
Convergences guarantees of GFlowNets

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Although they were introduced to approximate complex distributions defined up
2 to normalization, Generative Flow Networks (GFlowNets) only provide strong
3 guarantees once idealized conditions are matched. However, these conditions
4 are never satisfied exactly in practice when they are trained using gradient-based
5 methods. In this paper, we prove that minimizing the Trajectory Balance loss, a
6 popular GFlowNet objective, does lead to an induced distribution getting closer
7 to the target distribution of interest, confirming theoretically this long-standing
8 intuition from the GFlowNet literature. We ultimately show that the KL divergence
9 between both distributions is upper-bounded by the quantity being minimized, and
10 we further verify this theoretical statement on a simple sampling task.

11 1 Introduction

12 Unlike more standard variational methods that minimize some divergence between an approximation
13 and the target distribution of interest (typically, the reverse KL-divergence [Jordan et al., 1999]),
14 *Generative Flow Networks* [GFlowNets; Bengio et al., 2021, 2023] instead minimize an objective
15 quantifying “inconsistencies” with some conditions which, if they were to be satisfied, would ensure
16 that we can sample from the target distribution. This provides strong guarantees in the idealized
17 setting where those conditions are fully satisfied, but those are never met in practice since models are
18 inevitably bound by their finite capacity and optimization. In this paper, we show that the divergence
19 between the distribution induced by the GFlowNet and the target distribution can be theoretically
20 bounded, confirming the intuition in this community that minimizing losses derived from the ideal
21 conditions does yield better approximations. This sheds light on the behavior of GFlowNets in those
22 inexact settings encountered in practice. We also validate these theoretical results empirically on
23 some basic generative tasks.

24 2 Background

25 We first recall some basic notions of Generative Flow Networks (GFlowNets), and refer readers
26 to [Bengio et al., 2023] for a more comprehensive treatment. GFlowNets were introduced as a
27 framework for sampling from complex distributions over combinatorial objects [Bengio et al., 2021],
28 where the sample space is often so large that computing the normalizing constant (partition function)
29 is intractable. They cast the generation process as a sequential decision-making problem in a directed
30 acyclic graph (DAG) of states $\bar{\mathcal{S}} = \mathcal{S} \cup \perp$, starting from a unique source s_0 and always terminating
31 at a designated sink \perp . The objects of interest are the terminating states $x \in \mathcal{X} \subseteq \mathcal{S}$, defined as the
32 states encountered immediately before reaching \perp .

33 A GFlowNet is parameterized by forward transition probabilities $P_F(s' | s)$, which specify how
34 states are chosen along trajectories from s_0 to \perp . This induces a distribution over terminating states,

35 called the *terminating state distribution*:

$$P_F^\top(x) \triangleq \sum_{\tau: s_0 \rightsquigarrow x} P_F(\tau) = \sum_{\tau: s_0 \rightsquigarrow x} \prod_{t=0}^{T_\tau} P_F(s_{t+1} | s_t), \quad (1)$$

36 where $s_0 \rightsquigarrow x$ denotes all the trajectories terminating at x , with $s_{T_\tau} = x$ and $s_{T_\tau+1} = \perp$. Moreover,
 37 we also define a positive reward function $R(x)$, quantifying a notion of preference over terminating
 38 states x . The objective of a GFlowNet is to find a forward transition probability such that $P_F^\top(x)$
 39 matches the target distribution:

$$P^*(x) = \frac{R(x)}{Z^*}, \quad (2)$$

40 Among the various conditions introduced in the GFlowNet literature to characterize this P_F [Bengio
 41 et al., 2021, 2023, Madan et al., 2023, Zhang et al., 2023], the *Trajectory Balance* (TB) conditions
 42 [Malkin et al., 2022] stand out as being the most widely used in practice. It states that if there exists a
 43 positive scalar $Z > 0$ and a backward transition probability P_B (i.e., $P_B(\cdot | s)$ is a distribution over
 44 the parents of s) such that for all trajectories $\tau = (s_0, s_1, \dots, s_T, \perp)$,

$$Z \prod_{t=0}^T P_F(s_{t+1} | s_t) = R(s_T) \prod_{t=1}^T P_B(s_{t-1} | s_t), \quad (3)$$

45 then the terminating state distribution $P_F^\top(x) \propto R(x)$, and $Z \equiv Z^*$. To learn this forward transition
 46 probability, the TB loss has been introduced as a way to quantify the ‘‘mismatch’’ in these conditions.
 47 It is a non-linear least-square objective of the form $\mathcal{L}_{\text{TB}}(\phi, \psi) = \frac{1}{2} \mathbb{E}_{\pi_b} [\Delta_{\text{TB}}^2(\tau; \phi)]$, where π_b is an
 48 arbitrary distribution over trajectories τ , and the residual is

$$\Delta_{\text{TB}}(\tau; \phi) = \log \frac{Z_\phi \prod_{t=0}^T P_F^\phi(s_{t+1} | s_t)}{R(s_T) \prod_{t=1}^T P_B^\phi(s_{t-1} | s_t)}, \quad (4)$$

49 where Z_ϕ and P_B^ϕ are learned alongside the forward transition probability P_F^ϕ . It is clear that if this
 50 loss reaches its global minimum (leaving aside discussions of expressivity of P_F^ϕ), then the residuals
 51 are all zeros and (3) are satisfied, leading to the guarantees about $P_F^{\phi^\top}(x)$.

