# METAGFN: EXPLORING DISTANT MODES WITH ADAPTED METADYNAMICS FOR CONTINUOUS GFLOWNETS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Generative Flow Networks (GFlowNets) are a class of generative models that sample objects in proportion to a specified reward function through a learned policy. They can be trained either on-policy or off-policy, needing a balance between exploration and exploitation for fast convergence to a target distribution. While exploration strategies for discrete GFlowNets have been studied, exploration in the continuous case remains to be investigated, despite the potential for novel exploration algorithms due to the local connectedness of continuous domains. Here, we introduce Adapted Metadynamics, a variant of metadynamics that can be applied to arbitrary black-box reward functions on continuous domains. We use Adapted Metadynamics as an exploration strategy for continuous GFlowNets. We show several continuous domains where the resulting algorithm, MetaGFN, accelerates convergence to the target distribution and discovers more distant reward modes than previous off-policy exploration strategies used for GFlowNets.

## 1 INTRODUCTION

027 028

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

029 Generative Flow Networks (GFlowNets) are a type of generative model that samples from a discrete space  $\mathcal{X}$  by sequentially constructing objects via actions taken from a learned policy  $P_F$  (Bengio et al., 2021a). The policy  $P_F(s, s')$  specifies the probability of transitioning from some state s to 031 some other state s'. The policy is parameterised and trained so that, at convergence, the probability of sampling an object  $x \in \mathcal{X}$  is proportional to a specified reward function R(x). GFlowNets offer 033 advantages over more traditional sampling methods, such as Markov chain Monte Carlo (MCMC), 034 by learning an amortised sampler, capable of single-shot generation of samples from the desired distribution. Since GFlowNets learn a parametric policy, they can generalise across states, resulting in higher performance across various tasks (Bengio et al., 2021a; Malkin et al., 2022; Zhang et al., 2022; 037 Jain et al., 2022; Deleu et al., 2022; Jain et al., 2023; Hu et al., 2023; Zhang et al., 2023; Shen et al., 038 2023b) and applications to conditioned molecule generation (Shen et al., 2023b), maximum likelihood estimation in discrete latent variable models (Hu et al., 2023), structure learning of Bayesian networks (Deleu et al., 2022), scheduling computational operations (Zhang et al., 2023), and discovering 040 reticular materials for carbon capture (Cipcigan et al., 2023). 041

Although originally conceived for discrete state spaces, GFlowNets have been extended to more general state spaces, such as entirely continuous spaces, or spaces that are hybrid discrete-continuous (Lahlou et al., 2023). In the continuous setting, given the current state, the policy specifies a continuous probability distribution over subsequent states, and the probability density over states  $x \in \mathcal{X}$  sampled with the policy is proportional to a reward density function r(x). The continuous domain unlocks more applications for GFlowNets, such as molecular conformation sampling (Volokhova et al., 2023) and continuous control problems (Luo et al., 2024).

GFlowNets are trained like reinforcement learning agents. Trajectories of states are generated either on-policy or off-policy, with the terminating state  $x \in \mathcal{X}$  providing a reward signal for informing a gradient step on the policy parameters. GFlowNets therefore suffer from the same training pitfalls as reinforcement learning. One such issue is slow temporal credit assignment, which has thus far been addressed by designing more effective loss functions, such as detailed balance (Bengio et al., 2021b), trajectory balance (Malkin et al., 2022) and sub-trajectory balance (Madan et al., 2022). This latter approach has recently been extended by providing an energy function inductive bias at the early stages of training (Pan et al., 2023).

Besides loss functions, another aspect of GFlowNet training is the exploration strategy for acquiring training samples. Exclusively on-policy learning is generally inadequate as it leads to inefficient exploration of new modes. More successful strategies therefore rely on off-policy exploration. For the discrete setting, numerous exploration strategies have been proposed including  $\epsilon$ -noisy with a uniform random policy, tempering, Generative Augmented Flow Networks (GAFN) (Pan et al., 2022), Thompson sampling (Rector-Brooks et al., 2023) and Local Search GFlowNets (Kim et al., 2024). While these approaches can be generalised to the continuous domain, there is limited literature benchmarking their effectiveness in this setting.

064 Sampling in the continuous setting is a common occurrence in various domains such as molecular 065 modelling (Hawkins, 2017; Yang et al., 2019) and Bayesian inference (Shahriari et al., 2016). 066 The local connectedness of a continuous domain allows for novel exploration strategies that are 067 not directly applicable in the discrete setting. Sendera et al. (2024) compared several exploration 068 strategies in the context of diffusion samplers, and proposed alternating between on-policy sampling 069 and backward sampling from a fixed set of MCMC samples when training continuous GFlowNets. These samples were selected prior to training, thereby inducing a fixed pre-training cost but keeping 071 the cost of each training trajectory constant. The authors showed effective exploration, even in high-dimensional problems. However, their MCMC approach requires access to cheap gradients 072 of the reward landscape, which might not always be readily available (Rengarajan et al., 2022). 073 In addition, there are some settings, such as in molecular conformation sampling, where MCMC 074 approaches are known to require significantly longer timescales to overcome energy barriers than 075 methods based on molecular dynamics (MD) simulations (Abrams and Bussi, 2014), thereby making 076 the pre-training cost prohibitively expensive. 077

In this work, we propose MetaGFN, a novel exploration algorithm for continuous GFlowNets inspired by metadynamics, an enhanced sampling method widely used for molecular modelling (Laio and 079 Parrinello, 2002). Unlike MCMC-based sampling methods, MetaGFN operates in a general black-box setting, relying solely on an oracle for reward values, without requiring access to reward gradients. In 081 common with standard metadynamics, MetaGFN converges quickly when a reduced-dimensional representation of the reward measure, specified by a so-called collective variable (CV) basis, is known. 083 While our method is domain-agnostic, this feature makes MetaGFN particularly well-suited for small 084 molecular systems, where such representations are known and often low-dimensional (Fiorin et al., 085 2013). MetaGFN requires no pre-training and only introduces a small, constant cost for each training 086 trajectory, which becomes negligible as the cost of reward evaluation increases.

