

Trajectory Balance: Improved Credit Assignment in GFlowNets

Nikolay Malkin^{1,2} Moksh Jain^{1,2} Emmanuel Bengio^{1,3} Chen Sun^{1,2} Yoshua Bengio^{1,2}

Abstract

Generative Flow Networks (GFlowNets) are a method for learning a stochastic policy for generating compositional objects, such as graphs or strings, from a given unnormalized density by sequences of actions, where many possible action sequences may lead to the same object. Prior temporal difference-like learning objectives for training GFlowNets, such as *flow matching* and *detailed balance*, are prone to inefficient credit propagation across action sequences, particularly in the case of long sequences. We propose a new learning objective for GFlowNets, *trajectory balance*, as a more efficient alternative to previously used objectives. We prove that any global minimizer of the trajectory balance objective can define a policy that samples exactly from the target distribution. In experiments on four distinct domains, we empirically demonstrate the benefits of the trajectory balance objective for GFlowNet convergence, diversity of generated samples, and robustness to long action sequences and large action spaces.

1. Introduction

Generative Flow Networks (GFlowNets; Bengio et al., 2021a;b) are models that exploit generalizable structure in an energy function \mathcal{E} to amortize sampling from the corresponding probability density function on a space of compositional objects \mathcal{X} , for example, graphs composed of nodes and edges. GFlowNets learn a stochastic policy that models generation of such structured objects in \mathcal{X} by producing a stochastic sequence of *actions* that incrementally modify a partial object (*state*), e.g., by adding an edge or a node to a graph, starting from a universal initial state (like an empty graph). A special ‘exit’ discrete action signals that the construction of the object $x \in \mathcal{X}$ is completed, and the policy is trained so as to make the likelihood of generat-

ing x proportional to the given unnormalized probability or reward $R(x) = e^{-\mathcal{E}(x)}$.

Like other models in deep reinforcement learning (RL; Sutton & Barto, 2018), GFlowNets are trained with a parametric policy that can be given desired inductive biases and allows generalization to states not seen in training. Natural domains for applying GFlowNets are those where exact sampling is intractable and local exploration (MCMC) methods perform poorly, but diversity of samples is desired (so local optimization is also not a good option). For example, GFlowNets have been used (Bengio et al., 2021a) to generate molecules (as graphs) by incremental addition of simple building blocks, where the non-negative reward $R(x)$ is the estimated strength of binding the constructed hypothetical molecule to some protein target: the number of possible candidates grows rapidly with the allowed molecule size, and the reward has many separated modes (Bengio et al., 2021a). Like all RL models that iteratively sample action sequences for training, GFlowNets pose the learning challenges of exploration/exploitation and credit assignment, i.e., propagation of a reward signal over a sequence of actions (van Hasselt et al., 2018; Bengio et al., 2020; Kumar et al., 2021), and the *efficiency of credit assignment and training in a GFlowNet is the focus of the present paper*.

The learning problem solved by GFlowNets also has two fundamental differences with the standard reward-maximization paradigm of RL. First, a GFlowNet aims to make the likelihood of reaching a terminating state proportional to the reward, not to concentrate it at a maximal-reward state. Thus, a GFlowNet must model the diversity in the target distribution, not only its dominant mode. Reward maximization in RL can be turned into sampling proportionally to the reward with appropriate entropy maximization regularization, if there is only one way to reach every state (Bengio et al., 2021a). The second difference with reward-maximization in RL is indeed that the GFlowNet training objectives still lead to correct sampling even when multiple action sequences lead to the same terminating state. Note that the likelihood of reaching a state is the sum of likelihoods of all action sequences leading to it, and that there may be an exponentially large number of such paths.

The set of all achievable sequences of actions and corresponding states can be organized in a directed graph

¹Mila ²Université de Montréal ³McGill University. Correspondence to: Nikolay Malkin <nikolay.malkin@mila.quebec>.

$G = (\mathcal{S}, \mathcal{A})$ in which the vertices \mathcal{S} are states (some of them designated as terminal states, in bijection with \mathcal{X}) and the edges $u \rightarrow v$ in \mathcal{A} each correspond to applying an action while in a state $u \in \mathcal{S}$ and landing in state v . Bengio et al. (2021a) describe a GFlowNet by a nonnegative function on the edges, called the *edge flow* $F : \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$, where $F(u \rightarrow v)$ is an *unnormalized* likelihood of taking the action that modifies state u to state v . The *GFlowNet policy* samples the transition $u \rightarrow v$ from state u with probability $F(u \rightarrow v) / \sum_{v'} F(u \rightarrow v')$ where the denominator sums over the outgoing edges from u . By analogy and following the classical notion of flows in networks (Ford & Fulkerson, 1956), one can think of this flow like the amount of water flowing through an edge (like a pipe) or a state (like a tee, where pipes meet). Because of a conservation law, the flow in a state should be the sum of flows into or out of it, and cycles are not permitted. It is shown that this GFlowNet policy samples x proportionally to $R(x)$ if F satisfies a set of linear *flow matching* constraints (the sum of incoming flows equals to the sum of outgoing flows), which can be converted into a temporal difference-like objective that can be optimized with respect to the parameters of a neural net that approximates F . An alternative objective based on *detailed balance* constraints was proposed in Bengio et al. (2021b). These objectives, however, like temporal-difference learning, can suffer from slow credit assignment (van Hasselt et al., 2018; Bengio et al., 2020; Kumar et al., 2021), as illustrated in Figure 1.

The main contribution of this work (§3) is a new parametrization and objective for GFlowNets. This objective, which we call *trajectory balance*, is computed on sampled full action sequences (trajectories) from the initial state to a terminal state, unlike the flow matching and detailed balance objectives. We prove that global minimization of trajectory balance implies that the learned action policy samples proportionally to R . We also empirically show that the trajectory balance objective accelerates training convergence relative to previously proposed objectives, improves the learned sampling policy with respect to metrics of diversity and divergence from the reward function, and allows learning GFlowNets that generate sequences far longer than was possible before. Comparative experimental validation is performed on four domains illustrating different features of the reward landscape:

- **Hypergrid** (§5.1), an illustrative synthetic environment with modes separated by wide troughs;
- **Molecule synthesis** (§5.2), a practical graph generation problem with reward given by a deep model;
- **Sequence generation** (§5.3), where we demonstrate robustness of trajectory balance to large action space sizes and long action sequences on synthetic data and illustrate the algorithms on real AMP sequence data.

2. Preliminaries

2.1. Markovian flows

We give some essential definitions, following §2 of Bengio et al. (2021b). Fix a directed acyclic graph $G = (\mathcal{S}, \mathcal{A})$ with state space \mathcal{S} and action space \mathcal{A} . Let $s_0 \in \mathcal{S}$ be the special *initial (source) state*, the only state with no incoming edges, and designate vertices with no outgoing edges as *terminal (sinks)*¹. We call the vertices *states*, the edges *actions*, the states reachable through outgoing edges from a state its *children*, and the sources of its incoming edges its *parents*.

