})$ and $\hat{A}_F = \hat{Q}_F(s_t, a_t) - \tilde{V}_F(s_t)$ which can further be generalized to $\sum_{l=0}^{T-1-t} \lambda^l (\hat{Q}_F(s_t, a_t) - \tilde{V}_F(s_t))$, allowing flexible bias-variance trade-off for gradient estimation (Schulman et al., 2016) (Appendix B.6).

B.5. Connection between policy-based training and Soft-Q learning

In the following text, we discuss the relationship between our policy-based method and Soft-Q learning (Haarnoja et al., 2018), one of the most representative MaxEnt RL methods.

First, we introduce their connection when the total estimator $\log Z$ is **fixed**. We can expand $-J_F$ as:

$$\begin{aligned} -J_F &= T\mathbb{E}_{d_{F,\mu}(s), \pi_F(s,a)} [\log \pi_B(s', a) - \log \pi_F(s, a)] \\ &= T\mathbb{E}_{d_{F,\mu}(s), \pi_F(s,a)} [\log \pi_B(s', a) + \mathbb{E}_{\pi_F(s,a)} [-\log \pi_F(s, a)]] \\ &= T\mathbb{E}_{d_{F,\mu}(s), \pi_F(s,a)} [\log \pi_B(s', a) + \mathcal{H}(\pi_F(s, \cdot))], \end{aligned} \quad (59)$$

where $a = (s \rightarrow s')$, and \mathcal{H} denotes the entropy of a distribution. It implies that fixing the total flow estimator $\log Z$, maximizing $-J_F$ w.r.t. π_F can be interpreted as a Maximum-Entropy (MaxEnt) RL problem, where $\log \pi_B(s', a)$ is the static reward w.r.t. $\pi_F(s, a)$. We define $Q_F^S(s, a) := -Q_F(s, a) + \log \pi_F(s, a)$ and $V_F^S(s, a) := -V_F(s, a)$, which implies that

$$Q_F^S(s, a) = \log \pi_B(s', a) + V_F^S(s'). \quad (60)$$

We use a parameterized function \tilde{Q}_F^S as the approximation to Q_F^S , and define $\tilde{V}_F^S(s) := \log \sum_a \exp\{\tilde{Q}_F^S(s, a)\}$ as the approximation to V_F^S . We parameterize π_F by $\pi_F(s, a; \theta) := \frac{\exp\{\tilde{Q}_F^S(s, a; \theta)\}}{\exp\{\tilde{V}_F^S(s; \theta)\}}$, which implies that:

$$\tilde{Q}_F^S(s, a) = \tilde{V}_F^S(s) + \log \pi_F(s, a). \quad (61)$$

In Soft-Q learning, π_F is updated by the gradient of the following objective:

$$\frac{T}{2} \mathbb{E}_{d_{\mathcal{D}}(s), \pi_{\mathcal{D}}(s,a)} \left[\left(\tilde{Q}_F^S(s, a; \theta) - \hat{Q}_F^S(s, a) \right)^2 \right], \quad (62)$$

where $\hat{Q}_F^S(s_t, a_t) = \log \pi_B(s_{t+1}, a_t) + \tilde{V}_F^S(s_{t+1})$. It has been shown that the policy gradients for static rewards plus the gradients of the policy entropy (or KL divergence from some reference policy) are equivalent to the gradients of the corresponding soft Q function (Schulman et al., 2017). Likewise, we demonstrate our policy-based method with $\lambda = 0$ is equivalent to Soft-Q learning with $\pi_D = \pi_F$ without any adaption. When choosing $\pi_D = \pi_F$, the gradients of equation (62) w.r.t. θ can be written as:

$$\begin{aligned} &\frac{T}{2} \nabla_{\theta} \mathbb{E}_{d_{F,\mu}(s) \pi_F(s,a)} \left[\left(\tilde{Q}_F^S(s, a; \theta) - (\log P_B(s', a) + \tilde{V}_F^S(s')) \right)^2 \right] \\ &= T\mathbb{E}_{d_{F,\mu}(s) \pi_F(s,a)} \left[\nabla_{\theta} \tilde{Q}_F^S(s, a; \theta) \left(\tilde{Q}_F^S(s, a; \theta) - \log P_B(s', a) - \tilde{V}_F^S(s') \right) \right] \\ &= T\mathbb{E}_{d_{F,\mu}(s) \pi_F(s,a)} \left[\nabla_{\theta} \left(\log \pi_F(s, a) + \tilde{V}_F^S(s) \right) \left(\tilde{V}_F^S(s) + R_F(s, a) - \tilde{V}_F^S(s') \right) \right] \\ &= T\mathbb{E}_{d_{F,\mu}(s) \pi_F(s,a)} \left[\nabla_{\theta} \left(\log \pi_F(s, a) + \tilde{V}_F^S(s) \right) \hat{\delta}_F(s, a) \right] \\ &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_F(s_t, a_t; \theta) \hat{\delta}_F(s_t, a_t) \right] + \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \tilde{V}_F^S(s_t; \theta) \hat{\delta}_F(s_t, a_t) \right], \end{aligned} \quad (63)$$

where $\hat{\delta}_F(s, a) := -\tilde{V}_F(s) + R_F(s, a) + \tilde{V}_F(s')$, and $\tilde{V}_F := -\tilde{V}_F^S$. Comparing to (66) with $\lambda = 0$, a clear equivalence can be established.

We further connect the Soft-Q learning objective (62) to the Flow Matching (FM) objective (Bengio et al., 2021) and explain the role that $\log Z$ plays during training. When \tilde{Q}_F^S achieves the optimal point, $\tilde{Q}_F^S = Q_F^S$. Consequently, $\tilde{Q}_F^S(x, (x \rightarrow s^f)) = \log \pi_B(x, (x \rightarrow s^f)) = \log \frac{F(x \rightarrow s^f)}{Z}$, $\tilde{V}_F^S(x) = \log \sum_a \exp\{Q_F^S(x, a)\} = \log \frac{R(x)}{Z} = \log \frac{F(x)}{Z}$ for the desired flow F . Accordingly, $\tilde{Q}_F^S(s_{T-2}, (s_{T-2} \rightarrow x)) = \tilde{V}_F^S(x) + \log P_B(s_{T-2}|x) = \log \frac{F(x)}{Z} + \log \frac{F(s_{T-2} \rightarrow x)}{F(x)} = \log \frac{F(s_{T-2} \rightarrow x)}{Z}$, and $\tilde{V}_F^S(s_{T-2}) = \log \sum_a \exp\{Q_F^S(s_{T-2}, a)\} = \log \frac{F(s_{T-2})}{Z}$. By induction, we can conclude that $\tilde{Q}_F^S(s, (s \rightarrow s')) = \log \frac{F(s \rightarrow s')}{Z}$, $V_F^S(s) = \log \frac{F(s)}{Z}$ and $\hat{Q}_F^S(s, (s \rightarrow s')) = \log(P_B(s|s')F(s'))$ when they achieve the optima. This further implies that objective (62) can be rewritten as:

$$\mathbb{E}_{P_{\mathcal{D}}(\tau)} \left[\sum_{t=1}^{T-1} \left(F^{\log}(s_{t-1} \rightarrow s_t) - \log \left(P_B(s_{t-1}|s_t) \sum_{s_{t+1}} \exp F^{\log}(s_t \rightarrow s_{t+1}) \right) \right)^2 \right], \quad (64)$$

where $F^{\log}(s_t \rightarrow s_{t+1}) := \tilde{Q}(s_t, a_t) + \log Z$ is the estimator for the logarithm of the desired edge flow, $\log F(s_t \rightarrow s_{t+1})$. This objective is similar to the FM objective, which can be written as:

$$\mathbb{E}_{P_{\mathcal{D}}(\tau)} \left[\sum_{t=1}^T \left(\log \left(\sum_{s_{t-1}} \exp F^{\log}(s_{t-1} \rightarrow s_t) \right) - \log \left(\sum_{s_t} \exp F^{\log}(s_t \rightarrow s_{t+1}) \right) \right)^2 \right]. \quad (65)$$

The reason is that the optimal solution of the objective (64) satisfies $\exp F^{\log}(s_t \rightarrow s_{t+1}) = P_B(s_t | s_{t+1}) \sum_{s_{t+1}} \exp F^{\log}(s_t \rightarrow s_{t+1})$. Taking summation over s_t of this equation, we get $\sum_{s_t} \exp F^{\log}(s_t \rightarrow s_{t+1}) = \sum_{s_{t+1}} \exp F^{\log}(s_t \rightarrow s_{t+1})$, the flow matching condition (2). We can see that $\log Z$ serves as a baseline for modeling $\log F$. Without $\log Z$, we need to approximate $\log F$ by F^{\log} directly. This often leads to numerical issues. For example, let's suppose a small perturbation of $\log F$, denoted as ϵ . Then the flow difference is $\exp(\log F + \epsilon) - F = (\exp \epsilon - 1)F$. As values of F can be exponentially large, especially for nodes near the root, the flow difference can be large even if ϵ is small, making approximation to F by $\exp F^{\log}$ very difficult. By contrast, TB-based methods and our policy-based method allow the updating of $\log Z$ to dynamically scale down the value of \tilde{Q} during training. This also complements the claims by Malkin et al. (2022a), who show that the TB-based method is more efficient than flow-matching and DB-based methods.