- 087088 The main contributions of this work are:
  - Introducing MetaGFN, an algorithm that adapts metadynamics to black box rewards and continuous GFlowNets;
  - Proving that the Adapted Metadynamics formulation underlying MetaGFN is consistent and reduces to standard metadynamics in the appropriate limit;
  - Showing empirically that MetaGFN outperforms existing GFlowNets exploration strategies in various continuous environments, including alanine dipeptide conformation sampling.

The rest of the paper is as follows. In Section 2 we review the theory of discrete and continuous GFlowNets as well as metadynamics and collective variables. In Section 3, we present the Adapted Metadynamics and MetaGFN algorithms. In Section 4, we evaluate MetaGFN against other approaches, showing that MetaGFN generally outperforms existing exploration strategies in various continuous environments. We finish with limitations and conclusions in Sections 5 and 6. Code for MetaGFN is available at [link in camera-ready].

102 103

104

106

091

092

094

- 2 PRELIMINARIES
- 105 2.1 DISCRETE GFLOWNETS
- In a GFlowNet, the *network* refers to a directed acyclic graph (DAG), denoted as G = (S, A). Nodes represent *states*  $s \in S$ , and edges represent *actions*  $s \to s' \in A$  denoting one-way transitions between

states. The DAG has two distinguishable states: a unique *source state*  $s_0$ , that has no incoming edges, and a unique *sink state*  $\perp$ , that has no outgoing edges.

The set of states,  $\mathcal{X} \subset \mathcal{S}$ , that are directly connected to the sink state are known as *terminating* 111 states. GFlowNets learn forward transition probabilities, known as a forward policy  $P_F(s'|s)$ , along 112 the edges of the DAG so that the resulting marginal distribution over the terminal states complete 113 trajectories (the *terminal distribution*), denoted as  $P^{\perp}(x)$ , is proportional to a given *reward function* 114  $R: \mathcal{X} \to \mathbb{R}$ . GFlowNets also introduce additional learnable objects, such as a *backward policy* 115  $P_B(s|s')$ , which is a distribution over the parents of any state of the DAG, to create losses that 116 train the forward policy. Objective functions for GFlowNets include flow matching (FM), detailed 117 balance (DB), trajectory balance (TB) and subtrajectory balance (STB) (Bengio et al., 2021a;b; Malkin et al., 2022; Madan et al., 2022). During training, the parameters of the flow objects are 118 updated with stochastic gradients of the objective function applied to batches of trajectories. These 119 trajectory batches can be obtained either directly from the current forward policy or from an alternative 120 algorithm that encourages exploration. These approaches are known as *on-policy* and *off-policy* 121 training respectively. 122

122

137

138 139

140

141

142

143 144

145

146

147

## 124 2.2 CONTINUOUS GFLOWNETS

Continuous GFlowNets extend the generative problem to continuous spaces (Lahlou et al., 2023), where the analogous quantity to the DAG is a measurable pointed graph (MPG) (Nummelin, 1984).
MPGs can model continuous spaces (e.g., Euclidean space, spheres, tori), as well as hybrid spaces, with a mix of discrete and continuous components, as often encountered in robotics, finance, and biology (Bortolussi and Policriti, 2008; Swiler et al., 2012; Neunert et al., 2020).

**Definition 2.1** (Measurable pointed graph (MPG)). Let  $(\bar{S}, T)$  be a topological space, where  $\bar{S}$  is the *state space*, T is the set of open subsets of  $\bar{S}$ , and  $\Sigma$  is the Borel  $\sigma$ -algebra associated with the topology of  $\bar{S}$ . Within this space, we identify the *source state*  $s_0 \in \bar{S}$  and *sink state*  $\perp \in \bar{S}$ , both distinct and isolated from the rest of the space. On this space we define a *reference transition kernel*  $\kappa : \bar{S} \times \Sigma \to [0, +\infty)$  and a *backward reference transition kernel*  $\kappa^b : \bar{S} \times \Sigma \to [0, +\infty)$ . The support of  $\kappa(s, \cdot)$  are all open sets accessible from *s*. The support of  $\kappa^b(s, \cdot)$  are all open sets where *s* is accessible from. Additionally, these objects must be well-behaved in the following sense:

- (i) Continuity: For all  $B \in \Sigma$ , the mapping  $s \mapsto \kappa(s, B)$  is continuous.
- (ii) No way back from the source: The backward reference kernel has zero support at the source state, i.e. for all  $B \in \Sigma$ ,  $\kappa^b(s_0, B) = 0$ .
- (iii) No way forward from the sink: When at the sink, applying the forward kernel keeps you there, i.e.  $\kappa(\perp, \cdot) = \delta_{\perp}(\cdot)$ , where  $\delta_{\perp}$  is the Dirac measure of the sink state.
- (iv) A fully-explorable space: The number of steps required to possibly reach any measurable  $B \in \Sigma$  from the source state, and to guarantee to reach the sink state, with the forward reference kernel is bounded.
- 148 The set of objects  $(\bar{S}, T, \Sigma, s_0, \bot, \kappa, \kappa^b)$  then defines an MPG.