A *complete trajectory* is a sequence of transitions $\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_n)$ going from the initial state s_0 to a terminal state s_n with $(s_t \rightarrow s_{t+1}) \in \mathcal{A}$ for all t . Let \mathcal{T} be the set of complete trajectories. A *trajectory flow* is a nonnegative function $F : \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}$. With our water analogy, it could be the number of water molecules travelling along this path (the units don’t matter because the flow function can be scaled arbitrarily, since we normalize them to get probabilities). For any state s , define the state flow $F(s) = \sum_{\tau \in \mathcal{T}} F(\tau)$, and, for any edge $s \rightarrow s'$, the edge flow

$$F(s \rightarrow s') = \sum_{\tau = (\dots \rightarrow s \rightarrow s' \rightarrow \dots)} F(\tau). \quad (1)$$

As a consequence of this definition, the *flow matching* constraint (incoming flow = outgoing flow) is satisfied for all states s that are not initial or terminal:

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

A nontrivial² trajectory flow F determines a distribution P over trajectories,

$$P(\tau) = \frac{1}{Z} F(\tau), \quad Z = F(s_0) = \sum_{\tau \in \mathcal{T}} F(\tau). \quad (3)$$

The trajectory flow F is *Markovian* if there exist action distributions $P_F(-|s)$ over the children of each nonterminal state s such that the distribution P has a factorization

$$P(\tau = (s_0 \rightarrow \dots \rightarrow s_n)) = \prod_{t=1}^n P_F(s_t | s_{t-1}). \quad (4)$$

Equivalently (Bengio et al. (2021b), Prop. 3) there are distributions $P_B(-|s)$ over the parents of each noninitial state s such that for any terminal x ,

$$P(\tau = (s_0 \rightarrow \dots \rightarrow s_n) | s_n = x) = \prod_{t=1}^n P_B(s_{t-1} | s_t). \quad (5)$$

¹Bengio et al. (2021a) allowed terminal states with outgoing edges. The difference is easily overcome by augmenting every such state x by a new terminal state x^\top with a stop action $x \rightarrow x^\top$.

²Here and below, ‘nontrivial’ = ‘not identically zero’.

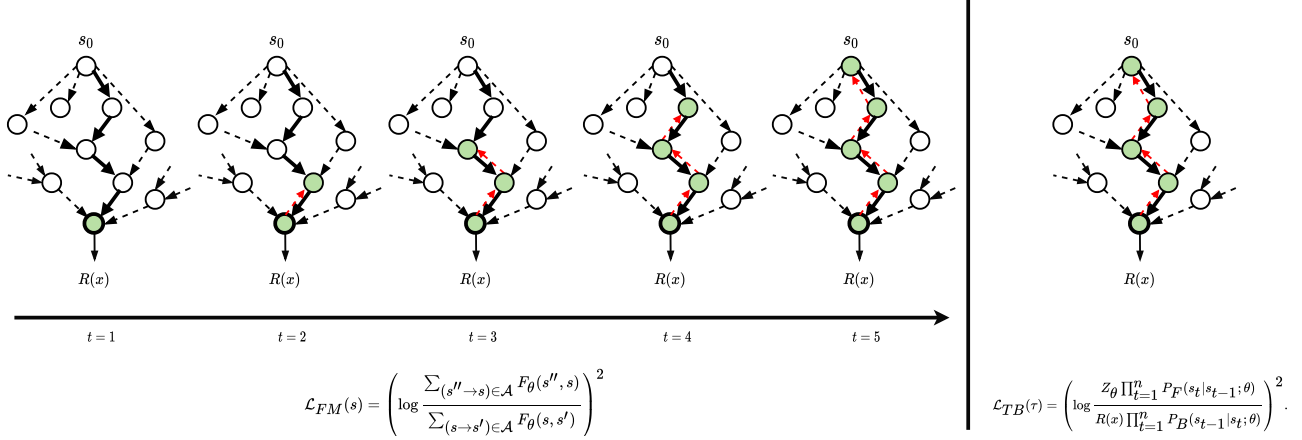


Figure 1. The hypothesized advantage in terms of speed of credit assignment using the trajectory balance (TB) loss (*right*) vs. flow matching (FM) loss (*left*). Similarly to temporal-difference methods, \mathcal{L}_{FM} slowly propagates flow matching errors backwards from the error at the terminal reward, whereas \mathcal{L}_{TB} immediately provides a stochastic training signal about the final reward to all steps of the generated action sequence. The contrast increases with longer action sequences.

If F is a Markovian flow, then P_F and P_B can be computed in terms of state and edge flows:

$$P_F(s' | s) = \frac{F(s \rightarrow s')}{F(s)}, \quad P_B(s | s') = \frac{F(s \rightarrow s')}{F(s')}, \quad (6)$$

supposing denominators do not vanish. We call P_F and P_B the *forward policy* and *backward policy* corresponding to F , respectively. These relations are summarized by the *detailed balance constraint*

$$F(s)P_F(s' | s) = F(s')P_B(s | s'). \quad (7)$$

Uniqueness properties. A Markovian flow is uniquely determined by an edge flow, i.e., a nontrivial choice of non-negative value on every edge satisfying the flow matching constraint (2). By [Bengio et al. \(2021b\)](#), Corollary 1, a Markovian flow is also uniquely determined by either of

- a positive constant $Z = F(s_0)$ and a choice of distribution $P_F(-|s)$ over children of every nonterminal state; or
- a nontrivial choice of nonnegative state flows $F(x)$ for every terminal state x and a choice of distribution $P_B(-|s)$ over parents of every noninitial state.

2.2. GFlowNets

Suppose that a nontrivial nonnegative reward function $R: \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ is given on the set of terminal states. GFlowNets ([Bengio et al., 2021a](#)) aim to approximate a Markovian flow F on G such that

$$F(x) = R(x) \quad \forall x \in \mathcal{X}. \quad (8)$$

We adopt the broad definition that a GFlowNet is any learning algorithm consisting of:

- a model capable of providing the initial state flow $Z = F(s_0)$ as well as the forward action distributions $P_F(-|s)$ for any nonterminal state s (and therefore, by the above, uniquely determining a Markovian flow F);
- an objective function, such that if the model is capable of expressing any action distribution and the objective function is globally minimized, then the constraint (8) is satisfied for the corresponding Markovian flow F .

The forward policy of a GFlowNet can be used to sample trajectories from the corresponding Markovian flow F by iteratively taking actions according to policy $P_F(-|s)$. If the objective function is globally minimized, then the likelihood of terminating at x is proportional to $R(x)$.

In general, an objective optimizing for (8) cannot be minimized directly because $F(x)$ is a sum over all trajectories leading to x , and computing it may not be practical. Therefore, two local objectives – *flow matching* and *detailed balance* have previously been proposed.

Flow matching objective ([Bengio et al., 2021a](#)). A model $F_\theta(s, s')$ ³ with learnable parameters θ approximates the edge flows $F(s \rightarrow s')$. The corresponding forward policy is given by $P_F(s' | s; \theta) \propto F_\theta(s, s')$ (Eq. (6)). Denote the corresponding Markovian flow by F_θ and distribution over trajectories by P_θ . The parameters are trained to minimize the error in the flow matching constraint (2) for all noninitial and nonterminal nodes s :

³In practice, it is convenient and more economical to provide a representation of s to the neural net, which simultaneously outputs the flows $F_\theta(s, s')$ for all s that are reachable by an action from s .

$$\mathcal{L}_{\text{FM}}(s) = \left(\log \frac{\sum_{(s'' \rightarrow s) \in \mathcal{A}} F_\theta(s'', s)}{\sum_{(s \rightarrow s') \in \mathcal{A}} F_\theta(s, s')} \right)^2 \quad (9)$$

and a similar objective \mathcal{L}'_{FM} pushing the inflow at $x \in \mathcal{X}$ to equal $R(x)$ at terminal nodes x ⁴. This objective is optimized for nonterminal states s and terminal states x from trajectories sampled from a training policy π_θ . Usually, π_θ is chosen to be a tempered (higher temperature) version of $P_F(-|s, \theta)$, which also helps exploration during training. The parameters are updated with stochastic gradient

$$\mathbb{E}_{\tau=(s_0 \rightarrow \dots \rightarrow s_n) \sim \pi_\theta} \nabla_\theta \left[\sum_{t=1}^{n-1} \mathcal{L}_{\text{FM}}(s_t) + \mathcal{L}'_{\text{FM}}(s_n) \right]. \quad (10)$$

As per Proposition 10 of Bengio et al. (2021b), if the training policy π_θ has full support, and a global minimum of the expected loss (9) over states on trajectories sampled from π_θ is reached, then the GFlowNet samples from the target distribution (i.e., F_θ satisfies (8)).