B.6. Model parameter updating rules

In the following context, we will explain the updating rules for P_F and μ by two methods, the vanilla policy-based method, also called Actor-Critic method, and the TRPO method. The updating rules for P_B follow analogically.

Actor-Critic First of all, we split the parameters θ into θ_F and θ_Z corresponding to π_F and μ ; since computing the exact V_F is usually intractable, we use \tilde{V}_F parametrized by η as the functional approximation. Given a batch of $\mathcal{D} = \{\tau^n\}_{n=1}^N$, we compute the sampling averaging approximation (SAA) of the following gradient estimators to update θ_F , θ_Z and η as proposed by Schulman et al. (2016) and Tsitsiklis & Van Roy (1996):

$$\begin{aligned} & \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} \hat{A}_F^\lambda(s_t, a_t) \nabla_{\theta_F} \log \pi_F(s_t, a_t; \theta_F) \right] + \mathbb{E}_{\mu(s_0)} [\hat{V}_F^\lambda(s_0) \log \mu(s_0; \theta_Z)], \\ & \mathbb{E}_{P_{F,\mu}(s_t)} \left[\sum_{t=0}^{T-1} \nabla_\eta (\hat{V}_F^\lambda(s_t) - \tilde{V}_F(s_t; \eta))^2 \right], \end{aligned} \quad (66)$$

where $\lambda \in (0, 1)$,

$$\begin{aligned} \hat{A}_F^\lambda(s_t, a_t) &:= \sum_{l=t}^{T-1} \lambda^{l-t} \hat{\delta}_F(s_l, a_l), \quad \hat{V}_F^\lambda(s_t) := \sum_{l=t}^{T-1} \lambda^{l-t} \hat{\delta}_F(s_l, a_l) + \tilde{V}_F(s_t), \\ \hat{\delta}_F(s_t, a_t) &:= R_F(s_t, a_t) + \tilde{V}_F(s_t) - \tilde{V}_F(s_{t+1}). \end{aligned} \quad (67)$$

\hat{A}_F is called **critic** and π_F is called **actor**. It can be verified that $\hat{A}_F^1(s_t, a_t) = \sum_{l=t}^{T-1} R_F(s_l, a_l) - \tilde{V}_F(s_t; \eta)$ renders an unbiased estimator of $\nabla_{\theta_F} J(\theta)$ as the first term is an unbiased estimation of Q_F and $\tilde{V}_F(\cdot; \eta)$ does not introduce estimation bias (see Remark B.7); $\hat{A}_F^0(s_t, a_t) = R(s_t, a_t) + \tilde{V}_F(s_{t+1}) - \tilde{V}_F(s_t)$ provided an direct functional approximation of $A_F(s_t, a_t)$, which usually render biased estimation with lower variance as \tilde{V}_F may not equal to V_F exactly. Thus, λ enables the **variance-bias trade-off** for robust gradient estimation. Likewise, $V_F^1(s_t) = \sum_{l=t}^{T-1} R(s_l, a_l)$ and $V_F^0(s_t) = R(s_t, a_t) + \tilde{V}_F(s_{t+1})$ for each τ , corresponding to unbiased and biased estimation of $V_F(s_t)$. Denoting the estimated gradient w.r.t. $(\theta_F, \theta_Z, \eta)$ as $(\hat{g}_F, \hat{g}_Z, \hat{g}_V)$, these parameters are updated by $(\theta'_F, \theta'_Z, \eta') \leftarrow (\theta_F - \alpha_F \hat{g}_F, \theta_Z - \alpha_Z \hat{g}_Z, \eta - \alpha_V \hat{g}_V)$.

TRPO Parameters θ_Z and η are updated in the same way as the actor-critic method. Parameter θ_F is updated by the optimization objective (6), which can be linearly approximated by:

$$\begin{aligned} & \min_{\theta'_F} \left(T \cdot \nabla_{\theta'_F} \mathbb{E}_{d_{F,\mu}(s;\theta), \pi(s,a;\theta'_F)} [A_F(s, a; \theta)] \right)^\top (\theta'_F - \theta_F) \\ & \text{s.t.} \quad (\theta'_F - \theta_F)^\top \left(\nabla_{\theta'_F}^2 D_{KL}^{d_{F,\mu}(\cdot;\theta)} (\pi_F(s, a; \theta_F), \pi_F(s, a; \theta'_F)) \right) (\theta'_F - \theta_F) \leq \delta_F. \end{aligned} \quad (68)$$

Using \hat{g}_F and \hat{H}_F to denote the estimations of the policy gradient in the objective and the Hessian of the KL divergence in the constraint, model parameters are updated by: $\theta'_F \leftarrow \theta_F - \left(\frac{2\delta_F}{\hat{g}_F^\top \hat{H}_F^{-1} \hat{g}_F} \right)^{0.5} \hat{H}_F^{-1} \hat{g}_F$. When the dimension of θ' is high, computing \hat{H}_F^{-1} is resource and time demanding. We have adopted the conjugate gradient method to estimate $\hat{H}_F^{-1} \hat{g}_F$ based on $\hat{H}_F \hat{g}_F$ (Hestenes et al., 1952).

C. Performance analysis

Lemma C.1. *Given two forward policies (π_F, π'_F) or backward (π_B, π'_B) , we have*

$$\begin{aligned} \frac{1}{T}(J_F - J'_F) &= \mathbb{E}_{d_{F,\mu}(s), \pi(a,s)}[A'_F(s, a)] + D_{KL}^{d_{F,\mu}}(\pi_F(s, a), \pi'_F(s, a)) \\ \frac{1}{T}(J_B - J'_B) &= \mathbb{E}_{d_{B,\rho}(s), \pi(a,s)}[A'_B(s, a)] + D_{KL}^{d_{B,\rho}}(\pi_B(s, a), \pi'_B(s, a)) \end{aligned} \quad (69)$$

Proof.

$$\begin{aligned} V_F(s_0) - V'_F(s_0) &= \mathbb{E}_{P_F(\tau|s_0)} \left[\sum_{t=0}^{T-1} R_F(s_t, a_t) \right] + \underbrace{V'_F(s_T) - V'_F(s_0)}_{:=0} \\ &= \mathbb{E}_{P_F(\tau|s_0)} \left[\sum_{t=0}^{T-1} R_F(s_t, a_t) + V'_F(s_t) - V'_F(s_{t+1}) \right] - V'_F(s_0) + V'_F(s_T) \\ &= \mathbb{E}_{P_F(\tau|s_0)} \left[\sum_{t=0}^{T-1} R_F(s_t, a_t) + V'_F(s_{t+1}) - V'_F(s_t) \right] \\ &= \mathbb{E}_{P_F(\tau|s_0)} \left[\sum_{t=0}^{T-1} A'_F(s_t, a_t) \right] + \mathbb{E}_{P_F(\tau|s_0)} \left[\sum_{t=0}^{T-1} (R_F(s_t, a_t) - R'_F(s_t, a_t)) \right] \end{aligned} \quad (70)$$

Thus,

$$\begin{aligned} J_F - J'_F &= \mathbb{E}_{\mu(s_0)}[V_F(s) - V'_F(s)] \\ &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} A'_F(s_t, a_t) \right] + \mathbb{E}_{P_{F,\mu}(\tau)} \left[\sum_{t=0}^{T-1} (R_F(s_t, a_t) - R'_F(s_t, a_t)) \right] \\ &= \mathbb{E}_{d_{F,\mu}(s), \pi_F(s, a)}[A'_F(s, a)] + \mathbb{E}_{d_{F,\mu}(s), \pi(a, s)}[R_F(s, a) - R'_F(s, a)] \\ &= T \cdot \mathbb{E}_{d_{F,\mu}(s), \pi_F(s, a)}[A'_F(s, a)] + T \cdot \mathbb{E}_{d_{F,\mu}(s)}[D_{KL}(\pi_F(\cdot, s), \pi'_F(\cdot, s))] \end{aligned} \quad (71)$$

The result for the backward case can easily be derived following a similar proof procedure as the forward case. \square

Lemma C.2. (Descent lemma (Beck, 2017)) *Supposing f is a β -smooth function, then for any θ and θ' :*

$$f(\theta') \leq f(\theta) + \langle \nabla_\theta f(\theta), \theta' - \theta \rangle + \frac{\beta}{2} \|\theta - \theta'\|_2^2 \quad (72)$$

