

# Revisiting Non-Acyclic GFlowNets in Discrete Environments

N. Morozov<sup>1,\*</sup>, I. Maksimov<sup>1,\*</sup>, D. Tiapkin<sup>2,3</sup>, S. Samsonov<sup>1</sup>

<sup>1</sup>HSE University <sup>2</sup>CMAF, École Polytechnique <sup>3</sup>LMO, Université Paris-Saclay



## GFlowNets

- **GFlowNets** are models designed to learn a sampler from a complex discrete space  $\mathcal{X}$  according to a probability mass function  $\mathcal{R}(x)$  (*GFlowNet reward*) given up to an unknown normalizing constant  $Z = \sum_{x \in \mathcal{X}} \mathcal{R}(x)$ .
- Introduce a directed acyclic graph  $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ . Non-terminal states describe "incomplete" objects, with an empty object denoted as  $s_0$ , edges — adding new components to them. Terminal states are "complete" objects and coincide with  $\mathcal{X}$ .
- Introduce probability distributions over complete trajectories in backward and forward directions

$$\mathcal{P}_B(\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_{n_\tau} \rightarrow s_f)) = \frac{\mathcal{R}(s_{n_\tau})}{Z} \prod_{t=1}^{n_\tau} \mathcal{P}_B(s_{t-1}|s_t), \quad \mathcal{P}_F(\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_{n_\tau} \rightarrow s_f)) = \prod_{t=1}^{n_\tau} \mathcal{P}_F(s_t|s_{t-1}),$$

where  $\mathcal{P}_B(s_{t-1}|s_t)$  is called a *backward policy* and  $\mathcal{P}_F(s_t|s_{t-1})$  is called a *forward policy*.

- **Main goal:** find a pair of policies such that  $\mathcal{P}_F(\tau) = \mathcal{P}_B(\tau)$  for all  $\tau$ . **Then by sampling**  $\tau \sim \mathcal{P}_F$ , **terminal states are sampled with probabilities proportional to**  $\mathcal{R}(x)$ .
- Parameterize policies by neural networks and introduce a training objective that will force this constraint. E.g., *Detailed Balance* objective:

$$\mathcal{L}_{DB}(s \rightarrow s') = \left( \log \frac{\mathcal{F}_\theta(s) \mathcal{P}_F(s'|s, \theta)}{\mathcal{F}_\theta(s') \mathcal{P}_B(s|s', \theta)} \right)^2, \quad \text{where } \mathcal{F}_\theta(s) = \mathcal{R}(s) \forall s \in \mathcal{X}.$$

## GFlowNets in Non-Acyclic Environments

We rigorously prove that the theory behind discrete GFlowNets can be directly extended to allow cycles in the environment graph.

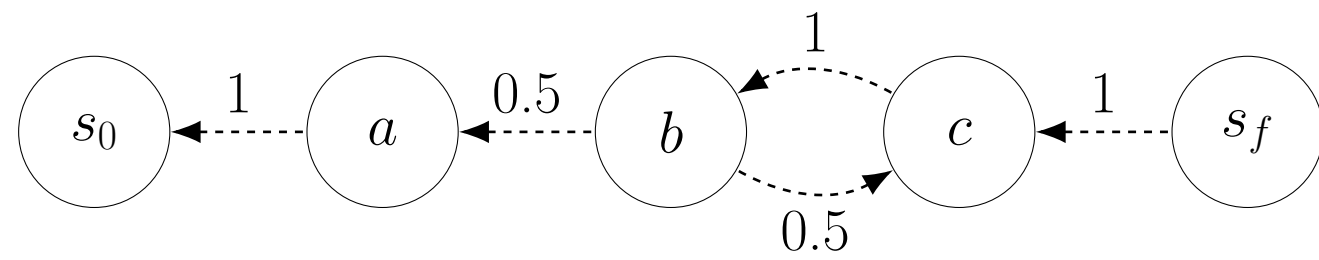


Figure 1: A graph with cycles, backward probabilities are labeled on the edges

### Main theoretical contributions and differences from the acyclic case:

- **Trajectories:** We prove expected trajectory length  $\mathbb{E}[n_\tau]$  is always finite. Under the assumption of non-zero backward probabilities  $\mathcal{P}_B(s|s') > 0$ .
- **Flows:** We redefine state flows  $\mathcal{F}(s)$  as **expected number of visits** instead of state visitation probabilities:

$$\mathcal{F}(s) = \mathcal{F}(s_f) \cdot \mathbb{E}_{\tau \sim \mathcal{P}(\tau)} \left[ \sum_{t=0}^{n_\tau+1} \mathbb{I}\{s_t = s\} \right], \quad \mathcal{F}(s_f) = Z$$

This allows us to rewrite key GFlowNet theory consistently with acyclic case

- **Expected Trajectory Length:** The expected length  $\mathbb{E}[n_\tau]$  is equal to the normalized total state flow:

$$\mathbb{E}[n_\tau] = \frac{1}{\mathcal{F}(s_f)} \sum_{s \in \mathcal{S} \setminus \{s_0, s_f\}} \mathcal{F}(s)$$