Detailed balance objective (Bengio et al., 2021b). A neural network model with parameters θ has input s and three kinds of outputs: an estimated state flow $F_\theta(s)$, an estimated distribution over children $P_B(-|s; \theta)$, and an estimated distribution over parents $P_F(-|s; \theta)$. The policy $P_F(-|s; \theta)$ and the initial state flow $F_\theta(s_0)$ uniquely determine a Markovian flow F_θ , which is not necessarily compatible with the estimated backward policy $P_B(-|s; \theta)$. The error in the detailed balance constraint (7) is optimized on actions $(s \rightarrow s')$ between nonterminal nodes seen along trajectories,

$$\mathcal{L}_{\text{DB}}(s, s') = \left(\log \frac{F_\theta(s) P_F(s'|s; \theta)}{F_\theta(s') P_B(s|s'; \theta)} \right)^2. \quad (11)$$

and a similar constraint $\mathcal{L}'_{\text{DB}}(s, x)$ is optimized at actions leading to terminal nodes. Similarly to flow matching, the parameters are updated with stochastic gradient

$$\mathbb{E}_{(s_0 \rightarrow \dots \rightarrow s_n) \sim \pi_\theta} \nabla_\theta \left[\sum_{t=1}^{n-1} \mathcal{L}_{\text{DB}}(s_{t-1}, s_t) + \mathcal{L}'_{\text{DB}}(s_{n-1}, s_n) \right]. \quad (12)$$

along trajectories sampled from a training policy π_θ . By Proposition 6 of Bengio et al. (2021b), a global minimum of the expected detailed balance loss under a π_θ with full support specifies a GFlowNet that samples from the target distribution, i.e., the flow F_θ satisfies (8).

Remarks. In some problems, the directed graph G is a tree, so each state has only one parent and there is a unique trajectory from s_0 to any state s . For instance, this is the

⁴The log in (9) is for numerical reasons, and as proposed by Bengio et al. (2021a) an additive smoothing coefficient can also be inserted inside the log.

Algorithm 1 Training a GFlowNet with trajectory balance

input Reward function $R : \mathcal{X} \rightarrow \mathbb{R}_{>0}$,

input Model and optimizer hyperparameters

1: Initialize models P_F, P_B, Z with parameters θ

2: **repeat**

3: Sample trajectory $\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_n)$ from policy $P_F(-|s; \theta)$ or a tempered version of it

4: $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{\text{TB}}(\tau)$ {gradient update on (14)}

5: **until** convergence monitoring on running $\mathcal{L}_{\text{TB}}(\tau)$

case for autoregressive sequence generation, where each action appends a word to the end of a partially constructed sequence (§5.3). If G is a tree, then P_B is trivial, and the detailed balance objective reduces to the flow matching objective, which in turn can be shown to be equivalent to to Soft Q-Learning (Haarnoja et al., 2017; Buesing et al., 2019) with temperature $\alpha = 1$, a uniform $q_{a'}$, and $\gamma = 1$.

3. Trajectory balance

Let F be a Markovian flow and P the corresponding distribution over complete trajectories, defined by (3), and let P_F and P_B be forward and backward policies determined by F . A direct algebraic manipulation of Eqs. (3,4,5) gives the *trajectory balance constraint* for any complete trajectory $\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_n = x)$:

$$Z \prod_{t=1}^n P_F(s_t | s_{t-1}) = F(x) \prod_{t=1}^n P_B(s_{t-1} | s_t), \quad (13)$$

where we have used that $P(s_n = x) = \frac{F(x)}{Z}$.

As explained in §A.2, the trajectory balance constraint (13) and the detailed balance constraint (7) are special cases of one general constraint, which is the detailed balance equation for a power of the stochastic transition operator.

Trajectory balance as an objective. We propose to convert (13) into an objective to be optimized along trajectories sampled from a training policy. Suppose that a model with parameters θ outputs estimated forward policy $P_F(-|s; \theta)$ and backward policy $P_B(-|s; \theta)$ for states s (just as for detailed balance above), as well as a global scalar Z_θ estimating $F(s_0)$. The scalar Z_θ and forward policy $P_F(-|s; \theta)$ uniquely determine a Markovian flow F_θ .

For a trajectory $\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_n = x)$, define the *trajectory loss*

$$\mathcal{L}_{\text{TB}}(\tau) = \left(\log \frac{Z_\theta \prod_{t=1}^n P_F(s_t | s_{t-1}; \theta)}{R(x) \prod_{t=1}^n P_B(s_{t-1} | s_t; \theta)} \right)^2. \quad (14)$$

If π_θ is a training policy – usually that given by $P_F(-|s; \theta)$ or a tempered version of it – then the trajectory loss is up-

dated along trajectories sampled from π_θ , i.e., with stochastic gradient

$$\mathbb{E}_{\tau \sim \pi_\theta} \nabla_\theta \mathcal{L}_{\text{TB}}(\tau). \quad (15)$$

The full algorithm, with batch size of 1, is presented as Algorithm 1 and its correctness is guaranteed by the following.

Proposition 1. *Let R be a positive reward function on \mathcal{X} .*

- (a) *If $P_F(-|-; \theta)$, $P_B(-|-; \theta)$, and Z_θ are the forward and backward policies and normalizing constant of a Markovian flow F satisfying (8), then $\mathcal{L}_{\text{TB}}(\tau) = 0$ for all complete trajectories τ .*
- (b) *Conversely, suppose that $\mathcal{L}_{\text{TB}}(\tau) = 0$ for all complete trajectories τ . Then the corresponding Markovian flow F_θ satisfies (8), and $P_F(-|-; \theta)$ samples proportionally to the reward.*

The proof is given in §A.1. In particular, if π_θ has full support and $\mathbb{E}_{\tau \sim \pi_\theta} \mathcal{L}_{\text{TB}}(\tau)$ is globally minimized over all forward and backward policies (P_F, P_B) and normalizing constants Z , then the corresponding Markovian flow F_θ satisfies (8) and $P_F(-|-; \theta)$ samples proportionally to the reward. (The positivity assumption on R is necessary to avoid division by 0 in (14), but can be relaxed by introduction of smoothing constants, just as was done for flow matching and detailed balance in Bengio et al. (2021a;b).)

Remarks. As discussed in Section 2, in the case of autoregressive generation, G is a directed tree, where each $s \in \mathcal{S}$ has a single parent state. In this case P_B is trivially $P_B = 1$, $\forall s \in \mathcal{S}$. The trajectory balance objective then simplifies to

$$\mathcal{L}_{\text{TB}}(\tau) = \left(\log \frac{Z_\theta \prod_{t=1}^n P_F(s_t | s_{t-1}; \theta)}{R(x)} \right)^2 \quad (16)$$

Training considerations. We have found it beneficial to parametrize Z in the logarithmic domain ($\log Z$ is the trainable parameter) and output logits for $P_F(-|s; \theta)$ and $P_B(-|s; \theta)$, so that all products in (14) become sums under the logarithm. This is consistent with the log-domain parametrization of flows in Bengio et al. (2021a). In addition, because all gradient updates affect the normalizing scalar Z , we found it helpful to set a higher learning rate for Z than for the parameters giving P_F and P_B .⁵

3.1. Canonical choice of reward-matching flow

The constraint (8), in general, does not have a unique solution: if the underlying undirected graph of G has cycles,

⁵Because the loss (14) is quadratic in $\log Z$, gradient updates on $\log Z$ are equivalent to setting it to a weighted moving average of the discrepancy between $\log \prod P_F$ and $\log(R(x) \prod P_B)$. Optimizers enhanced with momentum complicate things: we leave empirical investigation of these questions to future work.