Note that the support of  $\kappa(s, \cdot)$  and  $\kappa^b(s, \cdot)$  are analogous to the child and parent sets of a state s in a DAG. Similarly, a discrete GFlowNet's DAG satisfies discrete versions of (ii), (iii), and (iv).

The set of *terminating states*  $\mathcal{X}$  are the states that can transition to the sink, given by  $\mathcal{X} = \{s \in S : \kappa(s, \{\bot\}) > 0\}$ , where  $S := \overline{S} \setminus \{s_0\}$ . *Trajectories*  $\tau$  are sequences of states that run from source to sink,  $\tau = (s_0, \ldots, s_n, \bot)$ . The *forward Markov kernel*  $P_F : \overline{S} \times \Sigma \to [0, \infty)$  and *backward Markov kernel*  $P_B : \overline{S} \times \Sigma \to [0, \infty)$  have the same support as  $\kappa(s, \cdot)$  and  $\kappa^b(s, \cdot)$  respectively, where being a Markov kernel means states are mapped to probability measure, hence  $\int_{\overline{S}} P_F(s, ds') = \int_{\overline{S}} P_B(s, ds') = 1$ . A *flow* F is a tuple  $F = (f, P_F)$ , where  $f : \Sigma \to [0, \infty)$  is a *flow measure*, satisfying  $f(\{\bot\}) = f(s_0) = Z$ , where Z is the *total flow*.

The *reward measure* is a positive and finite measure R over the terminating states  $\mathcal{X}$ , we denote the density of this reward measure as r(x), for  $x \in \mathcal{X}$ . A flow F is said to satisfy the *reward-matching conditions* if

$$R(dx) = f(dx)P_F(x, \{\bot\})$$

162 If a flow satisfies the reward-matching conditions and trajectories are recursively sampled from 163 the Markov kernel  $P_F$  starting at  $s_0$ , the resulting measure over terminating states,  $P^{\perp}(B)$ , is 164 proportional to the reward:  $P^{\perp}(B) = \frac{R(B)}{R(\mathcal{X})}$  for any B in the  $\sigma$ -algebra of terminating states (Lahlou 165 et al., 2023). 166

Objective functions for discrete GFlowNets generalise to continuous GFlowNets. However, in the 167 continuous case, the forward policy  $\hat{p}_F : S \times \bar{S} \to [0, \infty)$ , backward policy  $\hat{p}_B : S \times \bar{S} \to [0, \infty)$ 168 and *parameterised flow*  $\hat{f}: S \to [0,\infty)$  parameterise the  $P_F, P_B$  transition kernels and flow measure f on an MPG. Discrete GFlowNets parameterise log transition probabilities and flows on a DAG. 170 In this work, we consider DB, TB and STB losses. For a complete trajectory  $\tau$ , the TB loss can be 171 written as 172  $L_{TB}(\tau) = \left(\log \frac{Z_{\theta} \prod_{t=0}^{n} \hat{p}_F(s_t, s_{t+1}; \theta)}{r(s_n) \prod_{t=0}^{n-1} \hat{p}_B(s_{t+1}, s_t; \theta)}\right)^2,$ 

173

where  $Z_{\theta}$  is the parameterised total flow (see Appendix A for the DB and STB loss functions).

176 177 178

2.3 EXPLORATION STRATEGIES FOR GFLOWNETS

179 GFlowNets can reliably learn using off-policy trajectories, a key advantage over hierarchical varia-180 tional models (Malkin et al., 2023). For optimal training, it is common to use a replay buffer and 181 alternate between on-policy and off-policy (exploration) batches (Shen et al., 2023a). In the discrete 182 case, proposed techniques to encourage exploration include  $\epsilon$ -noisy exploration, tempering and the incorporation of intermediate rewards (Bengio et al., 2021a; Pan et al., 2022). Exploration strategies 183 for continuous GFlowNets are less well studied in the literature but several methods designed for 184 discrete GFlowNets can be adapted (Sendera et al., 2024). In this work, we consider: 185

186 Local Search GFlowNets (Kim et al., 2024): Explores by backtracking and resampling on-policy 187 trajectories. Reconstructed trajectories with a higher reward than the original are used for training, 188 thus encouraging the learning of high-reward modes.

189 Thompson sampling. (Rector-Brooks et al., 2023): Explores high-uncertainty regions by using an 190 ensemble of policy heads with a shared torso. A random head generates the on-policy trajectory, 191 and the loss is computed by averaging contributions over heads, where each head is independently 192 included with probability p.

193 **Noisy exploration**: Explores by increasing policy uncertainty. A small constant is added to the policy 194 variance which is gradually reduced to zero over the course of training. 195

196 Nested sampling (Lemos et al., 2023): A Markov-Chain Monte Carlo (MCMC) algorithm is used to sample from the reward distribution. Backward sampling from these terminal states are used to 197 generate off-policy trajectories. 198

199 200

### 2.4 METADYNAMICS AND COLLECTIVE VARIABLES

201 Molecular dynamics (MD) simulates that dynamics of molecules using Langevin dynamics (LD) 202 (Pavliotis, 2014), a stochastic differential equation that models particle motion under friction and 203 random noise. LD trajectories ergodically sample the molecule's Gibbs measure,  $\rho_{\beta}(x) \propto e^{-\beta V(x)}$ , 204 where  $x \in \mathcal{X}$  represents atomic positions, V(x) is the molecular potential, and  $\beta$  is the thermo-205 dynamic beta.<sup>1</sup> However, when V(x) contains multiple, deep local minima, as is common in 206 biomolecules, then LD becomes inefficient, as it tends to get trapped in these minima, slowing 207 exploration of the full state space.