52 Related to our work, Malkin et al. [2023] showed that the KL divergence between $P_F^{\phi^\top}(x)$ and the
 53 target distribution can be bounded using the following data processing inequality [Zhang et al., 2019]

$$\text{KL}(P_F^{\phi^\top}(x) \| P^*(x)) \leq \text{KL}(P_F^\phi(\tau) \| P_B^\phi(\tau)), \quad (5)$$

54 with appropriately defined distributions over trajectories $P_F^\phi(\tau)$ and $P_B^\phi(\tau)$. However, the latter still
 55 remains intractable, since it depends on $P^*(x)$ itself.

56 3 Convergence guarantees of GFlowNets

57 The approximation of (3) by minimizing the TB loss inevitably leads to errors, the source of which is
 58 typically two-fold: (1) due to the choice of the variational family to parametrize P_F^ϕ , which may not
 59 be expressive enough, and (2) due to the finite nature of optimization. In those cases, that are always
 60 encountered in practice, we have no guarantee as to what $P_F^{\phi^\top}(x)$ could be, or how close to the target
 61 distribution $P^*(x)$ it is. The objective of this section is to build towards establishing a connection
 62 between the mismatch quantified by (4) and the divergence between these two distributions. All the
 63 proofs of the propositions in this section are available in Appendix A.

64 3.1 Estimation of the partition function

65 One advantage of TB compared to other GFlowNets conditions is that it provides an approximation
 66 of the partition function with the learned scalar Z_ϕ , in addition to an approximate sampler for P^* .
 67 However, it has been shown empirically that Z_ϕ alone may yield an unreliable estimation of the
 68 partition function [Malkin et al., 2022], which may be detrimental if used for model comparison.
 69 The following proposition offers a more precise relation between the true partition function Z^* , its
 70 learned counterpart Z_ϕ , and the residuals (4).

71 **Proposition 1** (Estimation of the partition function). *The log-partition function $\log Z^*$ of the target*
 72 *distribution (2) is related to $\log Z_\phi$ via*

$$\log Z^* = \log Z_\phi + \log \left(\mathbb{E}_{\tau \sim P_F^\phi} \left[\exp(-\Delta_{\text{TB}}(\tau; \phi)) \right] \right). \quad (6)$$

73 This result is reminiscent of the estimation of the log-marginal likelihood via a variational lower-
 74 bound in [Zhou et al., 2024], which can be obtained by applying Jensen’s inequality to (6). This
 75 proposition plays a key role in proving the results of the next section.

76 3.2 Convergence of the terminating state distribution

77 Using a similar argument as the one used to prove Proposition 1, we can show that the difference
 78 between the log-terminating state probability $\log P_F^{\phi\top}(x)$ and the target log-probability $\log P^*(x)$
 79 is also directly controlled by the residual (4) being minimized. This is the first evidence supporting
 80 the general wisdom in the GFlowNet literature that minimizing the TB loss leads to a more accurate
 81 approximation at a particular x .

82 **Proposition 2.** *For any $x \in \mathcal{X}$, the error between the log-terminating state probability associated*
 83 *with P_F^ϕ and the target log-probability is bounded by*

$$|\log P_F^{\phi\top}(x) - \log P^*(x)| \leq \max_{\tau \in \mathcal{T}} |\Delta_{\text{TB}}(\tau; \phi)| + \max_{\tau: s_0 \rightsquigarrow x} |\Delta_{\text{TB}}(\tau; \phi)|. \quad (7)$$

84 We used \mathcal{T} to denote the set of all trajectories. As an immediate consequence of this proposition, we
 85 can get a simpler (albeit looser, because it becomes independent of the terminating state x) bound

$$|\log P_F^{\phi\top}(x) - \log P^*(x)| \leq 2 \max_{\tau \in \mathcal{T}} |\Delta_{\text{TB}}(\tau; \phi)|. \quad (8)$$

86 The first term in the bound of Proposition 2, which is independent of x , has an intuitive interpretation.
 87 We could allocate all our training capacity to learn parameters ϕ so that the TB conditions match
 88 almost perfectly for all the trajectories leading to a certain terminating state x only (e.g., using a
 89 behavior policy that focuses almost exclusively on those trajectories). In that case, we would have
 90 $\max_{\tau: s_0 \rightsquigarrow x} |\Delta_{\text{TB}}(\tau; \phi)| \approx 0$. However, that does not mean that we perfectly recovered $\log P^*(x)$,
 91 because we have to take into account all the other terminating states as well in the computation of its
 92 normalization constant; this is precisely what $\max_{\tau \in \mathcal{T}} |\Delta_{\text{TB}}(\tau; \phi)|$ controls for.