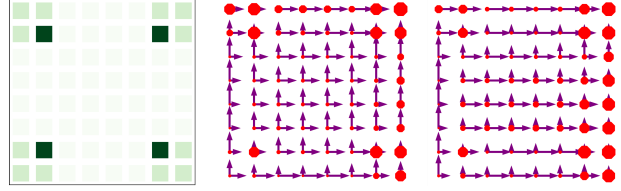


Figure 2. Left: The reward function on an 8×8 grid environment (§5.1) with $R_0 = 0.1$. Centre and right: Two forward action policies – with fixed uniform P_B and with a learned non-uniform P_B – that sample from this reward. The lengths of arrows pointing up and right from each state are proportional to the likelihoods of the corresponding actions under P_F , and the sizes of the red octagons are proportional to the termination action likelihoods.

there may be multiple Markovian flows whose corresponding action policies sample proportionally to the reward. However, by the uniqueness properties, for any choice of backward policy $P_B(-|-)$, there is a unique flow satisfying (8), and thus a unique corresponding forward policy $P_F(-|-)$ for states with nonzero flow. (See Fig. 2.)

In some settings, it may be beneficial to fix the backward policy P_B and train only the parameters giving P_F and Z_θ . For example, it may be difficult to construct a model that outputs a distribution over the parents of a given input state (e.g., for the molecule domain (§5.2), it is hard to force invariance to molecule isomorphism). A natural choice is to set $P_B(-|s)$ to be uniform over all the parents of a state s , i.e., $P_B(-|s) = 1/\#\{s' \mid (s' \rightarrow s) \in \mathcal{A}\}$.

4. Related work

GFlowNets on tree-structured DAGs (autoregressive generation) are equivalent to RL with appropriate entropy regularization or soft Q-learning and control as inference (Buesing et al., 2019; Haarnoja et al., 2017; 2018). The experiments of Bengio et al. (2021a), and accompanying discussion, show how these methods can fail badly in the general DAG case that is handled by GFlowNets. The manner of propagating a training signal over sequences of several actions in trajectory balance also has some similarities to losses used in RL computed on subtrajectories (Nachum et al., 2017).

GFlowNets are also related to MCMC methods for sampling from unnormalized densities. While there has been work on accelerating or partially amortizing sampling from unnormalized densities over discrete spaces when exact sampling is intractable (Grathwohl et al., 2021; Dai et al., 2020), some of it domain- or problem-specific (Xie et al., 2021), GFlowNets treat the compositional structure in data as a learning problem, not as a bias to build in to the sampler. Thus, the cost is borne by the learning of that structure through sampling, not through search at generation time.

5. Experiments

We evaluate the proposed trajectory balance objective against prior objectives for training GFlowNets as well as well as standard methods for learning policies that approximately sample objects proportionally to their rewards, like MCMC and techniques from RL. Our experiments span various domains, including the hypergrid and molecule synthesis tasks from Bengio et al. (2021a), as well as two new tasks – bit sequence generation and anti-microbial peptide generation – in which G is a directed tree.

5.1. Hypergrid environment

In this subsection, we study a synthetic hypergrid environment introduced by Bengio et al. (2021a). This task is easier than others we study, but we include it for completeness, and because it allows us to illustrate some interesting behaviours.⁶

In this environment, the nonterminal states \mathcal{S}° form a D -dimensional hypergrid with side length H :

$$\mathcal{S}^\circ = \{(s^1, \dots, s^D) \mid s^d \in \{0, 1, \dots, H-1\}, d = 1, \dots, D\},$$

and actions are operations of incrementing one coordinate in a state by 1 without exiting the grid, i.e.,

$$(s^1, \dots, s^d, \dots, s^D) \rightarrow (s^1, \dots, s^d + 1, \dots, s^D).$$

The initial state is $(0, \dots, 0)$. For every nonterminal state s , there is also a termination action that transitions to a corresponding terminal state s^\top (cf. footnote 1). The reward at a terminal state $s^\top = (s^1, \dots, s^D)^\top$ is given by

$$R(s^\top) = R_0 + 0.5 \prod_{d=1}^D \mathbb{I} \left[\left| \frac{s^d}{H-1} - 0.5 \right| \in (0.25, 0.5] \right] \\ + 2 \prod_{d=1}^D \mathbb{I} \left[\left| \frac{s^d}{H-1} - 0.5 \right| \in (0.3, 0.4) \right]$$

where \mathbb{I} is an indicator function and R_0 is a constant controlling the difficulty of exploration. This reward has peaks of height $2.5 + R_0$ near the corners of the hypergrid, surrounded by plateaux of height $0.5 + R_0$. These plateaux are separated by wide troughs with reward R_0 . An illustration with $H = 8$ and $D = 2$ is shown in the left panel of Fig. 2. This environment evaluates the ability of a GFlowNet to generalize from visited states to infer the existence of yet-unvisited modes.

We train GFlowNets with the detailed balance (DB) and trajectory balance (TB) objectives with different H , D , and R_0 , in addition to reproducing the flow matching (FM) experiments and non-GFlowNet baselines based upon Bengio et al.

(2021a)’s published code. Our GFlowNet policy model is a multilayer perceptron (MLP) that accepts as input a one-hot encoding of a state s (with the goal of enabling generalization) and outputs logits of the forward and backward policies $P_F(-|\cdot; \theta)$ and $P_B(-|\cdot; \theta)$ (as well as the estimated state flow $F_\theta(s)$ in the case of DB). The forward policy, backward policy, and state flow models share all but the last weight matrix of the MLP. This is consistent with Bengio et al. (2021a)’s model, where an identical architecture accepted s as input and output estimated flows $F_\theta(s, s')$ for all children s' of s . Details are given in §B.1.

We consider a 4-dimensional grid with $H = 8$ and a 2-dimensional grid with $H = 64$. The two grids have the same number of terminal states, but the 2-dimensional grid has longer expected trajectory lengths. For both grid sizes, we consider $R_0 = 0.1, 0.01, 0.001$, with smaller R_0 giving environments that are more difficult to explore due to the lower likelihood for models to cross the low-reward valley. For the models trained with DB and TB, we also explore the effect of fixing the backward policy to be uniform (§3.1).

Results. In Fig. 3, we plot the evolution over the course of training of the L_1 error between the true reward distribution (the reward $R(x)$ normalized over all terminal states) and the empirical distribution of the last $2 \cdot 10^5$ visited states for all settings. Although convergence to the same stable minimum is achieved by all models and settings, DB and TB training tend to converge faster than FM, with a slight benefit of TB over detailed balance in the 4-D environment. However, we caution against overgeneralizing from these conclusions, as we experimented with a single learning rate setting. In particular, while the DB and TB models share their parametrization of P_F and P_B and use the same learning rate, the FM model hyperparameters were copied from Bengio et al. (2021a). Because the FM model outputs flows, rather than policy logits, optimal learning rates may be different.