C.1. Proof of Theorem 3.5

Proof.

$$\begin{aligned} J_F^G &= J_F + (J_F^G - J_F) \\ &= J_F + \mathbb{E}_{P_{F,\mu}(\tau)} \left[\log \frac{P_F(\tau|s_0)Z}{P_G(\tau|x)R(x)} - \log \frac{P_F(\tau|s_0)Z}{P_B(\tau|x)R(x)} \right] \\ &= J_F + \mathbb{E}_{P_{F,\mu}(\tau)} \left[\log \frac{P_B(\tau|x)}{P_G(\tau|x)} \right] + \mathbb{E}_{P_{B,\rho}(\tau)} \left[\log \frac{P_B(\tau|x)}{\widetilde{P}_G(\tau|x)} \right] - \mathbb{E}_{P_{B,\rho}(\tau)} \left[\log \frac{P_B(\tau|x)}{\widetilde{P}_G(\tau|x)} \right] \\ &= J_F + J_B^G + \mathbb{E}_{\rho(x)}[\log Z_G(x)] + \sum_{\tau} (P_{F,\mu}(\tau) - P_{B,\rho}(\tau)) R_B^G(\tau|x) \end{aligned} \quad (73)$$

where $R_B^G(\tau|x) := \log \frac{P_B(\tau|x)}{P_G(\tau|x)} = \sum_{t=1}^T R_B^G(s_t, a_t)$, then

$$\begin{aligned} J_F^G &= J_F + J_B^G + \langle P_{F,\mu}(\cdot) - P_{B,\rho}(\cdot), R_B^G(\cdot) \rangle \\ &\leq J_F + J_B^G + \|P_{F,\mu}(\cdot) - P_{B,\rho}(\cdot)\|_1 \|R_B^G(\cdot)\|_\infty \\ &\leq J_F + J_B^G + T \cdot \|P_{F,\mu}(\cdot) - P_{B,\rho}(\cdot)\|_1 R_B^{G,max} \end{aligned} \quad (74)$$

where the first inequality follows from Hölder's inequality and the second inequality holds by $\max_{\tau} R_B^G(\tau) \leq T \cdot \max_{s,a} R_B^G(s, a) := T \cdot R_B^{G,max}$. By Pinsker's inequality:

$$\|P_{F,\mu}(\cdot) - P_{B,\rho}(\cdot)\|_1 \leq \sqrt{\frac{1}{2} D_{KL}(P_{F,\mu}(\tau), P_{B,\rho}(\tau))} \quad (75)$$

Besides,

$$\begin{aligned} D_{KL}(P_{F,\mu}(\tau), P_{B,\rho}(\tau)) &= \mathbb{E}_{P_{F,\mu}(\tau)} \left[\log \frac{P_F^\mu(\tau|x) P_{F,\mu}^\top(x)}{P_B(\tau|x) P_{F,\mu}^\top(x)} \right] \\ &\leq \mathbb{E}_{P_{F,\mu}(\tau)} \left[\log \frac{P_F^\mu(\tau|x)}{P_B(\tau|x)} \right] + \underbrace{\mathbb{E}_{P_{F,\mu}^\top(x)} \left[\log \frac{P_{F,\mu}^\top(x)}{R(x)/Z^*} \right]}_{\geq 0} \\ &= D_{KL}(P_{F,\mu}(\tau), P_B(\tau)) = D_{KL}^\mu(P_F(\tau), \tilde{P}_B(\tau)) - \log Z + \log Z^* \\ &= J_F + \log Z^* - \log Z \end{aligned} \quad (76)$$

Then we have:

$$J_F^G \leq J_F + J_B^G + \mathbb{E}_{\rho(x)}[\log Z_G(x)] + R_B^{G,max} \sqrt{\frac{1}{2} (J_F + \log Z - \log Z^*)} \quad (77)$$

□

C.2. Proof of Theorem 3.6

Proof. By Lemma C.1:

$$\frac{1}{T} (J'_F - J_F) = \mathbb{E}_{d'_{F,\mu}(s), \pi'_F(a,s)} [A_F(s, a)] + \underbrace{D_{KL}^{d'_{F,\mu}}(\pi'_F(\cdot, s), \pi_F(\cdot, s))}_{\leq \zeta_F} \quad (78)$$

Let $\bar{A}_F \in R^{|S|}$ denote the vector components of $\mathbb{E}_{\pi'_F(s,a)} [A_F(s, a)]$. Then we have

$$\begin{aligned} \mathbb{E}_{d'_{F,\mu}(s) \pi'_F(s,a)} [A_F(s, a)] &= \langle d'_{F,\mu}, \bar{A}_F \rangle \\ &= \langle d_{F,\mu}, \bar{A}_F \rangle + \langle d'_{F,\mu} - d_{F,\mu}, \bar{A}_F \rangle \\ &\leq \mathbb{E}_{d_{F,\mu}(s) \pi'_F(s,a)} [A_F(s, a)] + \|d'_{F,\mu} - d_{F,\mu}\|_1 \|\bar{A}_F\|_\infty \end{aligned} \quad (79)$$

By Lemma B.5:

$$\mathbb{E}_{d_{F,\mu}(s) \pi'_F(s,a)} [A_F(s, a)] \leq \mathbb{E}_{d_{F,\mu}(s) \pi'_F(a,s)} [A_F(s, a)] + 2 \mathbb{E}_{d_{F,\mu}(s)} [D_{TV}(\pi'_F(s, \cdot), \pi_F(s, \cdot))] \epsilon_F \quad (80)$$

By Pinsker's inequality:

$$\mathbb{E}_{d'_{F,\mu}(s) \pi'_F(s,a)} [A_F(s, a)] \leq \mathbb{E}_{d_{F,\mu}(s) \pi'_F(a,s)} [A_F(s, a)] + 2 \mathbb{E}_{d_{F,\mu}(s)} \left[\left(\frac{1}{2} D_{KL}(\pi'_F(\cdot, s), \pi_F(\cdot, s)) \right)^{0.5} \right] \epsilon_F \quad (81)$$

By Jensen's inequality:

$$\mathbb{E}_{d'_{F,\mu}(s) \pi'_F(s,a)} [A_F(s, a)] \leq \mathbb{E}_{d_{F,\mu}(s) \pi'_F(a,s)} [A_F(s, a)] + 2 \left(\frac{1}{2} \mathbb{E}_{d_{F,\mu}(s)} [D_{KL}(\pi'_F(\cdot, s), \pi_F(\cdot, s))] \right)^{0.5} \epsilon_F \quad (82)$$

Thus, we have

$$\frac{1}{T}(J'_F - J_F) \leq E_{d_{F,\mu}(s)\pi'_F(a,s)}[A_F(s,a)] + (2\zeta_F)^{0.5}\epsilon_F + \zeta_F \quad (83)$$

□

C.3. Proof of Theorem 3.7

Proof. By Lemma C.2,

$$J_F(\theta_{n+1}) \leq J_F(\theta_n) + \langle \nabla_{\theta_n} J_F(\theta_n), \theta_{n+1} - \theta_n \rangle + \frac{\beta}{2} \|\theta_{n+1} - \theta_n\|_2^2. \quad (84)$$

Thus,

$$\begin{aligned} -\langle \nabla_{\theta_n} J_F(\theta_n), \theta_{n+1} - \theta_n \rangle &\leq J_F(\theta_n) - J_F(\theta_{n+1}) + \frac{\beta}{2} \|\theta_{n+1} - \theta_n\|_2^2, \\ \alpha \langle \nabla_{\theta_n} J_F(\theta_n), \widehat{\nabla}_{\theta_n} J_F(\theta_n) \rangle &\leq J_F(\theta_n) - J_F(\theta_{n+1}) + \frac{\beta\alpha^2}{2} \|\widehat{\nabla}_{\theta_n} J_F(\theta_n)\|_2^2. \end{aligned}$$

Conditioning on θ_n , taking expectations over both sides and noting that $\mathbb{E}_{P(\cdot|\theta_n)} [\langle \nabla_{\theta_n} J_F(\theta_n), \widehat{\nabla}_{\theta_n} J_F(\theta_n) \rangle] = \langle \nabla_{\theta_n} J_F(\theta_n), \mathbb{E}_{P(\cdot|\theta_n)} [\widehat{\nabla}_{\theta_n} J_F(\theta_n)] \rangle = \|\nabla_{\theta_n} J_F(\theta_n)\|_2^2$, we have

$$\alpha \|\nabla_{\theta_n} J_F(\theta_n)\|_2^2 \leq J_F(\theta_n) - \mathbb{E}_{P(\theta_{n+1}|\theta_n)} [J_F(\theta_{n+1})] + \frac{\beta\alpha^2}{2} \mathbb{E}_{P(\cdot|\theta_n)} [\|\widehat{\nabla}_{\theta_n} J_F(\theta_n)\|_2^2]. \quad (85)$$