208 *Metadynamics* overcomes this by progressively modifying the potential landscape to discourage 209 visits to already-explored regions (Laio and Parrinello, 2002). It achieves this by regularly depositing 210 repulsive Gaussian biases, centered at the current state of the evolving LD trajectory. This modifies 211 the potential to  $V_{\text{total}}(x,t) = V(x) + V_{\text{bias}}(x,t)$ , where  $V_{\text{bias}}(x,t)$  is the cumulative bias at time t. 212 Intuitively, metadynamics thus progressively transforms the landscape into more level surface, in 213 which the system can diffuse more freely, thereby accelerating exploration (Figure 1). In the limit 214  $t \to \infty$ , the dynamics approximates free diffusion and uniformly samples the domain of V(x).

<sup>&</sup>lt;sup>1</sup>We review Langevin dynamics in Appendix B.



Figure 1: Illustration of metadynamics in a multi-well potential. Regular deposition of a bias leads to a total potential that gradually flattens, encouraging exploration.

230 Since biases are typically specified on a numerical grid, this gives rise to a memory cost that 231 grows exponentially with dimension. Therefore, biases are typically applied along low-dimensional 232 coordinates known as *collective variables* (CVs). A collective variable  $z(x) : \mathcal{X} \to \mathcal{Z}$  is any mapping 233 from the original (high-dimension) state space  $\mathcal{X}$  to a lower-dimensional space  $\mathcal{Z}$ . For potential V(x)234 and collective variables z, metadynamics is guaranteed to eventually uniformly sample the domain of 235 the marginal potential  $V(z') := \int_{\mathcal{X}} \delta(z' - z(x)) V(x) dx$ . In molecular contexts, ideal choices of 236 CVs correspond to the physical reaction coordinates that underly rare-event transitions (e.g. backbone dihedrals, interatomic distances) (Laio and Gervasio, 2008; De Vivo et al., 2016). In more general 237 settings, ideal CVs should resolve rare events dynamics on V(x), simplify the landscape (retaining 238 its essential features e.g. barriers and basins), and minimise the dimensionality of this representation. 239 For simple systems, CVs can be derived from known symmetries or macroscopic order parameters 240 that describe state changes. However, if the choice of CVs is not obvious, then they can learnt through 241 data-driven methods such as Time-Lagged Independent Component Analysis (TICA) (Molgedey and 242 Schuster, 1994), manifold learning algorithms, neural network autoencoders or variational methods 243 (Mardt et al., 2018; Bonati et al., 2021; Ramaswamy et al., 2021). For an in-depth review on these 244 CV-learning approaches, see Sidky et al. (2020). Metadynamics has seen numerous extensions (Bussi 245 and Laio, 2020). In the next section, we adapt the original metadynamics algorithm to the continuous 246 black box setting. We will assume suitable CVs are given, but the theory we present is agnostic to 247 their form-be it analytical expressions, neural networks, or non-parametric tabular mappings.

248

227

228 229

## 249 250

## 3 METAGFN: ADAPTED METADYNAMICS FOR GFLOWNETS

251 Metadynamics has two properties that make it well-suited as an exploration strategy. Firstly, it 252 eventually uniformly samples the domain, ensuring exploration of *all* local minima. In contrast, 253 techniques such as Local Search GFlowNets, Thompson sampling, and noisy exploration rely on 254 localised strategies that are prone to mode locking (Section 4). Secondly, its gradual diffusion from 255 the starting configuration allows incremental exploration, resulting in more stable training than global 256 approaches like nested sampling. However, metadynamics requires a potential and its gradients, 257 which are not always available in a black-box GFlowNet setting. We adapt metadynamics to this setting by interpreting the reward landscape as a potential energy surface, resulting in the exploration 258 algorithm Adapted Metadynamics (AM) (Section 3.1) and training algorithm MetaGFN (Section 3.2). 259

260 **Interpreting reward density as a potential:** We assume that  $\mathcal{X}$  is a manifold and that the reward 261 density r(x) is bounded and L1-integrable over  $\mathcal{X}$  with at most finitely many discontinuities. Thus, 262 the target density over terminal states,  $\rho(x) := r(x) / \int_{\mathcal{X}} r(x') dx'$ , can be expressed as a Gibbs 263 distribution:  $\rho(x) = \exp(-\beta' V(x)) / \int_{\mathcal{X}} r(x') dx'$ , where we identify  $V(x) = -\frac{1}{\beta'} \ln r(x)$  as the 264 corresponding potential, for some constant scalar  $\beta' > 0$ . Reward maxima thus correspond to 265 potential energy minima. From simplicity, we set  $\beta'$  equal to the thermodynamic beta ( $\beta$ ). Making 266 this choice means that the single parameter  $\beta$  uniquely controls the (unbiased) transition rates between minima of the potential<sup>2</sup>. Our goal is to use metadynamics to explore V(x), generating 267 268 high-reward terminal states that guide the GFlowNet to sample all reward density modes.

<sup>269</sup> 

<sup>&</sup>lt;sup>2</sup>From the Kramer formula; transition rate  $\propto \frac{1}{\beta} \exp(\beta \Delta V)$ , where  $\Delta V \propto \frac{1}{\beta'} = \frac{1}{\beta}$ .