Effect of uniform P_B . Note the difference in learning speed between models with fixed uniform backward policy P_B and models with learned P_B . As noted in §3.1, when P_B is fixed, there is a unique $P_F(-|\cdot; \theta)$ that globally minimizes the objective, and it may be approached slowly. However, if P_B and P_F are permitted to evolve jointly, they may more quickly approach one of the many optimal solutions. This is confirmed by the much faster convergence of models with learned P_B on the 64×64 grid.

We have observed that, especially for large grid sizes, when P_B and P_F are both learned, the model has a bias towards first taking all actions in one coordinate direction, then proceeding in the other direction until terminating (as in the right panel of Fig. 2), perhaps because a constant distribution over two actions (‘continue to the right’ and ‘terminate’)

⁶Example code appears at [this address](#).

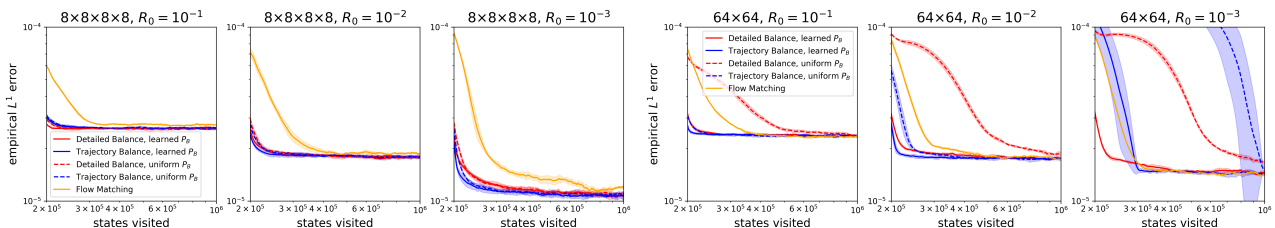


Figure 3. Empirical L^1 error between true and sampled state distributions on the grid environment with varying grid size and R_0 . Mean and standard error over 5 seeds. The curves for PPO and MCMC baseline would lie outside the plot bounds.

can be modeled with higher precision over a large portion of the grid than the complex position-dependent distribution as shown in the centre panel of Fig. 2.

5.2. Small drug molecule synthesis

Next, we consider the molecule generation task (Xie et al., 2021; Jin et al., 2020; Kumar et al., 2012; Gilmer et al., 2017; Shi et al., 2021) introduced for GFlowNets in Bengio et al. (2021a). We use their published code and task setup and extend it with an implementation of the TB objective.

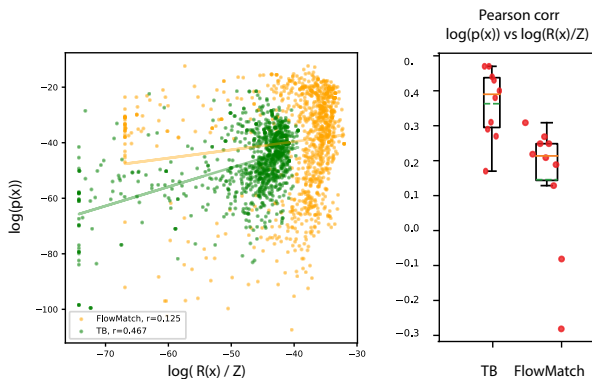


Figure 4. Pearson correlations between rewards and sampling probability. *Left*: Reward and sampling probability for two GFlowNets, trained with trajectory balance (TB) and flow matching. In y is the log-likelihood that a trajectory sampled from the learned policy $P_F(-|-; \theta)$ terminates at x . *Right*: Pearson correlations across 10 seeds, demonstrating consistently better fitting when using TB.

The goal is to generate molecules, in the form of graphs, with a low binding energy to the 4JNC inhibitor of the sEH (soluble epoxide hydrolase) protein. The graphs generated are junction trees (Jin et al., 2020) of a vocabulary of building blocks. The reward is defined as the normalized negative binding energy as predicted by a *proxy* model, itself trained to predict energies computed via docking simulations (Trott & Olson, 2010). The maximum trajectory length is 8, with the number of actions varying between around 100 and 2000 (the larger a molecule, the more possible additions exist),

making $|\mathcal{X}|$ about 10^{16} .

Results. We first plot in Fig. 4 the correlation of reward and sampling probability (the likelihood that a trajectory sampled from the learned policy $P_F(-|-; \theta)$ terminates at x) for GFlowNets trained using TB and using FM. We see a significantly better fitting by models trained with TB. The points used there are from the proxy’s training set (to which the GFlowNets do not have access in training).

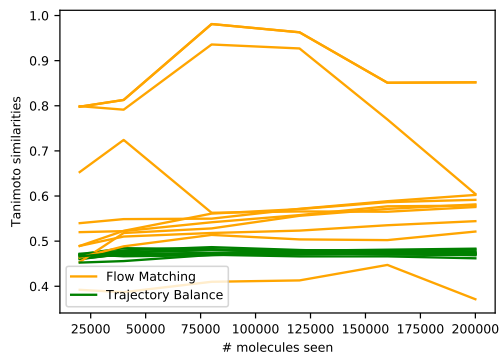


Figure 5. Average pairwise Tanimoto similarity for the top 1000 samples generated by GFlowNets trained with either TB or FM as a function of learning iterations. Each line is one of 10 seeds for each objective. Models trained with TB have significantly lower average similarity, hence greater diversity.

In Fig. 5 we plot the average pairwise Tanimoto similarity (Bender & Glen, 2004) for the top 1000 generated samples over the course of training. Models trained with TB consistently generate significantly more diverse molecules. These results showcase the benefits of TB, not only for faster temporal credit assignment, but for generalization and diversity. We provide additional details in §B.2.

5.3. Autoregressive sequence generation

Finally, we evaluate the trajectory balance (TB) objective in the context of autoregressive generation of sequences. We first consider a synthetic bit sequence generation, where we study the effect of trajectory length, as well as the size of the action space on the learning dynamics in GFlowNets. We then consider the more realistic task of generating peptides,

which are short protein sequences with anti-microbial properties, and evaluate the GFlowNets trained with TB against standard RL and MCMC baselines.

5.3.1. BIT SEQUENCES

Task. The goal is to generate bit sequences of a fixed length $n = 120$ ($\mathcal{X} = \{0, 1\}^n$), where the reward is designed to have modes at a given fixed set $M \subset \mathcal{X}$ that is unknown to the learner. The reward for a sequence x is defined as $R(x) = \exp(-\min_{y \in M} d(x, y))$, where d is the edit distance. We describe the procedure to generate M in §B.3.

For different integers k dividing n , we design action spaces for left-to-right generation of sequences in \mathcal{X} . For a fixed k , define the vocabulary V to be the set of k -bit sequences (e.g., for $k = 2$, $V = \{00, 01, 10, 11\}$). The state space \mathcal{S} is then defined as the set of partial sequences of length at most n that can be constructed using V , that is, the bit sequences with length a multiple of k . Actions consist of appending a word from V to the right of a state, so a forward policy on this state space is the same as an autoregressive sequence model over the vocabulary V . The sequences of length n are terminal, so any complete trajectory has exactly $\frac{n}{k}$ actions. Varying k while fixing \mathcal{X} and M allows us to study the effect of the tradeoff between trajectory lengths ($\frac{n}{k}$) and the action space sizes ($|V| = 2^k$) without changing the underlying learning problem.

We compare GFlowNets trained with the TB objective against GFlowNets trained with the FM objective and two non-GFlowNet baselines: A2C with Entropy Regularization (Williams & Peng, 1991; Mnih et al., 2016) and MARS (Xie et al., 2021). We use a Transformer-based architecture (Vaswani et al., 2017) across all the methods. We discuss further training details and hyperparameters in §B.3.