By the assumption that $\mathbb{E}_{P(\cdot|\theta)} [\|\widehat{\nabla}_{\theta} J_F(\theta) - \nabla_{\theta} J_F(\theta)\|_2^2] = \mathbb{E}_{P(\cdot|\theta)} [\|\widehat{\nabla}_{\theta} J_F(\theta)\|_2^2] - \|\nabla_{\theta} J_F(\theta)\|_2^2 \leq \sigma_F$, we have:

$$\alpha \|\nabla_{\theta_n} J_F(\theta_n)\|_2^2 \leq J_F(\theta_n) - \mathbb{E}_{P(\theta_{n+1}|\theta_n)} [J_F(\theta_{n+1})] + \frac{\beta\alpha^2}{2} \|\nabla_{\theta_n} J_F(\theta_n)\|_2^2 + \frac{\beta\alpha^2\sigma_F}{2}. \quad (86)$$

Consequently, we have:

$$\begin{aligned} \left(\alpha - \frac{\beta\alpha^2}{2} \right) \mathbb{E}_{P(\theta_{0:N-1})} \left[\sum_{n=0}^{N-1} \|\nabla_{\theta_n} J_F(\theta_n)\|_2^2 \right] &\leq \frac{N\beta\alpha^2\sigma_F}{2} + \mathbb{E}_{P(\theta_{0:N})} \left[\sum_{n=0}^{N-1} J_F(\theta_n) - J_F(\theta_{n+1}) \right], \\ \left(\alpha - \frac{\beta\alpha^2}{2} \right) \sum_{n=0}^{N-1} \mathbb{E}_{P(\theta_n)} [\|\nabla_{\theta_n} J_F(\theta_n)\|_2^2] &\leq \frac{N\beta\alpha^2\sigma_F}{2} + \mathbb{E}_{P(\theta_{0:N})} [J_F(\theta_0) - J_F(\theta_N)], \\ \left(\alpha - \frac{\beta\alpha^2}{2} \right) N \min_{n=0,\dots,N-1} \mathbb{E}_{P(\theta_n)} [\|\nabla_{\theta_n} J_F(\theta_n)\|_2^2] &\leq \frac{N\beta\alpha^2\sigma_F}{2} + \mathbb{E}_{P(\theta_0)} [J_F(\theta_0)] - \mathbb{E}_{P(\theta_N)} [J_F(\theta_N)]. \end{aligned} \quad (87)$$

Setting $\alpha = \sqrt{2/(\beta N)}$, we have:

$$\left(\sqrt{(2N)/\beta} - 1 \right) \min_{n=0,\dots,N-1} \mathbb{E}_{P(\theta_n)} [\|\nabla_{\theta_n} J_F(\theta_n)\|_2^2] \leq \sigma_F + \mathbb{E}_{P(\theta_0)} [J_F(\theta_0)] - \mathbb{E}_{P(\theta_N)} [J_F(\theta_N)]. \quad (88)$$

Since $J_F(\theta) + \log Z^* - \log Z(\theta) = D_{KL}^{\mu(\cdot|\theta)}(P_F(\tau|s_0; \theta), P_B(\tau|s_0))$ and $J_F(\theta^*) = 0$ with $\log Z^* - \log Z(\theta^*)$ for optimal parameter θ^* , then $J_F(\theta_N) + \log Z^* - \log Z(\theta_N) \geq J_F(\theta^*)$ and we have:

$$\begin{aligned} \min_{n=0,\dots,N-1} \mathbb{E}_{P(\theta_n)} [\|\nabla_{\theta_n} J_F(\theta_n)\|_2^2] &\leq \frac{\sigma_F + \mathbb{E}_{P(\theta_0)} [J_F(\theta_0)] - \mathbb{E}_{P(\theta_N)} [J_F(\theta_N) + \log Z^* - \log Z(\theta_N)]}{\left(\sqrt{(2N)/\beta} - 1 \right)} \\ &\quad + \frac{\mathbb{E}_{P(\theta_N)} [\log Z^* - \log Z(\theta_N)]}{\left(\sqrt{(2N)/\beta} - 1 \right)} \\ &\leq \frac{\sigma_F + \mathbb{E}_{P(\theta_0)} [J_F(\theta_0)] + \mathbb{E}_{P(\theta_N)} [|\log Z^* - \log Z(\theta_N)|]}{\left(\sqrt{(2N)/\beta} - 1 \right)}. \end{aligned} \quad (89)$$

By the assumption that $|\log Z - \log Z^*| \leq \sigma_Z$, we have:

$$\min_{n=0,\dots,N-1} \mathbb{E}_{P(\theta_n)} \left[\|\nabla_{\theta_n} J_F(\theta_n)\|_2^2 \right] \leq \frac{\sigma_F + \sigma_Z + \mathbb{E}_{P(\theta_0)}[J_F(\theta_0)]}{\left(\sqrt{(2N)/\beta} - 1 \right)}. \quad (90)$$

□

D. Additional experimental settings and results

In all experiments, we follow a regular way of sample policy design for TB-based and DB-based methods: sample policy is a mix of the learned forward policy and a uniform policy where the mix-in factor of the uniform policy starts at 1.0 and decays exponentially at a rate of 0.99 after each training iteration. For our policy-based methods, we set the value of hyper-parameter λ to 0.99. The values of the decay rate and λ are selected based on the results of the ablation study. Trust region hyper-parameter ζ_F is set to 0.01 select from $\{0.01, 0.02, 0.03, 0.04, 0.05\}$. In the set of biological and molecular sequence design experiments, we use a temperature hyper-parameter β for $R^\beta(x)$. following (Shen et al., 2023) and set it to 3, 3 and 5 for SIX6, PHO4 and QM9 dataset respectively. We use the Adam optimizer for model optimization. The learning rates of forward and backward policy are equal to 1×10^{-3} , which is selected from $\{5 \times 10^{-3}, 1 \times 10^{-3}, 5 \times 10^{-4}, 1 \times 10^{-4}\}$ by TB-U. The learning rates of value functions are set to 5×10^{-3} , which is selected from $\{1 \times 10^{-2}, 5 \times 10^{-3}, 1 \times 10^{-3}\}$ by RL-U. The learning rates of total flow estimator is 1×10^{-1} , which is selected from $\{1 \times 10^{-1}, 5 \times 10^{-2}, 1 \times 10^{-2}, 5 \times 10^{-3}\}$ by TB-U. The sample batch size is set to 128 for each optimization iteration. For all experiments, we report the performance with five different random seeds.

D.1. Evaluation metrics

The total variation D_{TV} between $P_F^\top(x)$ and $P^*(x)$ is defined as:

$$D_{TV}(P_F^\top, P^*) = \frac{1}{2} \sum_{x \in \mathcal{X}} |P_F^\top(x) - P^*(x)|. \quad (91)$$

The total variation is similar to the average l_1 -distance used in previous works, which can be computed by $\frac{1}{|\mathcal{X}|} \sum_x |P^*(x) - P_F^\top(x)|$. However, the average l_1 -distance may be inappropriate as $|\mathcal{X}|$ is usually large ($> 10^4$) and $\sum_x |P^*(x) - P_F^\top(x)| \leq 2$, resulting in the expected l_1 -distance being heavily scaled down by $|\mathcal{X}|$.

The Jensen–Shannon divergence D_{JSD} between $P_F^\top(x)$ and $P^*(x)$ is defined as:

$$D_{JSD}(P_F^\top, P^*) = \frac{1}{2} D_{KL}(P_F^\top, M) + \frac{1}{2} D_{KL}(P^*, M), \quad P^M = P_F^\top + P^*. \quad (92)$$

Following Shen et al. (2023), the mode accuracy Acc of $P_F^\top(x)$ w.r.t. $P^*(x)$ is defined as:

$$Acc(P_F^\top, P^*) = \min \left(\frac{\mathbb{E}_{P_F^\top(x)} R(x)}{\mathbb{E}_{P^*(x)} R(x)}, 1 \right). \quad (93)$$

D.2. Hyper-grid modeling

For value-based methods like TB-based training strategy, the policy P_D for training data sampling works in an offline way. The design of guided distribution follows the way proposed by Shen et al. (2023). Explicit designs are shown in the following sections.