93 Instead of considering what happens at the level of a single terminating state x , we can also get
 94 more global guarantees on the divergence between these two distributions, similar to what typical
 95 variational methods would minimize. The following proposition gives a bound on their KL-divergence
 96 that still depends exclusively on $\Delta_{\text{TB}}(\tau; \phi)$.

97 **Proposition 3.** *Let ϕ be the parameters of the forward transition probabilities P_F^ϕ , the backward*
 98 *transition probabilities P_B^ϕ , and the total flow Z_ϕ in the trajectory balance loss (4). The KL-divergence*
 99 *between the terminating state probability distribution associated with P_F^ϕ and the target distribution*
 100 *P^* is bounded by*

$$\text{KL}(P_F^{\phi\top}(x) \| P^*(x)) \leq \mathbb{E}_{x \sim P_F^{\phi\top}} \left[\max_{\tau: s_0 \rightsquigarrow x} \Delta_{\text{TB}}(\tau; \phi) \right] - \min_{\tau \in \mathcal{T}} \Delta_{\text{TB}}(\tau; \phi). \quad (9)$$

101 We should emphasize that since the bounds of Propositions 2 & 3 both involve the maximum error
 102 being incurred for any complete trajectory, they are often impractical for quantifying how close the
 103 distribution induced by the GFlowNet is to the target. Indeed, the number of trajectories is often
 104 much larger than the number of states, which itself is already combinatorially large. Another caveat is
 105 that the KL divergence in (9) does not measure if $P_F^{\phi\top}$ misses some of the modes of P^* , potentially
 106 only giving a partial view on the accuracy of the approximation.

107 Nevertheless, the goal of these results is to confirm the intuition that the approximation found by the
 108 GFlowNet gets better as the TB loss is minimized, which was missing in the community as little was
 109 known about the behavior of an “inexact” GFlowNet that doesn’t satisfy the TB conditions exactly.
 110 In Appendix C, we show that this is not limited to the TB loss, and a similar bound can be derived for
 111 another GFlowNet objective called the *Detailed Balance* (DB) loss [Bengio et al., 2023].

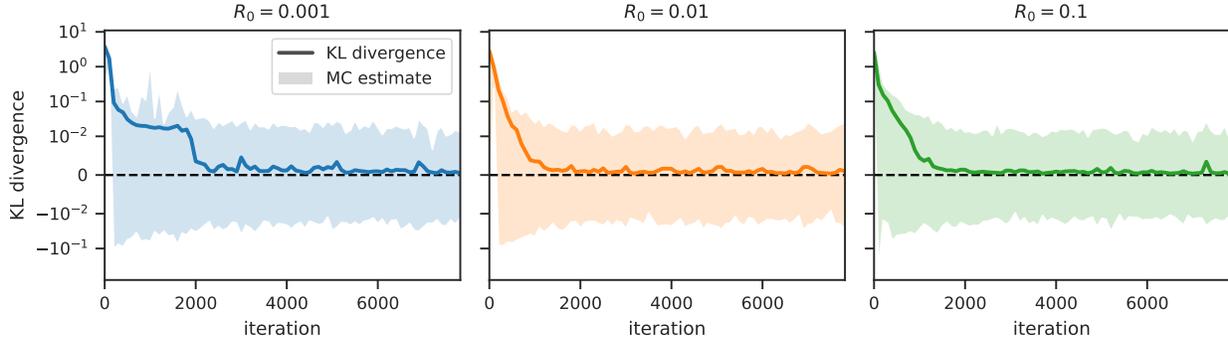


Figure 1: Empirical validation of Proposition 3 on the $8 \times 8 \times 8 \times 8$ hypergrid, with varying difficulty (controlled by R_0), during training. The true $\text{KL}(P_F^{\phi^\top}(x) \parallel P^*(x))$ eventually decreases to 0, and always fall in the confidence interval given by estimates of (10) (shaded area). The confidence interval is obtained with bootstrapping. See Figure 2 for similar results on the 64×64 hypergrid.

112 4 Empirical validation

113 The bounds established in Propositions 2 & 3 both depend on the maximum error across all complete
 114 trajectories $\tau \in \mathcal{T}$. This dependence on worst-case trajectories makes them less practical for
 115 accurately assessing how closely the terminating state distribution approximates the target, even in
 116 relatively simple generative tasks. Instead, these results serve primarily to confirm the intuition that
 117 minimizing the TB loss moves the learned distribution closer to the target.