To evaluate the methods we use (1) Spearman correlation between the probability of generating the sequence $p(x) = \frac{F(x)}{Z}$ and its reward $R(x)$ on a test set sampled approximately uniformly over the possible values of the reward, (2) number of modes captured (number of reference sequences from M for which a candidate within a distance δ has been generated). In our experiments, $n = 120$, $|M| = 60$, $k \in \{1, 2, 4, 6, 8, 10\}$, and $\delta = 20$.

Results. Fig. 6 presents the results for the Spearman correlation evaluation. We observe that GFlowNets trained with the TB objective learn policies that correlate best with the reward $R(x)$ across all action spaces. In particular, we observe the effect of inefficient credit assignment in GFlowNets trained with FM, as the correlation improves with increasing k , i.e., shorter trajectories. On the other hand, large action spaces also hurt GFlowNets trained with the FM objective, while the TB objective is robust to them. Additionally, we can observe in Fig. 7 that for fixed k , GFlowNets trained

with TB discover more modes faster than other methods.

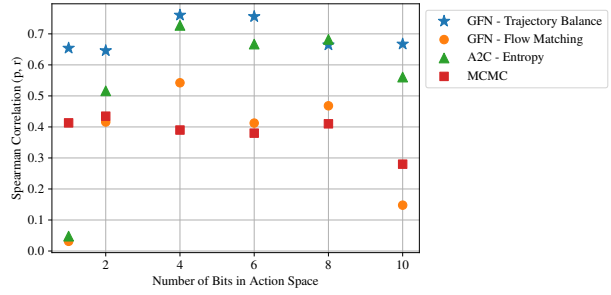


Figure 6. Spearman correlation of the sampling probability under different learned policies and reward on a test set, plotted against the number of bits k in the symbols in V in the bit sequence generation task. GFlowNets trained with trajectory balance learn policies that have the highest correlation with the reward $R(x)$ and are robust to length and vocabulary size.

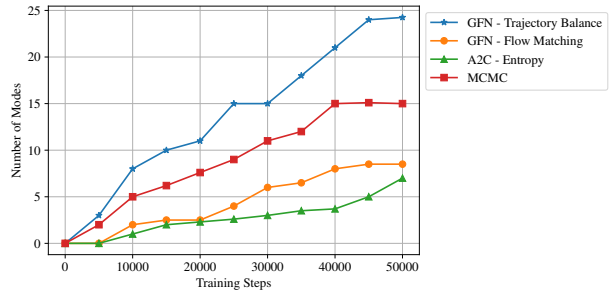


Figure 7. Number of modes discovered over the course of training on the bit sequence generation task with $k = 1$. GFlowNets trained with trajectory balance discover more modes faster.

5.3.2. ANTI-MICROBIAL PEPTIDE (AMP) GENERATION

In this section, we consider the practical task of generating peptide sequences that have anti-microbial activity. The goal is to generate a protein sequence (where the vocabulary consists of 20 amino acids and a special end-of-sequence action), with maximum length 60. We take 6438 known AMP sequences and 9522 non-AMP sequences from the DBAASP database (Pirtskhalava et al., 2021). We then train a classifier on this dataset, using 20% of the data as a validation set. The probability output by this model for a sequence to be classified as an AMP is used as the reward $R(x)$ in our experiments.

The state and actions are designed just as in the previous experiment, with each action appending a symbol to the right of a state. We again compare TB and FM, as well as A2C with entropy regularization and MCMC as baselines. We again use Transformers for all the experiments on this

Table 1. Results on the AMP generation task.

	Top 100 reward	Top 100 diversity
GFN- \mathcal{L}_{TB}	0.85 ± 0.03	18.35 ± 1.65
GFN- \mathcal{L}_{FM}	0.78 ± 0.05	12.61 ± 1.32
AAC-ER	0.79 ± 0.02	7.32 ± 0.76
MCMC	0.75 ± 0.02	12.56 ± 1.45

task; see further details in Appendix B.4. We generate 2048 sequences from each method, and pick the top 100 sequences ranked by their reward $R(x)$. As metrics, we use the mean reward for these 100 sequences and the average pairwise edit distance among them as a measure of *diversity*.

Results. We present the results in Table 1, where we observe that GFlowNets trained with TB outperform all baselines on both performance and diversity metrics.

6. Discussion and conclusion

We introduced a novel training loss for GFlowNets, trajectory balance (TB), which yields faster and better training than the previously proposed flow matching (FM) and detailed balance (DB) losses. We proved that this objective, when minimized, yields the desired GFlowNet property of sampling from the target distribution specified by an unnormalized reward function. This new loss was motivated by the observation that the FM and DB losses are local in the action sequence and, similarly to temporal-difference learning, may require many iterations for credit assignment to propagate to early actions. We found that TB discovered more modes of the energy function faster and was more robust than FM and DB to the exponential growth of the state space, due in part to the lengths of sequences and in part to the size of the action space.

A factor to consider when interpreting our experimental results is that because we use a neural net rather than a tabular representation of policies, the early states’ transitions are informed by downstream credit assignment via parameter sharing. Early states also get many more visits because there are more possible states near the ends of sequences than near the initial state. Finally, TB trades off the advantage of immediately providing credit to early states with the disadvantage of relying on sampling of long trajectories and thus a potentially higher variance of the stochastic gradient.

All in all, we found that trajectory balance is a superior training objective in a broad set of experiments, making it the default choice for future work on GFlowNets.

Acknowledgments

This research was enabled in part by computational resources provided by Compute Canada (www.computecanada.ca). All authors are funded by their primary academic institution. We also acknowledge funding from CIFAR, Samsung, IBM, Microsoft, and the Banting Postdoctoral Fellowship.

The authors are grateful in particular to Dinghui Zhang as well as the members of the Mila GFlowNet group for many fruitful research discussions.

References

- Bender, A. and Glen, R. C. Molecular similarity: a key technique in molecular informatics. *Organic & biomolecular chemistry*, 2(22):3204–3218, 2004.
- Bengio, E., Pineau, J., and Precup, D. Interference and generalization in temporal difference learning. *International Conference on Machine Learning (ICML)*, 2020.
- Bengio, E., Jain, M., Korablyov, M., Precup, D., and Bengio, Y. Flow network based generative models for non-iterative diverse candidate generation. *Neural Information Processing Systems (NeurIPS)*, 2021a.
- Bengio, Y., Deleu, T., Hu, E., Lahlou, S., Tiwari, M., and Bengio, E. GFlowNet foundations. *arXiv preprint 2111.09266*, 2021b.
- Buesing, L., Heess, N., and Weber, T. Approximate inference in discrete distributions with monte carlo tree search and value functions. *Artificial Intelligence and Statistics (AISTATS)*, 2019.
- Dai, H., Singh, R., Dai, B., Sutton, C., and Schuurmans, D. Learning discrete energy-based models via auxiliary-variable local exploration. *Neural Information Processing Systems (NeurIPS)*, 2020.
- Ford, L. R. and Fulkerson, D. R. Maximal flow through a network. *Canadian Journal of Mathematics*, 8:243–248, 1956.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. *International Conference on Machine Learning (ICML)*, 2017.
- Grathwohl, W., Swersky, K., Hashemi, M., Duvenaud, D. K., and Maddison, C. J. Oops I took a gradient: Scalable sampling for discrete distributions. *International Conference on Machine Learning (ICML)*, 2021.
- Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-based policies.