Environment In this hyper-grid environment, $\mathcal{S} \setminus \{s^f\}$ is equal to $\{s = ([s]_1, \dots, [s]_D) | \forall i \in [1, \dots, D], [s]_i \in [0, \dots, N-1]\}$ where the initial state $s^0 = (0, \dots, 0)$ and s^f can be represented by any invalid coordinate tuple of the hyper-grid, and is denoted as $(-1, \dots, -1)$ in our implementation. For each state $s \neq s^f$, we have $D+1$ possible actions in $\mathcal{A}(s)$: (1) increment the i_{th} coordinate by one, arriving at $s' = ([s]_0, \dots, [s]_i + 1, \dots)$; (2) choose stopping actions ($s \rightarrow s^f$), terminating the process and return $x = s$ as the terminating coordinate tuple. By definition, \mathcal{G} is not a

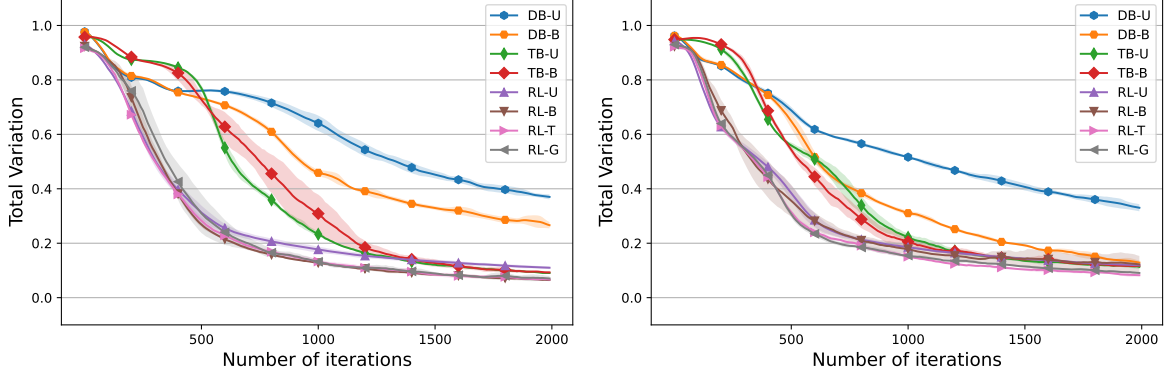


Figure 4. Training curves by D_{TV} between P_F^\top and P^* for $64 \times 64 \times 64$ (left) and $32 \times 32 \times 32 \times 32$ hyper-grids (right). The curves are plotted based on their mean and standard deviation values across five runs.

graded DAG, and $\mathcal{S} \setminus \{s^f\} = \mathcal{X}$ as all coordinate tuples can be returned as the terminating states. The reward $R(x)$ is defined as:

$$R(x) = R_0 + R_1 \prod_{d=1}^D \mathbb{I} \left[\left| \frac{[s]_d}{N-1} - 0.5 \right| \in (0.25, 0.5] \right] + R_2 \prod_{d=1}^D \mathbb{I} \left[\left| \frac{[s]_d}{N-1} - 0.5 \right| \in (0.3, 0.4] \right] \quad (94)$$

where $R_0 = 10^{-2}$, $R_1 = 0.5$ and $R_2 = 2$ in our experiment. For a trajectory $\tau = (s_0, \dots, s_i, \dots, s_f)$, the probability of the guided distribution can be written as:

$$P_G(\tau|x) = \prod_i P_G(s_i|s_{i-1}, x);$$

$$\forall s_i \neq s^f : P_G(s_i|s_{i-1}, x) = \begin{cases} \frac{P_F(s_i|s_{i-1})}{\sum_{s:s \neq s^f} P_F(s|s_{i-1}) + \epsilon^f}, & \text{if } R(s_{i-1}) \leq R_0 \\ P_F(s_i|s_{i-1}), & \text{otherwise} \end{cases}$$

$$P_G(s^f|s_{i-1}, x) = \begin{cases} \frac{\epsilon^f}{\sum_{s:s \neq s^f} P_F(s|s_{i-1}) + \epsilon^f}, & \text{if } R(s_{i-1}) \leq R_0 \\ P_F(s^f|s_{i-1}), & \text{otherwise} \end{cases} \quad (95)$$

where $s_i \neq s^f$ and $\epsilon^f = 10^{-5}$. As all the coordinate tuples can be terminating states, this is, s_f is the child of all the other states, the expression above means that for a state s_{i-1} with low reward, its probability of being a terminating state is replaced by a small value ϵ^f . In this way, we discourage the generative process from stopping early at low reward coordinate tuples.

Model Architecture Policy P_F is parametrized by a neural network with 4 hidden layers and the hidden dimension is 256. Policy P_B is fixed to be uniform over valid actions or parameterized in the same way as P_F . Coordinate tuples are transformed by K-hot encoding before being fed into Neural Networks.

Additional Experiment Results The obtained results of $64 \times 64 \times 64$ and $32 \times 32 \times 32 \times 32$ grids across five runs are summarized in Fig. 4 and Table 2. The graphical illustrations of $P_F^\top(x)$ are shown in Figs. 12 and Figs. 13. We can observe similar performance trends as in 256×256 and 128×128 grids, where both TB-based methods and our policy-based method are better than DB-based method, and our policy-based methods achieve much faster convergence than TB-based method. These trends are less obvious than in 256×256 and 128×128 grids. This phenomenon can be ascribed to that environment height H have more influence on the modeling difficulty than environment dimension N . The reason is that hyper-grids are homogeneous w.r.t. each dimension, and the minimum distance between modes only depends on H .

D.3. Biological and Molecular Sequence design

Environment In this environment, $\mathcal{S} = \{-1, 0, \dots, N-1\}^D$ with element s corresponds to a sequence composed of integers ranging from -1 and $N-1$ logits. The set $\{0, \dots, N-1\}$ denotes the nucleotide types or building blocks, and the

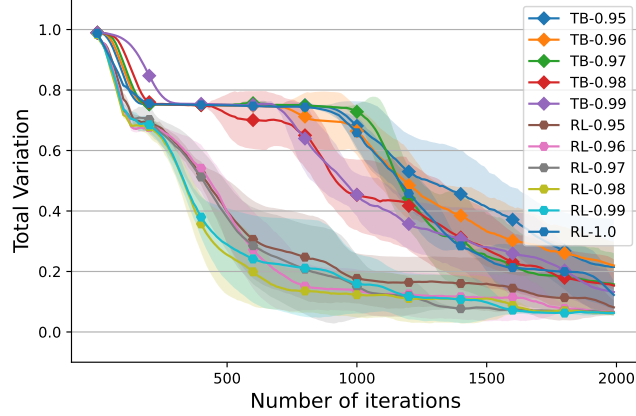


Figure 5. Performance comparison between RL-U with different λ values and TB-U with different decay rates of the mix-in factor. The curves are plotted based on their mean and standard deviation values.

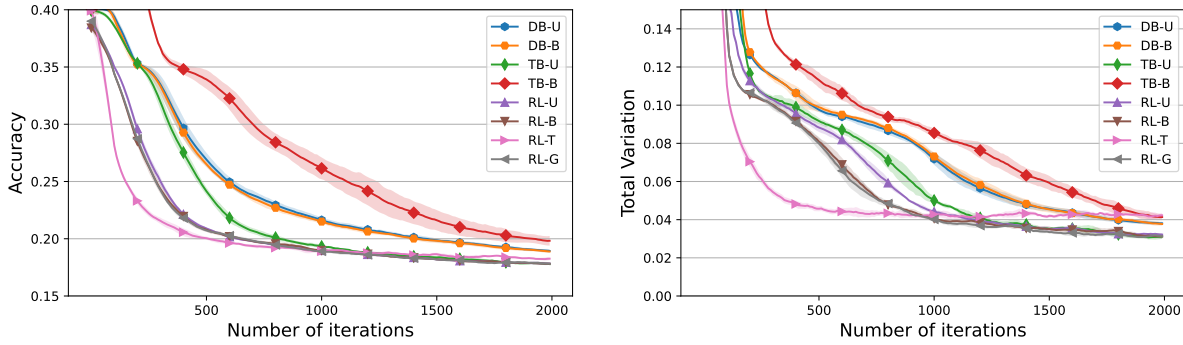


Figure 6. Training curves by D_{TV} between P_F^T and P^* for SIX6 (left) and QM9 (right). The curves are plotted based on their mean and standard deviation values.

integer -1 represents that the corresponding position within s is unfilled.

The initial empty sequence s^0 is represented by $\{-1\}^D$. For $s_t \in \mathcal{S}_t$, there are t elements in $\{0, \dots, N-1\}$ and the rest equal to -1 . There are $N \cdot (D-t)$ actions in $\mathcal{A}(s)$ that correspond to fill in one of the empty slots by one integer in $\{0, \dots, N-1\}$. The generative process will not stop until sequences are fulfilled. By definition, \mathcal{G} is a graded DAG and $\mathcal{S}_N = \mathcal{X} = \{0, \dots, N\}^D$. We use the reward values provided in the dataset directly, and following Shen et al. (2023), normalize rewards to $[1 \times 10^{-3}, 10]$, $[1 \times 10^{-3}, 10]$ and $[1 \times 10^{-3}, 100]$ for the SIX6, PHO4 and QM9 datasets respectively. The guided distribution design also follows Shen et al. (2023).

Model Architecture Policies are constructed in the same way as the hyper-grid modeling experiment.