293 294 295

296

323

Algorithm 1: Adapted Metadynamics : Manifold environment of terminating states  $\mathcal{X}$  with reward density  $r : \mathcal{X} \to \mathbb{R}$ . Input Initial state  $(x_t, p_t) \in \mathcal{X} \times T_{x_t}(\mathcal{X})$ . Collective variables  $z = (z_1, \ldots, z_d)$ . 273 **Parameters**: Gaussian width  $\sigma = (\sigma_1, \ldots, \sigma_d) \in \mathbb{R}^d$ . Gaussian height w > 0. Stride  $n \in \mathbb{Z}^+$ . 274 LD parameters:  $\gamma, \beta$ . Timestep  $\Delta t$ . 275  $\hat{N} \leftarrow 0$ 276  $\hat{R} \leftarrow 0$ 277  $\hat{V}(z) \leftarrow 0$ 278 4  $V_{\text{bias}}(z) \leftarrow 0$ 279 5 every timestep  $\Delta t$ : 280  $z_t \leftarrow z(x_t)$ 6 281 every n timesteps  $n\Delta t$ : 7 282  $\hat{N} \leftarrow \hat{N} + \exp\left(-\frac{1}{2}\sum_{i=1}^{d}\frac{(z_i - z_{i,t})^2}{\sigma_i^2}\right)$ 8  $\hat{R} \leftarrow \hat{R} + r(x_t) \cdot \exp\left(-\frac{1}{2} \sum_{i=1}^{d} \frac{(z_i - z_{i,t})^2}{\sigma_i^2}\right)$ 284 Q 285  $\hat{V} \leftarrow -\frac{1}{\beta} \log \left( \hat{R} / (\hat{N} + \epsilon) + \epsilon \right)$ 286 10 287  $V_{\text{bias}}(z) \leftarrow V_{\text{bias}}(z) + n \cdot \Delta t \cdot w \cdot \exp\left(-\frac{1}{2}\sum_{i=1}^{d} \frac{(z_i - z_{i,t})^2}{\sigma_i^2}\right)$ 11 288 compute forces: 289 12  $F \leftarrow -\left(\nabla_z \hat{V}(z)|_{z=z_t} + \nabla_z V_{\text{bias}}(z)|_{z=z_t}\right) \cdot \nabla_x z|_{x=x_t}$ 290 13 291 14 292 **propagate**  $x_t, p_t$  by  $\Delta t$  using Langevin dynamics with computed force F (Alg. 3, Appendix B).

#### 3.1 ADAPTED METADYNAMICS

297 Metadynamics requires the gradient of the total potential, where  $-\nabla V_{\text{total}}(x,t) = -\nabla (V(x) +$ 298  $V_{\text{bias}}(x,t)$ ). Using the above assumptions, we have  $\nabla V(x) = -\nabla r(x)/(\beta' r(x))$ . However, r(x) is 299 often a computationally expensive black-box function, and its gradient,  $\nabla r(x)$ , is unknown. While 300 finite differences can estimate  $\nabla r(x)$  for smooth, low-dimensional reward distributions, this approach 301 is impractical in high-dimensional spaces. Below, we explain how to adapt metadynamics to this black-box setting and avoid finite difference gradient estimates by storing a dynamically-updated 302 kernel density estimate (KDE) of the potential. We call the modified metadynamics algorithm 303 Adapted Metadynamics (Algorithm 1). 304

305 Let  $x_t$  denote the metadynamics sample at time t and  $z(x) = (z_1(x), \ldots, z_d(x))$  be a given set of 306 collective variables, where  $z_i$  is a one-dimensional coordinate, and  $z \in \mathcal{Z}$  is d-dimensional. Let  $z_{i,t} := z_i(x_t)$  denote the corresponding  $i^{th}$  CV coordinate at time t and assume that the Jacobian  $\nabla_x z$  is well-defined. We store a discretisation of the KDE marginal potential and bias potentials 307 308 in CV space  $\mathcal{Z}$ . Since these are stored in memory, gradient computations are cheap and require no 309 further evaluations of r(x). 310

311 Let  $\hat{V}(z,t)$  represent the KDE of the marginal potential at time t. To compute  $\hat{V}(z,t)$ , we maintain 312 two separate KDEs: N(z,t) for visited states (line 8) and R(z,t) for cumulative rewards (line 9). 313 We update these KDEs on-the-fly at the same time the bias potential is updated, which occurs every 314 integer n steps of Langevin dynamics (line 7). If  $\mathcal{Z} \cong \mathbb{R}^d$ , we use Gaussian kernels with kernel width 315  $\sigma = (\sigma_1, \ldots, \sigma_k) \in \mathbb{R}^k$ , matching the width of the Gaussian bias (line 11)<sup>3</sup>. This is a reasonable 316 as both are set by the variability length scale of V(x). Finally, the KDE potential  $\hat{V}(z,t)$  is then computed as  $\hat{V}(z,t) = -\frac{1}{\beta} \log \left( \frac{\hat{R}(z,t)}{\hat{N}(z,t)+\epsilon} + \epsilon \right)$ , for a fixed constant  $\epsilon > 0$  (line 10). We found 317 318 319 empirically that  $\epsilon$  ensured numerical stability by preventing division by zero and bounding V from 320 above. Defining  $\hat{V}$  through the ratio  $\hat{R}/\hat{N}$  ensures that it rapidly and smoothly adjusts whenever new 321 modes are discovered. Furthermore, we prove that  $\hat{V}$  eventually discovers all reward modes in the CV space. More precisely, 322

<sup>&</sup>lt;sup>3</sup>If  $\mathcal{Z} \cong \mathbb{T}^d$  (*d*-torus), we use von Mises distributions.