- International Conference on Machine Learning (ICML)*, 2017.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *International Conference on Machine Learning (ICML)*, 2018.
- Jin, W., Barzilay, R., and Jaakkola, T. Chapter 11. junction tree variational autoencoder for molecular graph generation. *Drug Discovery*, pp. 228–249, 2020. ISSN 2041-3211.
- Kumar, A., Voet, A., and Zhang, K. Fragment based drug design: from experimental to computational approaches. *Current medicinal chemistry*, 19(30):5128–5147, 2012.
- Kumar, A., Agarwal, R., Ghosh, D., and Levine, S. Implicit under-parameterization inhibits data-efficient deep reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2021.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. *International Conference on Machine Learning (ICML)*, 2016.
- Nachum, O., Norouzi, M., Xu, K., and Schuurmans, D. Bridging the gap between value and policy based reinforcement learning. *Neural Information Processing Systems (NeurIPS)*, 2017.
- Pirtskhalava, M., Amstrong, A. A., Grigolava, M., Chubinidze, M., Alimbarashvili, E., Vishnepolsky, B., Gabrielian, A., Rosenthal, A., Hurt, D. E., and Tartakovsky, M. Dbasp v3: Database of antimicrobial/cytotoxic activity and structure of peptides as a resource for development of new therapeutics. *Nucleic Acids Research*, 49(D1):D288–D297, 2021.
- Shi, Y., Huang, Z., Feng, S., Zhong, H., Wang, W., and Sun, Y. Masked label prediction: Unified message passing model for semi-supervised classification. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Trott, O. and Olson, A. J. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of computational chemistry*, 31(2):455–461, 2010.
- van Hasselt, H., Doron, Y., Strub, F., Hessel, M., Sonnerat, N., and Modayil, J. Deep reinforcement learning and the deadly triad. *arXiv preprint 1812.02648*, 2018.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *Neural Information Processing Systems (NeurIPS)*, 2017.
- Williams, R. J. and Peng, J. Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268, 1991.
- Xie, Y., Shi, C., Zhou, H., Yang, Y., Zhang, W., Yu, Y., and Li, L. MARS: Markov molecular sampling for multi-objective drug discovery. *International Conference on Learning Representations (ICLR)*, 2021.

A. Theoretical appendix

A.1. Proof of Proposition 1

Recall Proposition 1:

Proposition. *Let R be a positive reward function on \mathcal{X} .*

- (a) *If $P_F(-|-; \theta)$, $P_B(-|-; \theta)$, and Z_θ are the forward and backward policies and normalizing constant of a Markovian flow F satisfying (8), then $\mathcal{L}_{\text{TB}}(\tau) = 0$ for all complete trajectories τ .*
- (b) *Conversely, suppose that $\mathcal{L}_{\text{TB}}(\tau) = 0$ for all complete trajectories τ . Then the corresponding Markovian flow F_θ satisfies (8), and $P_F(-|-; \theta)$ samples proportionally to the reward.*

Proof. Part (a) is an elementary manipulation of the trajectory balance constraint (13), with $R(x)$ substituted for $F(x)$ by the reward matching assumption (8).

Conversely, if $\mathcal{L}_{\text{TB}}(\tau) = 0$ for all complete trajectories $\tau = (s_0 \rightarrow \dots \rightarrow s_n = x)$, then the policies $P_F(-|-; \theta)$ and $P_B(-|-; \theta)$ satisfy the constraint

$$Z \prod_{t=1}^n P_F(s_t | s_{t-1}; \theta) = R(x) \prod_{t=1}^n P_B(s_{t-1} | s_t; \theta). \quad (17)$$

Let x be a terminal state. By iterating the law of total probability, we have

$$\sum_{\tau=(s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_n=x)} \prod_{t=1}^n P_B(s_{t-1} | s_t; \theta) = 1. \quad (18)$$

(Each term in this sum is the conditional likelihood of τ conditioned on terminating at $s_n = x$ under the the Markovian flow F'_θ uniquely determined by setting terminal state flows $F'_\theta(x) = R(x)$ and backward policy $P_B(-|-; \theta)$, cf. the uniqueness properties.)

On the other hand, we have

$$\begin{aligned} F_\theta(x) &= \sum_{\tau=(s_0 \rightarrow \dots \rightarrow s_n)=x} F_\theta(\tau) && \text{(by definition of state flows)} \\ &= \sum_{\tau=(x_0 \rightarrow \dots \rightarrow s_n)=x} Z \prod_{t=1}^n P_F(s_t | s_{t-1}; \theta) && \text{(by (4))} \\ &= \sum_{\tau=(x_0 \rightarrow \dots \rightarrow s_n)=x} R(x) \prod_{t=1}^n P_B(s_{t-1} | s_t; \theta) && \text{(by (17))} \\ &= R(x) && \text{(by (18)).} \end{aligned}$$

We conclude that F_θ satisfies (8), as desired.

(We remark that one can show in a similar way that $F_\theta(s \rightarrow s') = F'_\theta(s \rightarrow s')$ for all actions $(s \rightarrow s') \in \mathcal{A}$, and thus, by the uniqueness properties, $F_\theta = F'_\theta$, i.e., the forward and backward policies determine the same Markovian flow.) \square

A.2. Generalizations

The trajectory balance constraint (13) can be generalized to partial (not complete) trajectories, i.e., those that do not start at s_0 and end in a terminal state. Generalizations such as those we present here could be useful for a future goal of modularized or hierarchical GFlowNets, where each module (or low-level GFlowNet) can apply them to just the subsequences they have access to (cf. §9.4 and §10.2 in Bengio et al. (2021b)).

Subtrajectory balance. If $\tau = (s_m \rightarrow s_{m+1} \rightarrow \dots \rightarrow s_n)$ is a partial trajectory (i.e., $(s_t \rightarrow s_{t+1}) \in \mathcal{A}$ for all t), then, for any Markovian flow F with forward and backward policies P_F and P_B ,

$$F(s_m) \prod_{t=m}^{n-1} P_F(s_{t+1} | s_t) = F(s_n) \prod_{t=m}^{n-1} P_B(s_t | s_{t+1}). \quad (19)$$

This can be derived by showing that both sides are equal to

$$\sum_{\tau=(\dots \rightarrow s_m \rightarrow s_{m+1} \rightarrow \dots \rightarrow s_n \rightarrow \dots) \in \mathcal{T}} F(\tau). \quad (20)$$

The trajectory balance constraint (13) is the special case of this for full trajectories, while the detailed balance constraint (7) is the special case of trajectories with only one edge. This subtrajectory balance constraint can be converted into a learning objective: a model can output estimated state flows $F_\theta(s)$ only for certain nonterminal states s (“hubs”), and the error in (19) optimized along segments of trajectories between these hubs. Thus the detailed balance loss corresponds to all nodes being hubs, and the trajectory balance loss corresponds to only the initial state s_0 being a hub.

Non-forward trajectories. Trajectory balance has a more general form for trajectories that have a mix of forward and backward steps. Here we describe just one example: terminal-terminal paths that take several backward steps, then take several forward steps.