Additional Experiment Results For PHO4 dataset, the exact computation of P_F^T by dynamic programming is expensive. Thus, we only plot the training curves by Acc as shown in Fig. 7, and the mean and standard deviation of Acc values at the last iteration are provided in Table 3 where P_F^T is empirically estimated by sampling. According to Fig. 7, we can observe similar performance trends as in the QM9 and SIX6 datasets. Our policy-based methods achieve faster convergence rates and better converged metric values than the TB-based and DB-based methods. While the converged Acc of RL-T is slightly worse than the other policy-based methods, it achieves the fastest convergence rate. RL-G and RL-B, both employing a parametrized π_B , demonstrate similar performance and converge faster than RL-U, which utilizes a uniform π_B .

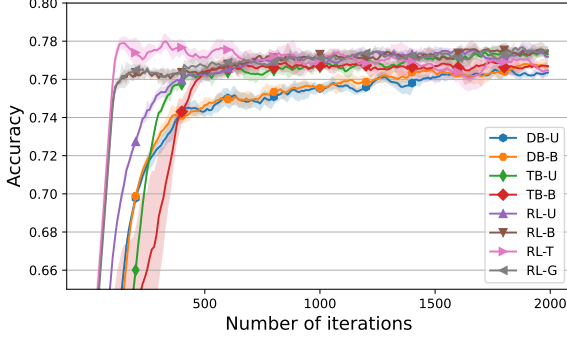


Figure 7. Training curves by Acc of P_F^\top w.r.t P^* for PHO4. The curves are plotted based on their mean and standard deviation values.

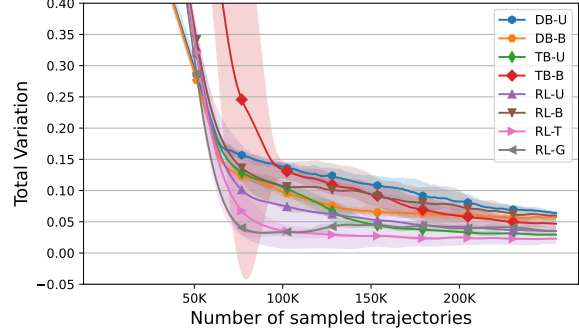


Figure 8. Training curves by D_{KL} of P_F^\top w.r.t P^* for the BN structure learning experiment. The curves are plotted based on their mean and standard deviation values.

D.4. Bayesian Network Structure Learning

In this experiment, we investigate GFlowNets for Bayesian Network (BN) structure learning following the settings adopted in Malkin et al. (2022b). The set \mathcal{X} in GFlowNets here corresponds to a set of BN structures, which are also DAGs. BN structure learning can be understood as approximating $P(x|\mathcal{D}) \propto R(x)$ given a dataset \mathcal{D} . With a set of nodes, the state space for GFlowNets is the set of all possible DAGs over the given nodes. The actions correspond to adding edges over a DAG without introducing a cycle. The generative process of a BN structure is interpreted as starting from an empty graph, an action is taken to decide to add an edge or terminate the generative process at the current graph structure.

Environment A Bayesian Network is a probabilistic model that represents the joint distribution of N random variable and the joint distribution factorizes according to the network structure x :

$$P(y_1, \dots, y_N) = \prod_{n=1}^N P(y_n | Pa_x(y_n)) \quad (96)$$

where $Pa_x(y_n)$ denote the set of parent nodes of y_n according to graph x . As the structure of any graph can be represented by its adjacency matrix, the state space can be defined as $\mathcal{S} := \{s | \mathcal{C}(s) = 0\} \in \{0, 1\}^{N \times N}$ where \mathcal{C} corresponds to the acyclic graph constraint (Deleu et al., 2022), the initial state $s^0 = \{0\}^{N \times N}$ and specially $s^f := \{-1\}^{N \times N}$ in our implementation. For each state s , $a \in \mathcal{A}(s)$ can be any action that turns one of 0 values of s to be 1 (i.e. adding an edge) while keeping $\mathcal{C}(s') = 0$ for the resulting graph s' , or equal to $(s \rightarrow s^f)$ that stopping the generative process and return $x = s$ as the terminating state. By definition, the corresponding \mathcal{G} is not a graded DAG.

Given observation dataset \mathcal{D}_y of $y_{1:N}$, the structure learning problem can be understood as approximating $P(x|\mathcal{D}_y) \propto P(x, \mathcal{D}_y) = P(\mathcal{D}_y|x)P(x)$. Without additional information about graph structure x , $P(x)$ is often assumed to be uniform. Thus, $P(x|\mathcal{D}_y) \propto P(\mathcal{D}_y|x)$ and the reward function is defined as $R(x) \propto P(\mathcal{D}_y|x)$. Distribution $P(\mathcal{D}_y|x)$ is also called graph score and we use *BGe* score (Kuipers et al., 2014) in our experiment. Following Malkin et al. (2022b), the ground-truth graph structure is generated from Erdős–Rényi model, and the observation dataset \mathcal{D}_y with 100 samples is simulated from the ground-truth graph. The guided distribution design follows the hyper-grid experiment. A low probability value, 10^{-5} is assigned to the transition probability of $(s \rightarrow s^f)$ if $|\log R^{\max} - \log R(s)| \geq 10$.

Model Architecture Policies are constructed in the same way as the hyper-grid modeling experiment, but adjacency matrices are fed into neural networks directly without encoding.

Experiment Results The number of possible DAGs grows exponentially with the number of nodes. Thus, we here test the same benchmark in Malkin et al. (2022b) with the number of nodes set to 5 and the corresponding total numbers of DAGs is about 2.92×10^4 . The experimental results across five runs are shown in Fig. 8 and Table. 4. The converged D_{KL} of all methods are close, indicating that these methods achieve similar performances. This can be ascribed to the fact that the

state space is relatively small. Nevertheless, RL-T achieves the best performance, and both RL-T and RL-G achieve fast convergence. These results further demonstrate the effectiveness of our policy-based methods for GFlowNet training.

D.5. Performance Tables of $P_F^\top(x)$

Method	256 × 256			128 × 128		
	$D_{TV}(\downarrow)$	$D_{JSD}(\downarrow)$	Time	$D_{TV}(\downarrow)$	$D_{JSD}(\downarrow)$	Time
DB-U	0.605 ± 0.013	0.246 ± 0.011	37.2	0.324 ± 0.012	0.084 ± 0.005	18.5
DB-B	0.446 ± 0.014	0.144 ± 0.006	46.8	0.262 ± 0.013	0.056 ± 0.005	21.3
TB-U	0.073 ± 0.009	0.008 ± 0.001	30.9	0.045 ± 0.007	0.003(3) ± 0.001	14.3
TB-B	0.210 ± 0.105	0.066 ± 0.050	32.5	0.044 ± 0.008	0.003(5) ± 0.001	16.4
RL-U	0.062 ± 0.006	0.008 ± 0.001	44.4	0.043 ± 0.008	0.004 ± 0.001	18.3
RL-B	0.070 ± 0.019	0.010 ± 0.003	68.5	0.039 ± 0.015	0.003(1) ± 0.002	35.2
RL-T	0.070 ± 0.006	0.008 ± 0.001	90.2	0.050 ± 0.004	0.004 ± 0.001	58.3
RL-G	0.043 ± 0.004	0.005 ± 0.000	69.8	0.049 ± 0.009	0.004 ± 0.001	35.4

Table 1. Converged metric values of different methods for the modeling of 256 × 256 and 128 × 128 grids. Time costs are provided in minutes.

Method	64 × 64 × 64			32 × 32 × 32 × 32		
	$D_{TV} \downarrow$	$D_{JSD} \downarrow$	Time	$D_{TV} \downarrow$	$D_{JSD} \downarrow$	Time
DB-U	0.369 ± 0.013	0.108 ± 0.006	19.3	0.325 ± 0.015	0.075 ± 0.007	11.9
DB-B	0.260 ± 0.008	0.066 ± 0.005	21.5	0.125 ± 0.004	0.017 ± 0.001	14.6
TB-U	0.087 ± 0.007	0.007 ± 0.001	16.1	0.108 ± 0.008	0.011 ± 0.001	17.3
TB-B	0.090 ± 0.007	0.008 ± 0.001	17.4	0.119 ± 0.011	0.013 ± 0.002	20.6
RL-U	0.108 ± 0.006	0.013 ± 0.001	17.4	0.118 ± 0.005	0.015 ± 0.001	11.5
RL-B	0.064(7) ± 0.011	0.004 ± 0.001	30.2	0.120 ± 0.032	0.013 ± 0.007	25.3
RL-T	0.065 ± 0.007	0.005 ± 0.001	49.0	0.083 ± 0.010	0.006 ± 0.001	59.3
RL-G	0.067 ± 0.009	0.005 ± 0.002	31.8	0.088 ± 0.007	0.007 ± 0.001	24.4

Table 2. Converged metric values of different methods for the modeling of 64 × 64 × 64 and 32 × 32 × 32 × 32 grids. Time costs are provided in minutes.