**Theorem 3.1.** If the collective variable z(x) is analytic with a bounded domain, then

$$\lim_{\epsilon \to 0} \left( \lim_{\sigma \to 0} \left( \lim_{t \to \infty} \hat{V}(z, t) \right) \right) = V, \tag{1}$$

where  $V = V(z') := \int_{\mathcal{X}} \delta(z' - z(x)) V(x) dx.$ 

The proof is in Appendix C.

Note that the algorithm can be extended to a batch of trajectories, where each metadynamics trajectory evolves independently, but with a shared  $\hat{V}$  and  $V_{\text{bias}}$  which receive updates from every trajectory in the batch. We found empirically that this accelerates exploration and reduces stochastic gradient noise during training and this is the version we use in our experiments (Section 4).

## 3.2 MetaGFN

324

325 326

327

331

332

333

334 335

336

349

350

351

352

357

375

337 Each Adapted Metadynamics sample  $x_i \in \mathcal{X}$  is an off-policy terminal state sample. To train a 338 GFlowNet, complete trajectories are required. We generate these by backward sampling from the 339 terminal state, giving a trajectory  $\tau = (s_0, s_1, \dots, s_n = x_i)$ , where each state  $s_{i-1}$  is sampled from 340 the current backward policy distribution  $\hat{p}_B(s_{i-1}|s_i;\theta)$ , for *i* from *n* to 1. This approach means that 341 the generated trajectory  $\tau$  has reasonable credit according to the loss function, thereby providing a 342 useful learning signal. However, since this requires a backward policy, this is compatible with DB, TB, and STB losses, but not FM loss. Given the superior credit assignment of the former losses, this 343 is not a limitation (Madan et al., 2022). 344

Additionally, we use a replay buffer. Due to the theoretical guarantee that Adapted Metadynamics
 will eventually sample all collective variable space (Theorem 3.1), AM samples are ideal candidates
 for storing in a replay buffer. When storing these trajectories in the replay buffer, there are two
 obvious choices:

- 1. Store the entire trajectory the first time it is generated;
- 2. Store only the Adapted Metadynamics sample and regenerate trajectories using the current backward policy when retrieving from the replay buffer.

We investigated both options in preliminary experiments (Appendix D.2). Option 2 was found to be uniformly superior and it is the version we use in our main experiments (Section 4). We call the overall training algorithm *MetaGFN*, with pseudocode presented in Algorithm 2, below.

Algorithm 2: MetaGFN

358 Input : Forward policy  $P_F$ . Backwards policy  $P_B$ . Loss function L. 359 Parameters: How often to run Adapted Metadynamics batches, freqMD. How often to run 360 replay buffer batches, freqRB. Batch size, b. Stride,  $n \in \mathbb{Z}^+$ . Time step,  $\Delta t > 0$ 361 1 for each episode do: 362 if episode number is divisible by freqMD: 2 Run Adapted Metadynamics (batch size b) for time  $n\Delta t$ , obtain samples  $\{x_1, \ldots, x_b\}$ 3 Push  $\{x_1, \ldots, x_b\}$  to the replay buffer 364 4 Backward sample from  $\{x_1, \ldots, x_b\}$  using current  $P_B$  to obtain trajectories  $\{\tau_1, \ldots, \tau_b\}$ 5 365 elif episode number is divisible by freqRB: 6 366 Random sample  $\{x_1, \ldots, x_b\}$  from the replay buffer 7 367 Backward sample from  $\{x_1, \ldots, x_b\}$  using current  $P_B$  to obtain trajectories  $\{\tau_1, \ldots, \tau_b\}$ 8 368 9 else: 369 Generate trajectories  $\{\tau_1, \ldots, \tau_b\}$  on-policy 10 370 Compute loss  $l = \sum_{i=1}^{b} L(\tau_i, P_F, P_B)$ 11 371 Take gradient step on loss l12 372 373 374

4 EXPERIMENTS

We compare MetaGFN with Thompson sampling, noisy, Local Search GFlowNets and nested sampling (Section 2.3). We run experiments in five continuous environments, summarised below. For

each exploration strategy, we use a replay buffer and alternate between exploration and replay buffer
batches. For MetaGFN, we use freqRB = 2, freqMD = 10. Forward and backward kernels are
Gaussian/von Mises mixture distributions, with distribution parameters specified at each state by an
MLP. In all environments, the additional computational expense of running Adapted Metadynamics
was negligible (<5%) compared to the training time of the models. Full experimental details can be</li>
found in Appendix D.

**Line environment**: A one-dimensional environment with state space  $\mathcal{S} = \mathbb{R} \times \{t \in \mathbb{N}, 1 \le t \le 3\}$ , where *t* indexes the position of a state in a trajectory. The source state is  $s_0 = (0, 0)$ . The terminal states are therefore  $\mathcal{X} = \mathbb{R} \times \{3\} \cong \mathbb{R}$ . The collective variable is the identity, i.e. z = x. The reward density, plotted in Figure 2, consists of an asymmetric bimodal peak near the origin and an additional distant lone peak. It is given by the Gaussian mixture distribution:

392

393 394 395

396

where  $\mathcal{N}(\mu, \sigma^2)$  is a Gaussian density with mean  $\mu$  and variance  $\sigma^2$ .