118 To provide a more practical validation of the theoretical results from the previous section, we turn to
 119 the following equality derived in the proof of Proposition 3, which connects the KL divergence to the
 120 residuals (the proposition being simply an upper bound of the RHS):

$$\begin{aligned} \text{KL}(P_F^{\phi^\top}(x) \parallel P^*(x)) &= \log \left(\mathbb{E}_{\tau \sim P_F^\phi} [\exp(-\Delta_{\text{TB}}(\tau; \phi))] \right) \\ &\quad + \mathbb{E}_{x \sim P_F^{\phi^\top}} \left[\log \left(\mathbb{E}_{\tau \sim P_B^\phi(\cdot | x)} [\exp(\Delta_{\text{TB}}(\tau; \phi))] \right) \right]. \end{aligned} \quad (10)$$

121 Here, $\tau \sim P_B^\phi(\cdot | x)$ denotes trajectories sampled according to the backward transition probability
 122 P_B^ϕ , starting from the terminating state x . To corroborate the theoretical result of Proposition 3, we
 123 monitor whether the true KL divergence computed throughout training falls within the confidence
 124 interval induced by Monte Carlo estimates of (10) (see Algorithm 1 for details), in cases where both
 125 $P_F^{\phi^\top}(x)$ and $P^*(x)$ can be computed analytically.

126 We consider the hypergrid task introduced by Bengio et al. [2021], which is a simple and widely
 127 studied environment in the GFlowNet literature. In this environment, the state space is organized
 128 as an N -dimensional grid where all the states are terminating. The initial state is $(0, \dots, 0)$, and
 129 transitions correspond to incrementing one of the coordinates by one, up to a maximum size H . The
 130 reward function in defined is such a way that it peaks around the corners of the hypergrid, with valleys
 131 in between controlled by a parameter R_0 determining the relative difficulty of finding all the modes
 132 of the target distribution (hence, the difficulty of approximating it); see Appendix B for details. Both
 133 $P^*(x)$ and $P_F^{\phi^\top}(x)$ can be evaluated in closed form here (the former by exhaustive enumeration of
 134 the terminating states and their rewards, the latter using dynamic programming [Malkin et al., 2023]).

135 In Figure 1, we show the evolution of $\text{KL}(P_F^{\phi^\top}(x) \parallel P^*(x))$ throughout training on a 4-dimensional
 136 grid of size 8, with varying levels of difficulty, along with the confidence intervals derived from esti-
 137 mating (10). We observe that the true KL divergence goes down to 0, showing that the approximation
 138 found by the GFlowNet gets better, and this divergence stays steadily within the confidence interval
 139 during training, confirming our theoretical result in Proposition 3. Note that even on this small task,
 140 enumerating exhaustively all the trajectories in \mathcal{T} is infeasible, making the exact evaluation of the
 141 bound (9) impossible.

142 **References**

- 143 E. Bengio, M. Jain, M. Korablyov, D. Precup, and Y. Bengio. Flow Network based Generative Models
144 for Non-Iterative Diverse Candidate Generation. *Advances in Neural Information Processing*
145 *Systems (NeurIPS)*, 2021.
- 146 Y. Bengio, S. Lahlou, T. Deleu, E. J. Hu, M. Tiwari, and E. Bengio. GFlowNet Foundations. *Journal*
147 *of Machine Learning Research (JMLR)*, 2023.
- 148 M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An Introduction to Variational Methods
149 for Graphical Models. *Machine learning*, 1999.
- 150 K. Madan, J. Rector-Brooks, M. Korablyov, E. Bengio, M. Jain, A. Nica, T. Bosc, Y. Bengio, and
151 N. Malkin. Learning GFlowNets from partial episodes for improved convergence and stability.
152 *International Conference on Machine Learning (ICML)*, 2023.
- 153 N. Malkin, M. Jain, E. Bengio, C. Sun, and Y. Bengio. Trajectory Balance: Improved Credit
154 Assignment in GFlowNets. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 155 N. Malkin, S. Lahlou, T. Deleu, X. Ji, E. Hu, K. Everett, D. Zhang, and Y. Bengio. GFlowNets and
156 variational inference. *International Conference on Learning Representations (ICLR)*, 2023.
- 157 D. Zhang, C. Rainone, M. Peschl, and R. Bondesan. Robust Scheduling with GFlowNets. *Interna-*
158 *tional Conference on Learning Representations (ICLR)*, 2023.
- 159 M. Zhang, T. Bird, R. Habib, T. Xu, and D. Barber. Variational f-divergence Minimization. *Informa-*
160 *tion Theory and Machine Learning Workshop (NeurIPS)*, 2019.
- 161 M. Zhou, Z. Yan, E. Layne, N. Malkin, D. Zhang, M. Jain, M. Blanchette, and Y. Bengio. PhyloGFN:
162 Phylogenetic inference with generative flow networks. *International Conference on Learning*
163 *Representations (ICLR)*, 2024.