Method	SIX6			QM9			PHO4
	$Acc \uparrow$	$D_{TV} \downarrow$	$D_{JSD} \downarrow$	$Acc \uparrow$	$D_{TV} \downarrow (\times 10^{-2})$	$D_{JSD} \downarrow (\times 10^{-3})$	$Acc \uparrow$
DB-U	0.862 ± 0.006	0.189 ± 0.002	0.032 ± 0.001	0.964 ± 0.006	3.80 ± 0.08	1.77 ± 0.16	0.763 ± 0.002
DB-B	0.864 ± 0.009	0.188 ± 0.001	0.032 ± 0.001	0.965 ± 0.003	3.77 ± 0.08	1.74 ± 0.16	0.766 ± 0.003
TB-U	0.941 ± 0.004	0.178(2) ± 0.000	0.029 ± 0.000	0.988 ± 0.003	3.14 ± 0.12	1.11 ± 0.16	0.773 ± 0.002
TB-B	0.904 ± 0.007	0.198 ± 0.004	0.035 ± 0.001	0.984 ± 0.002	4.13 ± 0.16	1.89 ± 0.25	0.767 ± 0.002
RL-U	0.940 ± 0.003	0.178 ± 0.001	0.030 ± 0.000	0.986 ± 0.002	3.21 ± 0.06	1.11 ± 0.09	0.772 ± 0.002
RL-B	0.944 ± 0.003	0.178(2) ± 0.001	0.030 ± 0.000	0.988 ± 0.002	3.12 ± 0.10	1.07 ± 0.13	0.774 ± 0.003
RL-T	0.934 ± 0.002	0.198 ± 0.001	0.030 ± 0.001	0.989(5) ± 0.005	3.72 ± 0.12	2.26 ± 0.03	0.765 ± 0.003
RL-G	0.946 ± 0.0017	0.178(1) ± 0.001	0.029 ± 0.000(4)	0.989(2) ± 0.002	3.14 ± 0.13	1.10 ± 0.13	0.775 ± 0.002

Table 3. Converged metric values of different methods for the SIX6, QM9 and PHO4 datasets.

Method	$D_{TV}(\times 10^{-2})$	$D_{JSD}(\times 10^{-3})$	Method	$D_{TV}(\times 10^{-2})$	$D_{JSD}(\times 10^{-3})$
DB-U	6.41 ± 0.42	6.62 ± 0.32	DB-B	5.81 ± 0.52	6.01 ± 0.12
TB-U	2.95 ± 0.26	2.84 ± 0.07	TB-B	4.69 ± 0.80	5.22 ± 1.23
RL-U	3.52 ± 2.03	4.33 ± 3.13	RL-B	6.10 ± 0.67	6.31 ± 1.08
RL-T	2.35 ± 0.13	2.15 ± 0.16	RL-G	3.72 ± 0.49	4.60 ± 1.61

Table 4. Converged metric values of different methods for the BN structure experiment.

D.6. Graphical Representation of P_F^\top

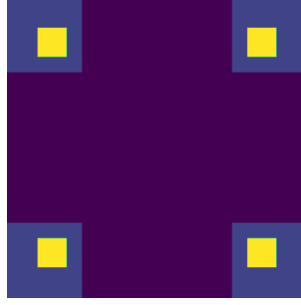


Figure 9. Graphical representation of P^* for a two-dimensional hyper-grid.

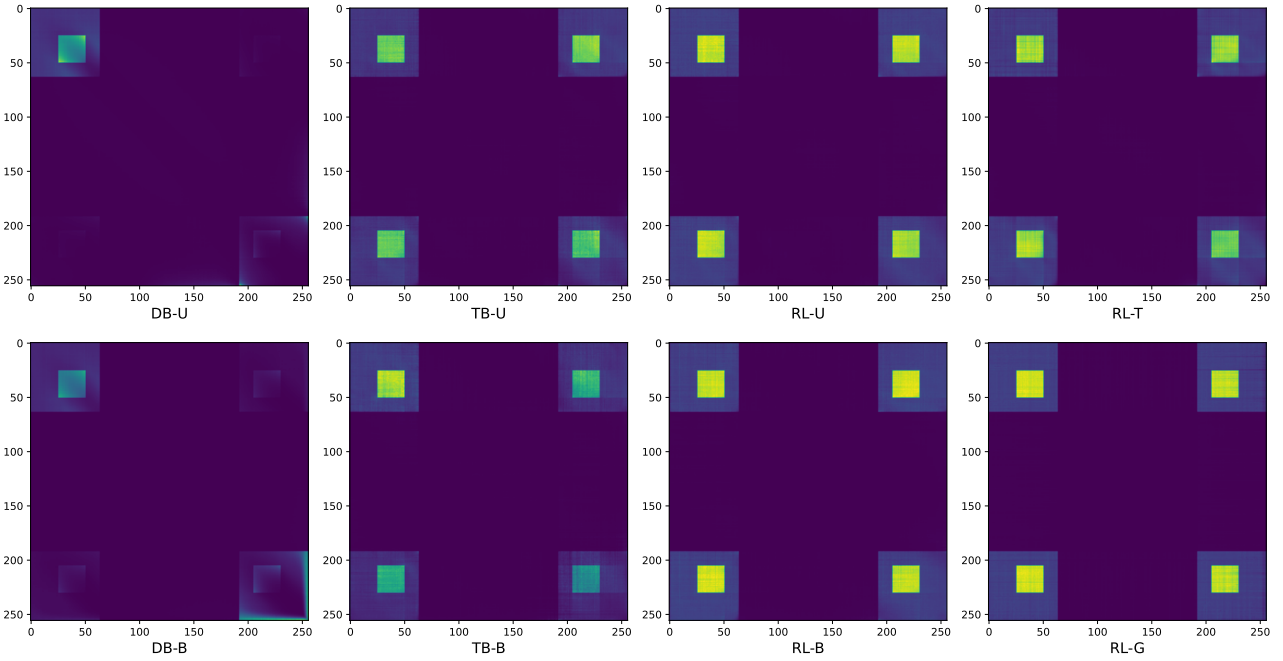


Figure 10. Graphical illustrations of $P_F^\top(x)$ averaged across 5 runs of different training strategies for a 256×256 hyper-grid.

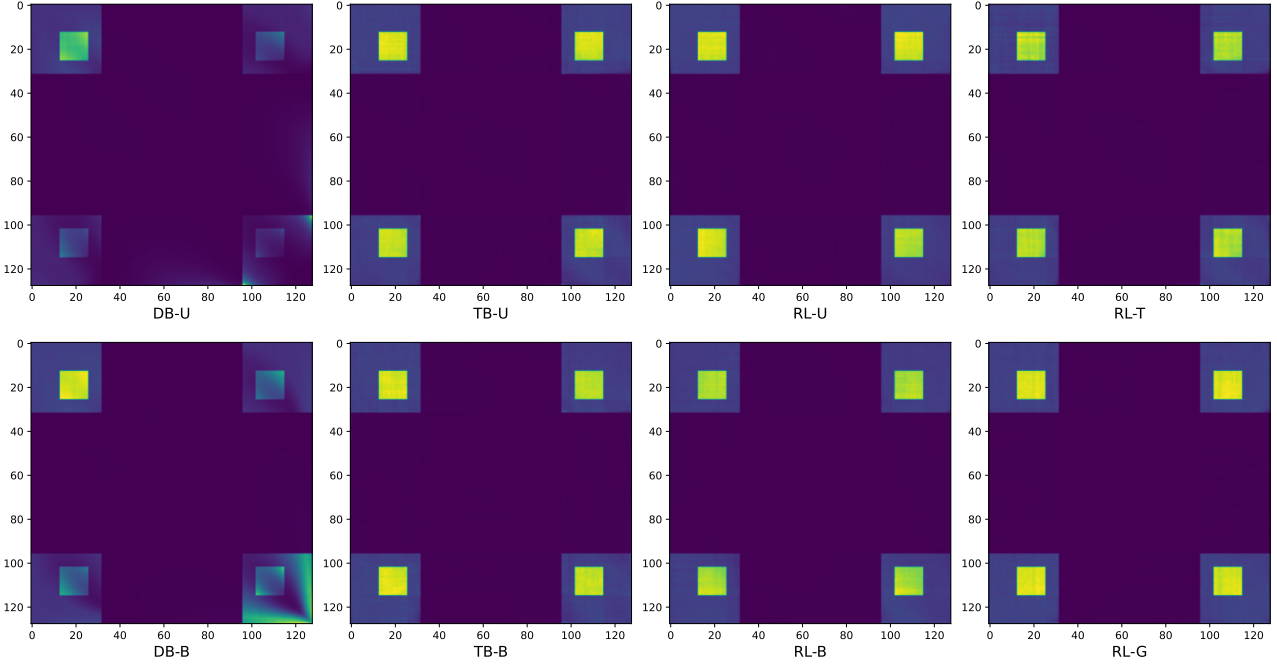


Figure 11. Graphical illustrations of $P_F^T(x)$ averaged across 5 runs of different training strategies for a 128×128 hyper-grid.