 $r(x) = \begin{cases} \mathcal{N}(-2.0, 1.0) + \mathcal{N}(-2.0, 0.4) + \\ \mathcal{N}(2.0, 0.6) + \mathcal{N}(20.0, 0.1); & -5 \le x \le 23 \\ 0; & \text{otherwise}, \end{cases}$ (2)

Alanine dipeptide environment: One application of continuous 397 GFlowNets is molecular conformation sampling (Volokhova et al., 398 2023). Here, we train a GFlowNet to sample conformational states of 399 alanine dipeptide (AD), a a small biomolecule of 23 atoms that plays 400 a key role in modelling the dynamics of proteins (Hermans, 2011). 401 The metastable states of AD are distinguished in a 2D CV space 402 defined by the backbone dihedral angles  $\phi$  and  $\psi$ . The resulting free 403 energy surface  $V(\phi, \psi)$  in explicit water, obtained after extensive 404 sampling with long molecular dynamics simulations, is shown in Figure 3. The metastable states, in increasing energy, are  $P_{||}$ ,  $\alpha_R$ ,  $C_5$ , 405  $\alpha', \alpha_L$ , and  $\alpha_D$ . The state space is  $S = \mathbb{T}^2 \times \{t \in \mathbb{N}, 1 \le t \le 3\}$ , with source state  $s_0 = P_{||} = (-1.2, 2.68)$ . Terminal states are  $\mathcal{X} = \mathbb{T}^2 \times \{3\} \cong \mathbb{T}^2$ . The reward function is the Boltzmann weight,  $r(\phi, \psi) = \frac{1}{Z} \exp(-\beta V(\phi, \psi))$ , where Z is the normalisation con-406 407 408 409 stant. 410

411 Grid environments: We consider (hyper)grids in d = 2, 3 and 4 412 dimensions. The state space is  $S = [-15, 15]^d \times \{t \in \mathbb{N}, 1 \le t \le 3\}.$ 413 The k-dimensional (hyper)grid consists of  $2^k$  modes, located at the 414 corners of the (hyper)cube  $[-10, 10]^d$ . The (hyper)grid is centered on the origin with an edge width of 20. The reward modes are 415 Gaussians with variance  $\sigma^2 = 2$  (Figure 4). We use trajectory 416 lengths of 3, 5, and 6 for dimensions 2, 3, and 4 respectively. The 417 collective variable is the identity, i.e. z = x. The source state is the 418 origin and terminal states are in  $\mathcal{X} \cong \mathbb{R}^k$ . 419



```
4.1 RESULTS
```

In each environment, we run experiments with three different loss functions: Detailed Balance (DB),
Trajectory Balance (TB) and Subtrajectory Balance (STB). We evaluate performance by computing
the L1 error between the known reward distribution and the empirical on-policy distribution during
training resulting from 10<sup>4</sup> independent samples. The results, averaged over 10 repeats, are shown in
Figure 5. In 12 out of 15 of the experiments, MetaGFN outperformed all other exploration strategies.
We give further analysis below.

*Line Environment*: Among the losses, TB shows the lowest variance, and for all losses, MetaGFN always converges to a lowest error <sup>4</sup>. Indeed, MetaGFN is the only method that consistently samples



Figure 2: Line environment reward density.



Figure 3: Free energy surface of alanine dipeptide.



Figure 4: Grid environment reward density in dimension 2.

<sup>&</sup>lt;sup>4</sup>We do not show results for nested sampling on the line environment as the implementation used did not support one-dimensional environments, see Appendix D.



Figure 5: The L1 difference between on-policy and reward distribution during training for different loss functions and exploration strategies. The mean is plotted with standard error over 10 repeats.
DB - Detailed Balance loss, TB - Trajectory Balance loss, STB - Subtrajectory Balance loss.

452

453

456 457

the distant reward peak at x = 20 (Appendix D.2). Although other strategies occasionally sampled 458 the distant peak, MetaGFN alone converges because it *continuously* samples this peak, even if the 459 forward policy starts to lock onto the central modes, thus ensuring that the replay buffer is always 460 populated with diverse samples. The slight increase in the loss of MetaGFN around batch number  $5 \times 10^3$  occurs as the on-policy distribution widens when Adapted Metadynamics first discovers 461 the distant peak. Appendix D.2 provides further analysis of Adapted Metadynamics and compares 462 different MetaGFN variants, with and without noise, and with and without trajectory regeneration. 463 We confirm that the variant of MetaGFN presented in Figure 5 (no added noise, always regenerate 464 trajectories) is the most robust. 465

Alanine Dipeptide Environment: For each loss, MetaGFN generally converges to a lower error than 466 all other exploration strategies. However, for TB loss, the average L1 error is marginally higher than 467 on-policy training, but this conceals the fact that the best-case error is smaller. To better understand 468 this result, we examined the best and worst training runs (as measured by L1 error) for TB on-policy 469 and TB MetaGFN, shown in Figure 6. Note that, unlike on-policy, the best MetaGFN run learns to 470 sample the rare  $\alpha_L$  mode. In the worst case, MetaGFN fails to converge (although this is rare; only 471 one of 10 runs failed). In Table 1, we quantify how often the different AD modes were sampled over 472 the different repeats (a mode is considered sampled if the on-policy distribution has a mode within 473 the correct basin of attraction). TB loss with MetaGFN is the only combination that consistently 474 samples the majority of modes. The only mode not sampled by any method is  $\alpha_D$ , which has a 475 natural abundance approximately 10 times less frequent than  $\alpha_L$ .