- **Constrained Optimization:** Rewriting the task and using the equation above we achieve solution with **smallest possible total flow** which allows us to control mean trajectory length of the final sampler.

$$\begin{aligned} & \min_{\mathcal{F}, \mathcal{P}_F, \mathcal{P}_B} \sum_{s \in \mathcal{S} \setminus \{s_0, s_f\}} \mathcal{F}(s) \\ & \text{subject to } \left( \log \frac{\mathcal{F}(s) \mathcal{P}_F(s'|s)}{\mathcal{F}(s') \mathcal{P}_B(s|s')} \right)^2 = 0, & \forall s \rightarrow s' \in \mathcal{E}, \\ & \mathcal{F}(s_f) \mathcal{P}_B(x|s_f) = \mathcal{R}(x), & \forall x \rightarrow s_f \in \mathcal{E}. \end{aligned}$$

- **Equivalence to RL:** We generalize the connection to RL, showing that training a non-acyclic GFlowNet with a fixed backward policy is equivalent to solving an **entropy-regularized RL problem** in a corresponding MDP.

## State Flow Regularization

Use  $\mathcal{F}_\theta(s)$  as a regularizer in the DB loss to control the expected trajectory length:

$$\left( \log \frac{\mathcal{F}_\theta(s) \mathcal{P}_F(s'|s, \theta)}{\mathcal{F}_\theta(s') \mathcal{P}_B(s|s', \theta)} \right)^2 + \lambda \mathcal{F}_\theta(s), \quad \text{where } \mathcal{F}_\theta(s) = \mathcal{R}(s) \forall s \in \mathcal{X}.$$

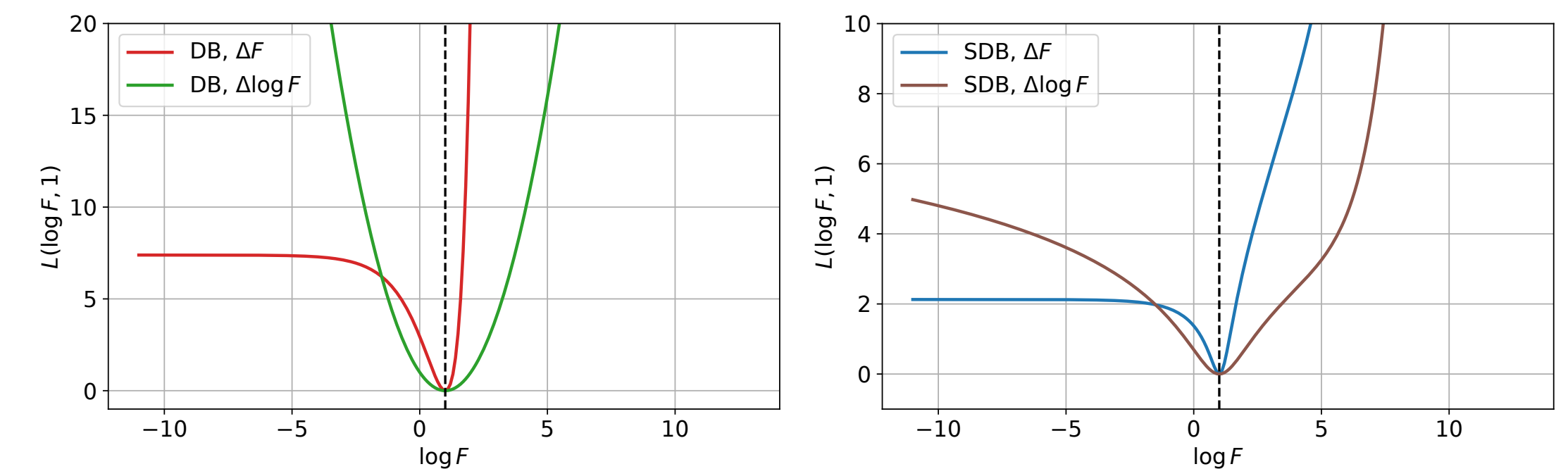
## Loss Stability and Scaling Hypothesis

(Brunswic et al., 2024) introduced the concept of **stable losses**: losses that cannot be decreased by increasing the flow on a cycle, which means *expected trajectory length is bounded during training*.

$$\mathcal{L}_{SDB}(s \rightarrow s') = \log \left( 1 + \varepsilon (\mathcal{F}_\theta(s) \mathcal{P}_F(s'|s, \theta) - \mathcal{F}_\theta(s') \mathcal{P}_B(s|s', \theta))^2 \right) (1 + \eta \mathcal{F}_\theta(s))$$

We put out the **scaling hypothesis**:

The main factor that plays a crucial role in loss stability in practice is the scale in which the error between flows is computed: log-flow scale  $\Delta \log \mathcal{F}$  vs flow scale  $\Delta \mathcal{F}$ . We hypothesize that using log-flow scale losses without regularization can lead to arbitrarily large trajectory length, while flow scale losses are biased towards solutions with smaller flows and thus do not suffer from this issue.