164 Appendix

165 A Proofs of the convergence guarantees for the TB loss

166 Before proving the proposition from the main text, we will first state a lemma that will be useful in
167 the proofs later

168 **Lemma 1** (Bound on log-expectation-exp). *Let $\mathbf{p} = (p_1, \dots, p_n)$ be a vector of probabilities*
169 *(i.e., $p_i \geq 0$ and $\sum_i p_i = 1$), and $\mathbf{x} = (x_1, \dots, x_n)$ be an arbitrary vector. Then*

$$\left| \log \sum_{i=1}^n p_i \exp(x_i) \right| \leq \max_i |x_i|. \quad (11)$$

170 *Proof.* For any vector $\mathbf{y} = (y_1, \dots, y_n)$, and since $p_i \geq 0$, we have the following inequalities

$$\sum_{i=1}^n p_i \left(\min_j y_j \right) \leq \sum_{i=1}^n p_i y_i \leq \sum_{i=1}^n p_i \left(\max_j y_j \right). \quad (12)$$

171 Since $\sum_i p_i = 1$, both sides of these inequalities can be further simplified, only involving the
172 minimum and maximum of \mathbf{y} . We can apply these inequalities to $y_i = \exp(x_i)$, and observe that
173 $\min_j \exp(x_j) = \exp(\min_j x_j)$ (and similarly for the max), because the exponential is a monotoni-
174 cally increasing function

$$\exp \left(\min_j x_j \right) \leq \sum_{i=1}^n p_i \exp(x_i) \leq \exp \left(\max_j x_j \right). \quad (13)$$

175 Taking the logarithm of the inequalities above, and again using the fact that log is monotonically
176 increasing

$$\min_j x_j \leq \log \sum_{i=1}^n p_i \exp(x_i) \leq \max_j x_j. \quad (14)$$

177 Another way to write (14) is

$$\log \sum_{i=1}^n p_i \exp(x_i) \leq \max_i x_i \leq \max_i |x_i| \quad (15)$$

$$-\log \sum_{i=1}^n p_i \exp(x_i) \leq -\min_i x_i \leq \max_i |x_i|, \quad (16)$$

178 which concludes the proof. \square

179 **Proposition 1** (Estimation of the partition function). *The log-partition function $\log Z^*$ of the target*
180 *distribution (2) is related to $\log Z_\phi$ via*

$$\log Z^* = \log Z_\phi + \log \left(\mathbb{E}_{\tau \sim P_F^\phi} \left[\exp(-\Delta_{\text{TB}}(\tau; \phi)) \right] \right). \quad (6)$$

181 *Proof.* Based on (4), for some complete trajectory $\tau = (s_0, s_1, \dots, s_T, \perp)$, we can find an expression
182 for the product of the reward at the terminating state s_T and the backward probability of τ as a
183 function of the residual $\Delta_{\text{TB}}(\tau; \phi)$

$$R(s_T) P_B^\phi(\tau | s_T) = R(s_T) \prod_{t=1}^T P_B^\phi(s_{t-1} | s_t) = Z_\phi \exp(-\Delta_{\text{TB}}(\tau; \phi)) \prod_{t=0}^T P_F^\phi(s_{t+1} | s_t), \quad (17)$$

184 where we used the notation $P_B^\phi(\tau | s_T) = \prod_{t=1}^T P_B^\phi(s_{t-1} | s_t)$. Since $P_B^\phi(\cdot | x)$ is a distribution
185 over the complete trajectories terminating at x [Bengio et al., 2023], the reward $R(x)$ is given by

$$R(x) = \sum_{\tau: s_0 \rightsquigarrow x} R(x) P_B^\phi(\tau | x) = Z_\phi \sum_{\tau: s_0 \rightsquigarrow x} P_F^\phi(\tau) \exp(-\Delta_{\text{TB}}(\tau; \phi)) \quad (18)$$

186 By definition of the partition function Z , and using the fact that P_F^ϕ induces a probability distribution
 187 over all the complete trajectories [Bengio et al., 2023]

$$Z^* = \sum_{x \in \mathcal{X}} R(x) \quad (19)$$

$$= Z_\phi \sum_{x \in \mathcal{X}} \sum_{\tau: s_0 \rightsquigarrow x} P_F^\phi(\tau) \exp(-\Delta_{\text{TB}}(\tau; \phi)) \quad (20)$$

$$= Z_\phi \sum_{\tau \in \mathcal{T}} P_F^\phi(\tau) \exp(-\Delta_{\text{TB}}(\tau; \phi)) \quad (21)$$

$$= Z_\phi \mathbb{E}_{\tau \sim P_F^\phi} [\exp(-\Delta_{\text{TB}}(\tau; \phi))]. \quad (22)$$