Let $s_1 = s'_1$ be any state (not necessarily a child of the GFlowNet’s initial state s_0) and $(s_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_n)$ and $(s'_1 \rightarrow s'_2 \rightarrow \dots \rightarrow s'_{n'})$ two trajectories from s_0 to terminal states. Then the following must hold for any Markovian flow F :

$$R(s'_{n'}) \prod_{t=1}^{n'-1} P_B(s'_t | s'_{t+1}) \prod_{t=1}^{n-1} P_F(s_{t+1} | s_t) = R(s_n) \prod_{t=1}^{n-1} P_B(s_t | s_{t+1}) \prod_{t=1}^{n'-1} P_F(s'_{t+1} | s'_t). \quad (21)$$

That is, the path that goes “backward, then forward” from s_n to $s'_{n'}$, must have the same likelihood no matter in which direction it is traversed, up to the ratio of rewards at the endpoints. A simple way to derive (21) is by writing the trajectory balance constraint for two paths from the GFlowNet’s initial state to s_n and $s'_{n'}$, that are identical until s_1 and then diverge, then dividing one constraint by the other. Notice that the flow $F(s_1)$ is not present here. Thus, (21) can be converted into an learning objective does not require a model to output any state flows (even the initial state flow Z). Such terminal-terminal paths could also be used for exploration of \mathcal{X} with MCMC-like local search algorithms.

B. Experimental appendix

B.1. Hypergrid

For the GFlowNet policy model, we use an MLP of the same architecture as Bengio et al. (2021a), with 2 hidden layers of 256 hidden units each. We train all models with a learning rate of 0.001 (P_F and P_B policy model) and 0.1 (Z_θ) with up to 10^6 sampled trajectories with a batch size of 16, using the Adam optimizer with all other parameters at their default values.

To reproduce the flow matching and non-GFlowNet baselines, we used the code published by Bengio et al. (2021a) out of the box.

B.2. Molecule synthesis

We use the dataset and proxy model provided by Bengio et al. (2021a). We also train GFlowNet using the same architecture and hyperparameters but using the trajectory balance loss presented in this paper, using fixed uniform backward policy P_B . The binding scores in the provided dataset were computed with AutoDock (Trott & Olson, 2010).

In all experiments done for GFlowNets trained with either trajectory balance or flow matching, we used reward exponent $\beta = 8$. In contrast to Bengio et al. (2021a), we used a more exploratory training policy: with probability 0.25 (instead of the original 0.05) trajectories are set to stop at some length k , which is chosen uniformly between 3 and 8, the minimum and maximum allowed trajectory length respectively.

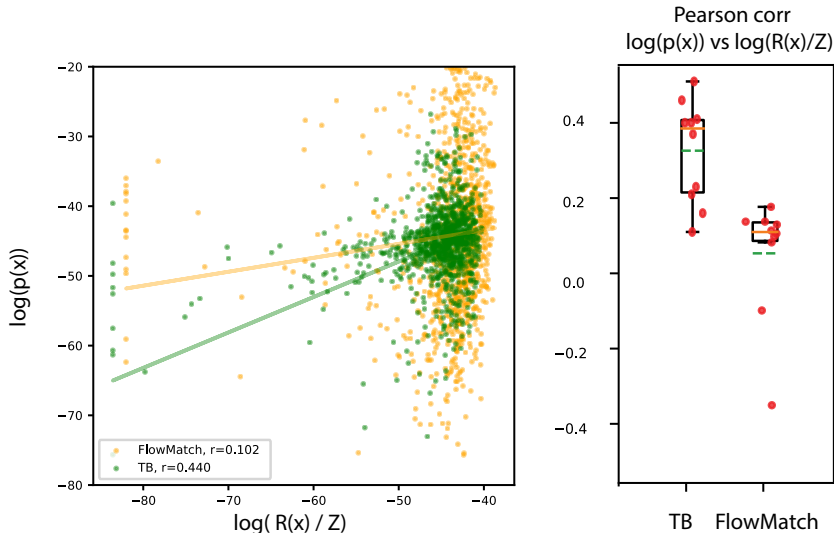


Figure 8. Correlation of reward and sampling probability under learned GFlowNets. *Left*: Reward and sampling probability for two GFlowNets, trained with trajectory balance (TB) and flow matching. The y -coordinate is the log-likelihood that a trajectory sampled from the learned policy $P_F(-|-; \theta)$ terminates at x . *Right*: Pearson correlations of reward and sampling probability across 10 seeds. Unlike Fig. 4 where the tested set of molecules was from the *proxy*’s training set, here we use a set of points generated by a flow-matching-trained GFlowNet. Presumably this model should be better at its own distribution, but even here we see that trajectory balance makes better predictions.

B.3. Bit sequence generation

Generating reference sequences. Let H be a set of symbols (short bit sequences of length b), for instance $H = \{0110, 1100, 1111, 0000, 0011\}$. Sequences in S are then constructed by randomly combining m symbols from H , for instance, 0011110000000011 where $m = 4$. This construction imposes a structure on $R(x)$. In our experiments we set $m = \frac{n}{b}$, $b = 8$, $H = \{‘00000000’, ‘11111111’, ‘11110000’, ‘00001111’, ‘00111100’\}$.

Generating the test set. Since the reward is defined based on the edit distance from the sequences in set M , we generate a test set sampled approximately uniformly over the possible values of $R(x)$ as follows: (1) pick a mode $s \in M$, (2) modify i bits randomly $\forall i < n$ and we repeat this for all the modes.

Implementation. We implement GFlowNets with TB and FM in PyTorch for autoregressive generation tasks, along with the A2C baseline. For the MARS (MCMC) baseline we modify the implementation released by (Bengio et al., 2021a).

Hyperparameters. We use a Transformer (Vaswani et al., 2017) as the neural network architecture for all the methods. We use 3 hidden layers with hidden dimension 64 with 8 attention heads. All methods were trained for 50,000 iterations, with a minibatch size of 16. We set the the random action probability to 0.001, the reward exponent β to 2, and the sampling temperature for P_F to 1 for the GFlowNets. For trajectory balance we use a learning rate of 5×10^{-4} for the flow parameters and 5×10^{-3} for $\log Z$. For flow matching we use a learning rate of 5×10^{-4} with leaf loss coefficient $\lambda_T = 10$. For A2C with entropy regularization we share parameters between the actor and critic networks, and use learning rate of 10^{-4} with entropy regularization coefficient 10^{-3} . For the MARS baseline we set the learning rate to 5×10^{-4} . For all the methods we use the Adam optimizer.

B.4. AMP generation

Vocabulary. The vocabulary of the 20 amino acids is defined as: [‘A’, ‘C’, ‘D’, ‘E’, ‘F’, ‘G’, ‘H’, ‘I’, ‘K’, ‘L’, ‘M’, ‘N’, ‘P’, ‘Q’, ‘R’, ‘S’, ‘T’, ‘V’, ‘W’, ‘Y’]

Reward Model. We use a Transformer-based classifier, with 4 hidden layers, hidden dimension 64, and 8 attention heads. We train it with a minibatch of size 256, with learning rate 10^{-4} , with early stopping on the validation set.

Hyperparameters. As with the bit sequences, we use a Transformer (Vaswani et al., 2017) as the neural network architecture for all the methods. We use 3 hidden layers with hidden dimension 64 with 8 attention heads. All method were trained for 20,000 iterations, with a mini batch size of 16. We set the the random action probability to 0.01, the reward exponent $\beta : R(x)^\beta$ to 3, and the sampling temperature for P_F to 1 for the GFlowNets. For trajectory balance we use a learning rate of 5×10^{-3} for the flow parameters and 1×10^{-2} for $\log Z$. For flow matching we use a learning rate of 5×10^{-4} with leaf loss coefficient $\lambda_T = 25$. For A2C with entropy regularization we share parameters between the actor and critic networks, and use learning rate of 5×10^{-4} with entropy regularization coefficient 10^{-2} . For the MARS baseline we set the learning rate to 5×10^{-4} . We run the experiments on 3 seeds and report the mean and standard deviation over the three runs in Table 1.