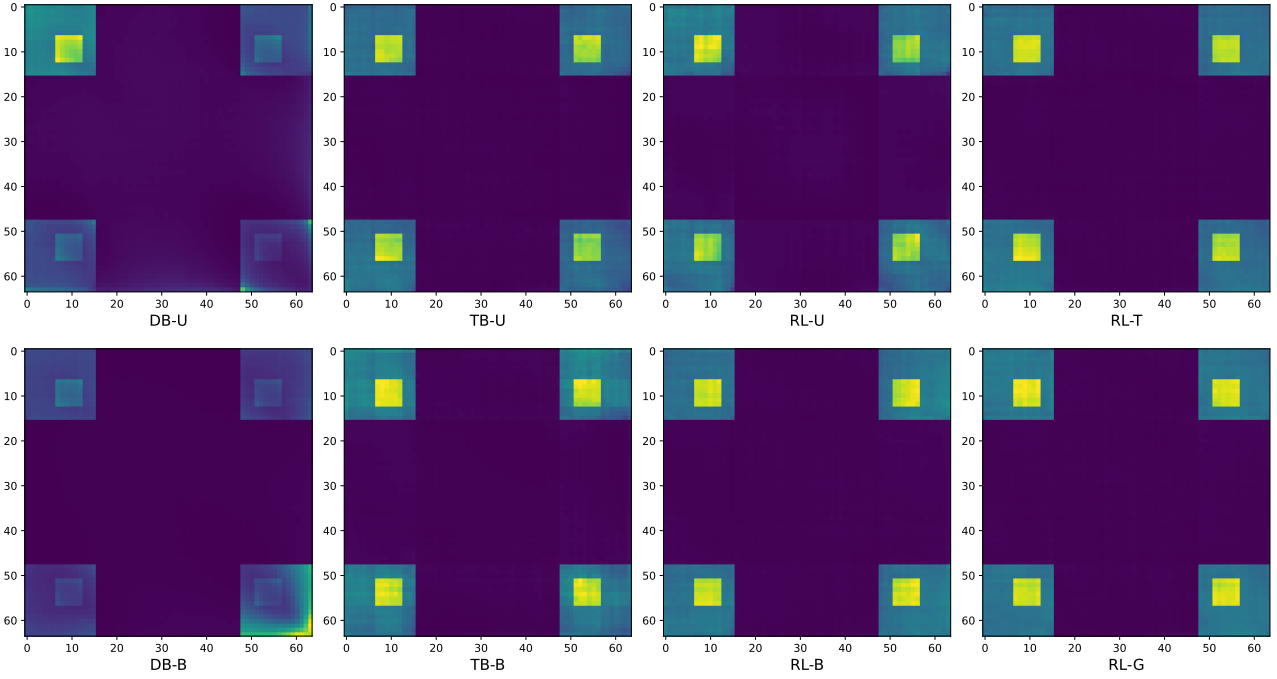


Figure 12. Graphical illustrations of $P_F^T(x)$ averaged across 5 runs of different training strategies for a $64 \times 64 \times 64$ hyper-grid. For visualization easiness, only the marginals of 2 dimensions are plotted.

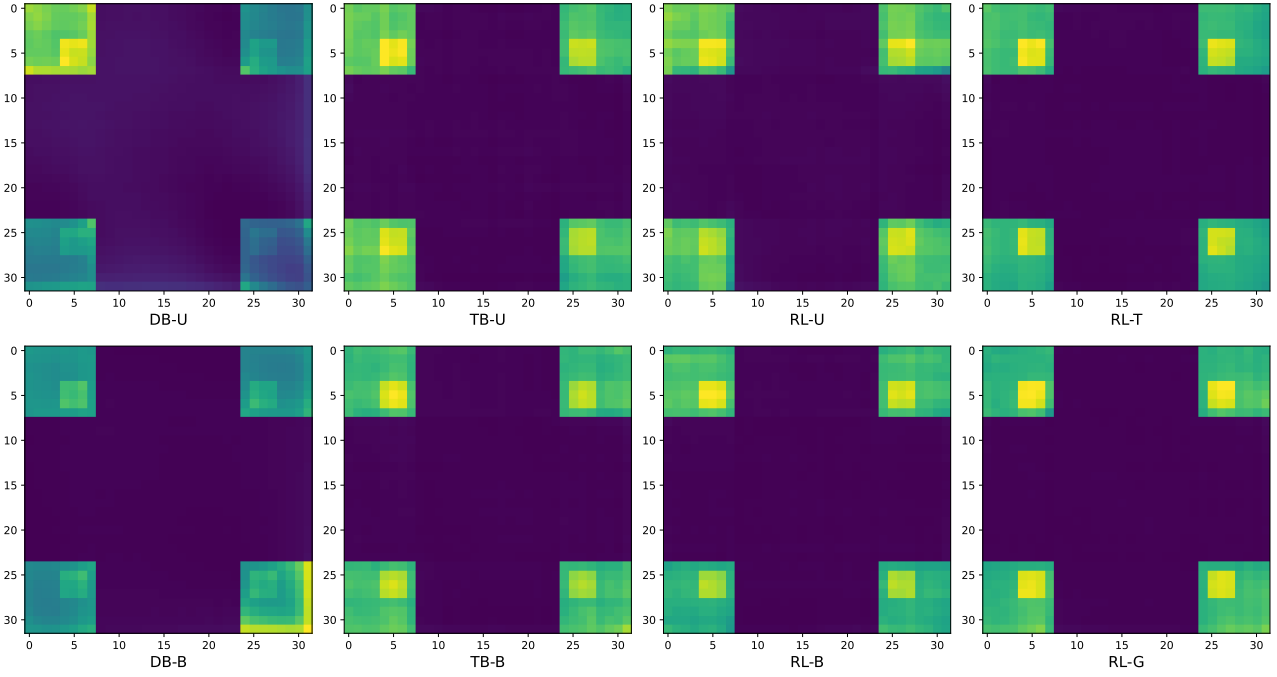


Figure 13. Graphical illustrations of $P_F^T(x)$ averaged across 5 runs of different training methods for a $32 \times 32 \times 32 \times 32$ hyper-grid. For visualization easiness, only the marginals of 2 dimensions are plotted.

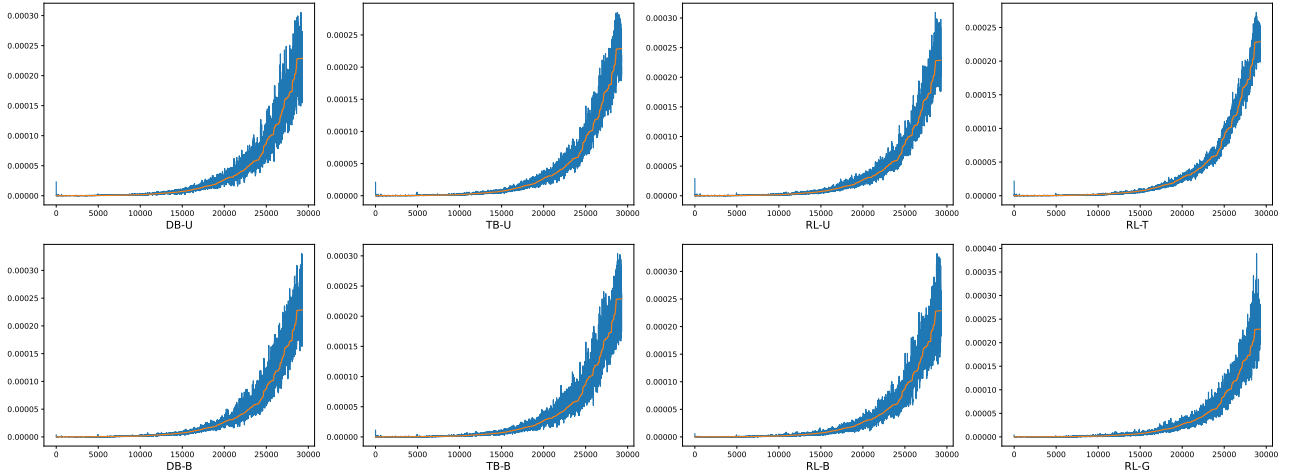


Figure 14. In each plot, the blue lines are the graphical illustrations of $P_F^T(x)$ averaged across five runs of different training methods for the BN structure learning experiment. The orange line is the ground-truth distribution and its values are plotted in an increasing order.