188 We can conclude by taking the log of the equality above. \square

189 **Proposition 2.** For any $x \in \mathcal{X}$, the error between the log-terminating state probability associated
 190 with P_F^ϕ and the target log-probability is bounded by

$$|\log P_F^{\phi^\top}(x) - \log P^*(x)| \leq \max_{\tau \in \mathcal{T}} |\Delta_{\text{TB}}(\tau; \phi)| + \max_{\tau: s_0 \rightsquigarrow x} |\Delta_{\text{TB}}(\tau; \phi)|. \quad (7)$$

191 *Proof.* We can use (4) to find an expression for $P_F^\phi(\tau)$ as a function of the residual $\Delta_{\text{TB}}(\tau; \phi)$:

$$\prod_{t=0}^T P_F^\phi(s_{t+1} | s_t) = \exp(\Delta_{\text{TB}}(\tau; \phi)) \frac{R(s_T)}{Z_\phi} \prod_{t=1}^T P_B^\phi(s_{t-1} | s_t). \quad (23)$$

192 For any $x \in \mathcal{X}$, using the notation $P_B^\phi(\tau | x) = \prod_{t=1}^T P_B^\phi(s_{t-1} | s_t)$ (with $s_T = x$), which is
 193 a properly defined distribution over the complete trajectories terminating in x , and by (1) of the
 194 terminating state probability distribution associated with P_F^ϕ , we get

$$P_F^{\phi^\top}(x) = \sum_{\tau: s_0 \rightsquigarrow x} \prod_{t=0}^{T_\tau} P_F^\phi(s_{t+1} | s_t) \quad (24)$$

$$= \frac{R(x)}{Z_\phi} \sum_{\tau: s_0 \rightsquigarrow x} \exp(\Delta_{\text{TB}}(\tau; \phi)) P_B^\phi(\tau | x) \quad (25)$$

$$= \frac{R(x)}{Z^*} \frac{Z^*}{Z_\phi} \mathbb{E}_{\tau \sim P_B^\phi(\cdot | x)} [\exp(\Delta_{\text{TB}}(\tau; \phi))] \quad (26)$$

$$= P^*(x) \frac{Z^*}{Z_\phi} \mathbb{E}_{\tau \sim P_B^\phi(\cdot | x)} [\exp(\Delta_{\text{TB}}(\tau; \phi))] \quad (27)$$

195 Although the normalization constant Z^* of the target distribution is still unknown, we can fortunately
 196 write it as a function of the residual as well thanks to Proposition 1. Combining it with (27), we
 197 get an expression of the difference in log-probabilities as a function of $\Delta_{\text{TB}}(\tau; \phi)$ only, where the
 198 expectations over complete trajectories are respectively taken wrt. P_F^ϕ and $P_B^\phi(\cdot | x)$

$$\log P_F^{\phi^\top}(x) - \log P^*(x) = \log \left(\mathbb{E}_{\tau \sim P_F^\phi} [\exp(-\Delta_{\text{TB}}(\tau; \phi))] \right) + \log \left(\mathbb{E}_{\tau \sim P_B^\phi(\cdot | x)} [\exp(\Delta_{\text{TB}}(\tau; \phi))] \right). \quad (28)$$

199 Using the triangle inequality, we can conclude by applying Lemma 1 to both terms of the RHS

$$\left| \log \left(\mathbb{E}_{\tau \sim P_F^\phi} [\exp(-\Delta_{\text{TB}}(\tau; \phi))] \right) \right| \leq \max_{\tau \in \mathcal{T}} |\Delta_{\text{TB}}(\tau; \phi)| \quad (29)$$

$$\left| \log \left(\mathbb{E}_{\tau \sim P_B^\phi(\cdot | x)} [\exp(\Delta_{\text{TB}}(\tau; \phi))] \right) \right| \leq \max_{\tau: s_0 \rightsquigarrow x} |\Delta_{\text{TB}}(\tau; \phi)|. \quad (30)$$

200 \square

201 **Proposition 3.** Let ϕ be the parameters of the forward transition probabilities P_F^ϕ , the backward
 202 transition probabilities P_B^ϕ , and the total flow Z_ϕ in the trajectory balance loss (4). The KL-divergence
 203 between the terminating state probability distribution associated with P_F^ϕ and the target distribution
 204 P^* is bounded by

$$\text{KL}(P_F^{\phi^\top}(x) \| P^*(x)) \leq \mathbb{E}_{x \sim P_F^{\phi^\top}} \left[\max_{\tau: s_0 \rightsquigarrow x} \Delta_{\text{TB}}(\tau; \phi) \right] - \min_{\tau \in \mathcal{T}} \Delta_{\text{TB}}(\tau; \phi). \quad (9)$$

Algorithm 1 KL-divergence estimation

- 1: **for** $m = 1, \dots, M$ **do**
- 2: Sample a trajectory τ_m following P_F^ϕ , terminating at state x_m
- 3: Starting at x_m , sample a trajectory τ'_m following P_B^ϕ
- 4: Construct a Monte Carlo estimate of (10)

$$\widehat{\text{KL}} = \log \left(\frac{1}{M} \sum_{m=1}^M \exp(-\Delta_{\text{TB}}(\tau_m; \phi)) \right) + \frac{1}{M} \sum_{m=1}^M \Delta_{\text{TB}}(\tau'_m; \phi)$$