## Experimental Results and Observations

We run experiments on two non-acyclic environments. Namely hypergrids and permutations. The table below presents results on the **permutation environment**, where permutations of length  $n$  are generated. Action space consists of 1) swapping adjacent elements, 2) cyclic shift to the right. **Blue** indicates the best metric, **red** indicates the smallest expected trajectory length.

Loss	$n = 8$					$n = 20$				
	$C(k)$	$L^1 \downarrow$	$\Delta \mathcal{R} \downarrow$	$\Delta \log \mathcal{Z} \downarrow$	$\mathbb{E}[n_\tau]$	$C(k)$	$L^1 \downarrow$	$\Delta \mathcal{R} \downarrow$	$\Delta \log \mathcal{Z} \downarrow$	$\mathbb{E}[n_\tau]$
DB, $\Delta \mathcal{F}$	0.215 $\pm$ 0.198	0.214 $\pm$ 0.086	0.814 $\pm$ 0.826	<b>2.43 <math>\pm</math> 0.28</b>	4.31 $\pm$ 0.05	0.453 $\pm$ 0.002	0.343 $\pm$ 0.000	42.98 $\pm$ 0.000	<b>2.00 <math>\pm</math> 0.00</b>	7.55 $\pm$ 0.50
SDB, $\Delta \mathcal{F}$	0.031 $\pm$ 0.012	0.046 $\pm$ 0.023	0.074 $\pm$ 0.025	3.32 $\pm$ 0.15	4.31 $\pm$ 0.05	0.452 $\pm$ 0.001	0.343 $\pm$ 0.000	42.98 $\pm$ 0.000	<b>2.01 <math>\pm</math> 0.00</b>	7.55 $\pm$ 0.50
DB, $\Delta \log \mathcal{F}$ , $\lambda = 10^{-3}$	0.036 $\pm$ 0.015	0.056 $\pm$ 0.024	0.018 $\pm$ 0.010	2.80 $\pm$ 0.04	4.31 $\pm$ 0.05	0.041 $\pm$ 0.002	0.064 $\pm$ 0.000	0.023 $\pm$ 0.005	3.23 $\pm$ 0.00	7.55 $\pm$ 0.50
SDB, $\Delta \log \mathcal{F}$ , $\lambda = 10^{-3}$	0.037 $\pm$ 0.013	0.056 $\pm$ 0.019	0.020 $\pm$ 0.015	2.79 $\pm$ 0.04	4.31 $\pm$ 0.05	0.041 $\pm$ 0.002	0.064 $\pm$ 0.000	0.026 $\pm$ 0.003	3.22 $\pm$ 0.00	7.55 $\pm$ 0.50
DB, $\Delta \log \mathcal{F}$ , $\lambda = 10^{-5}$	<b>0.005 <math>\pm</math> 0.001</b>	<b>0.001 <math>\pm</math> 0.000</b>	<b>0.005 <math>\pm</math> 0.004</b>	4.31 $\pm$ 0.05	4.31 $\pm$ 0.05	0.017 $\pm$ 0.002	0.035 $\pm$ 0.002	<b>0.003 <math>\pm</math> 0.003</b>	7.55 $\pm$ 0.50	7.55 $\pm$ 0.50
SDB, $\Delta \log \mathcal{F}$ , $\lambda = 10^{-5}$	<b>0.005 <math>\pm</math> 0.001</b>	0.002 $\pm$ 0.000	<b>0.006 <math>\pm</math> 0.006</b>	4.36 $\pm$ 0.09	4.36 $\pm$ 0.09	<b>0.014 <math>\pm</math> 0.001</b>	<b>0.025 <math>\pm</math> 0.001</b>	<b>0.005 <math>\pm</math> 0.005</b>	7.31 $\pm$ 0.07	7.31 $\pm$ 0.07

- Our empirical results support the scaling hypothesis: even the standard DB in  $\Delta \mathcal{F}$  scale is stable; however, non-acyclic GFlowNets trained with  $\Delta \mathcal{F}$  scale losses often fail to accurately match the reward distribution;
- Both DB and SDB in  $\Delta \log \mathcal{F}$  scale result in better matching the reward distribution but need to be utilized with state flow regularization to ensure small expected trajectory length  $\mathbb{E}[n_\tau]$ .
- Learning with a fixed  $\mathcal{P}_B$  is possible without stable losses and regularization, however, manually picking  $\mathcal{P}_B$  with small  $\mathbb{E}[n_\tau]$  is challenging;
- $\lambda$  is crucial, the optimization task is very sensitive to the choice of this parameter, changing the order of  $\lambda$  can significantly impact the result by extending or shrinking of mean trajectory length.