- 5: **return** the estimate $\widehat{\text{KL}}$
-

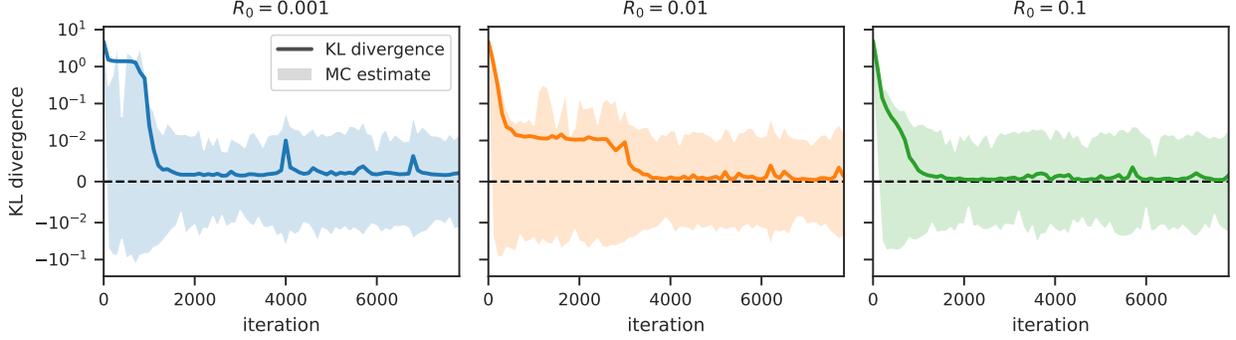


Figure 2: Empirical validation of Proposition 3 on the 64×64 hypergrid, with varying difficulty (controlled by R_0), during training. The true $\text{KL}(P_F^{\phi^\top}(x) \| P^*(x))$ eventually decreases to 0, and always fall in the confidence interval given by estimates of (10) (shaded area). The confidence interval is obtained with bootstrapping.

205 *Proof.* This is a direct consequence of (28) in the proof of Proposition 2 above:

$$\text{KL}(P_F^{\phi^\top}(x) \| P^*(x)) = \mathbb{E}_{x \sim P_F^{\phi^\top}} [\log P_F^{\phi^\top}(x) - \log P^*(x)] \quad (31)$$

$$= \log \left(\mathbb{E}_{\tau \sim P_F^\phi} [\exp(-\Delta_{\text{TB}}(\tau; \phi))] \right) + \mathbb{E}_{x \sim P_F^{\phi^\top}} \left[\log \left(\mathbb{E}_{\tau \sim P_B^\phi(\cdot|x)} [\exp(\Delta_{\text{TB}}(\tau; \phi))] \right) \right] \quad (32)$$

$$\leq \left[\max_{\tau \in \mathcal{T}} -\Delta_{\text{TB}}(\tau; \phi) \right] + \mathbb{E}_{x \sim P_F^{\phi^\top}} \left[\max_{\tau: s_0 \rightsquigarrow x} \Delta_{\text{TB}}(\tau; \phi) \right] \quad (33)$$

$$= \mathbb{E}_{x \sim P_F^{\phi^\top}} \left[\max_{\tau: s_0 \rightsquigarrow x} \Delta_{\text{TB}}(\tau; \phi) \right] - \min_{\tau \in \mathcal{T}} \Delta_{\text{TB}}(\tau; \phi), \quad (34)$$

206 where the inequality in (33) is a consequence of (14) in the proof of Lemma 1. \square

207 B Details of the empirical evaluation

208 The reward function of the hypergrid environment is defined as

$$R(x) = R_0 + \frac{1}{2} \prod_{n=1}^N \mathbf{1} \left(\left| \frac{x_n}{H-1} - 0.5 \right| \in (0.25, 0.5] \right) + 2 \prod_{n=1}^N \mathbf{1} \left(\left| \frac{x_n}{H-1} - 0.5 \right| \in (0.3, 0.4) \right) \quad (35)$$

209 In Algorithm 1 we describe a way to obtain a Monte Carlo estimate of the KL divergence using (10).
 210 In Figure 2, we show a similar experiment on the 2-dimensional hypergrid of size 64; the conclusions
 211 we made in the main text largely hold in this new setting.

212 C Proofs of the convergence guarantees for the DB loss

213 Bengio et al. [2023] showed that if P_F , P_B , and a “flow” function $F(s)$ defined over the all the states
 214 in the state space satisfy the following *Detailed Balance* (DB) conditions at every transition $s \rightarrow s'$
 215 where $s' \neq \perp$

$$F(s)P_F(s' | s) = F(s')P_B(s | s'), \quad (36)$$

216 with the boundary condition $F(x)P_F(\perp | x) = R(x)$ at any terminating state $x \in \mathcal{X}$, then $P_F^\top(x) \propto$
 217 $R(x)$. Taking inspiration from Section 2, we can convert this condition into a loss that can be
 218 minimized of the form $\mathcal{L}_{\text{DB}}(\phi) = \frac{1}{2}\mathbb{E}_{\pi_b}[\Delta_{\text{DB}}^2(s, s'; \phi)]$, where

$$\Delta_{\text{DB}}(s, s'; \phi) = \log \frac{F_\phi(s)P_F^\phi(s' | s)}{F_\phi(s')P_B^\phi(s | s')} \quad \Delta_{\text{DB}}(x, \perp; \phi) = \log \frac{F_\phi(x)P_F^\phi(\perp | x)}{R(x)} \quad (37)$$

219 Throughout this section, for a trajectory $\tau = (s_0, s_1, \dots, s_T, \perp)$, we will use the notations $s_{T+1} = \perp$,
 220 so that we can write $\Delta_{\text{DB}}(s_T, s_{T+1}; \phi)$. The following proposition is the counterpart of Proposition 2
 221 in the case of Detailed Balance.

222 **Proposition 4.** For any $x \in \mathcal{X}$, the error between the log-terminating state probability associated
 223 with P_F^ϕ and the target log-probability is bounded by

$$|\log P_F^{\phi\top}(x) - \log P^*(x)| \leq \max_{\tau \in \mathcal{T}} \left| \sum_{t=0}^{T_\tau} \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right| + \max_{\tau: s_0 \rightsquigarrow x} \left| \sum_{t=0}^{T_\tau} \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right|. \quad (38)$$

224 *Proof.* The proof is similar to the proof of Proposition 2. From (37), we can derive that

$$P_F^\phi(s' | s) = P_B^\phi(s | s') \frac{F_\phi(s')}{F_\phi(s)} \exp(\Delta_{\text{DB}}(s, s'; \phi)) \quad (39)$$

$$P_F^\phi(s, \perp) = \frac{R(x)}{F_\phi(x)} \exp(\Delta_{\text{DB}}(x, \perp; \phi)) \quad (40)$$

225 For some trajectory $\tau = (s_0, s_1, \dots, s_T, \perp)$,

$$P_F^\phi(\tau) = \prod_{t=0}^T P_F^\phi(s_{t+1} | s_t) = P_F^\phi(\perp | s_T) \prod_{t=0}^{T-1} P_F^\phi(s_{t+1} | s_t) \quad (41)$$

$$= R(s_T) \frac{1}{F_\phi(s_T)} \exp(\Delta_{\text{DB}}(x, \perp; \phi)) \prod_{t=0}^{T-1} P_B^\phi(s_t | s_{t+1}) \frac{F_\phi(s_{t+1})}{F_\phi(s_t)} \exp(\Delta_{\text{DB}}(s_t, s_{t+1}; \phi)) \quad (42)$$

$$= \frac{R(s_T)}{F_\phi(s_0)} \exp \left[\sum_{t=0}^T \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right] P_B^\phi(\tau | s_T) \quad (43)$$

226 We can also prove a similar result as Proposition 1, and show that we can write the true partition
 227 function Z^* as a function of $F_\phi(s_0)$ and the residuals:

$$Z^* = \sum_x R(x) = F_\phi(s_0) \mathbb{E}_{\tau \sim P_F^\phi} \left[\exp \left[- \sum_{t=0}^{T_\tau} \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right] \right] \quad (44)$$

228 Using (43) & (44), we can write the terminating state distribution as

$$P_F^{\phi\top}(x) = \sum_{\tau: s_0 \rightsquigarrow x} \frac{R(x)}{F_\phi(s_0)} P_B^\phi(\tau | x) \exp \left[\sum_{t=0}^{T_\tau} \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right] \quad (45)$$

$$= \frac{R(x)}{Z^*} \frac{Z^*}{F_\phi(s_0)} \mathbb{E}_{\tau \sim P_B^\phi(\cdot | x)} \left[\exp \left[\sum_{t=0}^{T_\tau} \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right] \right] \quad (46)$$

229 Writing this in terms of the difference of log-probabilities, we obtain

$$\begin{aligned} \log P_F^{\phi^\top}(x) - \log P^*(x) &= \log \left(\mathbb{E}_{\tau \sim P_B^\phi(\cdot|x)} \left[\exp \left[\sum_{t=0}^{T_\tau} \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right] \right] \right) \\ &\quad + \log \left(\mathbb{E}_{\tau \sim P_F^\phi} \left[\exp \left[- \sum_{t=0}^{T_\tau} \Delta_{\text{DB}}(s_t, s_{t+1}; \phi) \right] \right] \right) \end{aligned} \quad (47)$$

230 We can conclude using [Lemma 1](#).

□