476 Grid Environments: MetaGFN achieves the lowest error in 7 of 9 experiments, although note that 477 the performance of all methods generally decreases with increasing dimension. This is unsurprising, 478 since the learning task becomes more difficult. For MetaGFN, this is also an expected outcome 479 of the curse of dimensionality of the replay buffer grid; in high-dimensional spaces metadynamics 480 samples encounter high-reward modes less frequently, leading to a replay buffer with less sample 481 diversity. This is a known limitation of metadynamics (Laio and Parrinello, 2002). Nested Sampling 482 outperforms MetaGFN in one case (Grid 3D, DB loss) but relies on costly MCMC-based pre-training that scales exponentially with dimension and requires reward gradients, unlike MetaGFN. However, 483 the generally favourable performance of MetaGFN could be due to the method discovering an overall 484 larger diversity of off-policy samples, that are refreshed and updated during training, in contrast to 485 the static samples when using Nested Sampling approach.



Table 1: Number of correct samples of AD modes in trained GFlowNets over 10 independent repeats for DB, STB, and TB loss functions. OP - On-policy and MD - MetaGFN. The  $\alpha_D$  mode wasn't sampled in any model due to its low natural frequency.

Figure 6: Learned on-policy distribution for TB on-policy and TB MetaGFN training runs. The colour bar shows the probability density. Red histograms show the marginal distribution along the angular coordinates. Black curves show the marginal distributions of the ground truth. In the best case, MetaGFN is able to learn the  $\alpha_L$  mode. In the worst case, MetaGFN fails to converge. On-policy training, although more consistent, fails to learn to sample from the  $\alpha_L$  mode.

## 5 LIMITATIONS

 For metadynamics to be an effective sampler, the CVs must be low-dimensional and bounded. Either these CVs should be known analytically (as in our experiments) or learnt from data (Sidky et al., 2020). Alternatively, CVs could be learnt adaptively by parameterising  $z(x; \theta)$  by a neural network and updating its parameters by back-propagating through the GFlowNet loss during training. The metadynamics algorithm could also be replaced with a variant with smoother convergence properties, such as well-tempered metadynamics (Barducci et al., 2008) or on-the-fly probability enhanced sampling (OPES) (Invernizzi, 2021). We leave these as extensions for future work.

- 6 CONCLUSIONS

While exploration strategies for discrete Generative Flow Networks (GFlowNets) have received extensive attention, the methodologies for continuous GFlowNets remain relatively underexplored. To address this gap, we illustrated how metadynamics, a widely used enhanced sampling technique in molecular dynamics, can be adapted as an effective exploration strategy for continuous GFlowNets.

In molecular dynamics, atomic forces can be computed as the gradient of the potential, whereas continuous GFlowNets tackle problems where the reward function is a black box and gradients are inaccessible. We demonstrated how the method could be adapted by updating a kernel density estimate of the reward function on-the-fly, and proved that this is guaranteed to explore the space in an appropriate limit. Our empirical investigations show that MetaGFN offers a computationally efficient means to explore new modes in environments where prior knowledge of collective variables exists. Importantly, this work advocates an approach wherein techniques derived from molecular modelling can be adapted for machine learning tasks. Looking ahead, we anticipate that this could be a fruitful area of cross-disciplinary research, where existing ideas from the enhanced sampling literature can find further applications in generative modelling and reinforcement learning.

# 540 REFERENCES

 Cameron Abrams and Giovanni Bussi. Enhanced Sampling in Molecular Dynamics Using Metadynamics, Replica-Exchange, and Temperature-Acceleration. *Entropy*, 16(1):163–199, January 2014. ISSN 1099-4300. doi: 10.3390/e16010163. URL https://www.mdpi.com/1099-4300/ 16/1/163. Number: 1 Publisher: Multidisciplinary Digital Publishing Institute.

- Alessandro Barducci, Giovanni Bussi, and Michele Parrinello. Well-Tempered Metadynamics: A
  Smoothly Converging and Tunable Free-Energy Method. *Physical Review Letters*, 100(2):020603,
  January 2008. ISSN 0031-9007, 1079-7114. doi: 10.1103/PhysRevLett.100.020603. URL
  https://link.aps.org/doi/10.1103/PhysRevLett.100.020603.
- Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Y. Bengio. Flow Network
   based Generative Models for Non-Iterative Diverse Candidate Generation. *ArXiv*, June 2021a.
- Yoshua Bengio, T. Deleu, J. E. Hu, Salem Lahlou, Mo Tiwari, and Emmanuel Bengio. GFlowNet
   Foundations. *ArXiv*, November 2021b.
- Luigi Bonati, GiovanniMaria Piccini, and Michele Parrinello. Deep learning the slow modes for rare events sampling. *Proceedings of the National Academy of Sciences*, 118(44):e2113533118, November 2021. doi: 10.1073/pnas.2113533118. URL https://www.pnas.org/doi/ full/10.1073/pnas.2113533118. Publisher: Proceedings of the National Academy of Sciences.
- Massimiliano Bonomi, Giovanni Bussi, Carlo Camilloni, Gareth A. Tribello, Pavel Banáš, Alessandro 561 Barducci, Mattia Bernetti, Peter G. Bolhuis, Sandro Bottaro, Davide Branduardi, Riccardo Capelli, 562 Paolo Carloni, Michele Ceriotti, Andrea Cesari, Haochuan Chen, Wei Chen, Francesco Colizzi, 563 Sandip De, Marco De La Pierre, Davide Donadio, Viktor Drobot, Bernd Ensing, Andrew L. Ferguson, Marta Filizola, James S. Fraser, Haohao Fu, Piero Gasparotto, Francesco Luigi Gervasio, 565 Federico Giberti, Alejandro Gil-Ley, Toni Giorgino, Gabriella T. Heller, Glen M. Hocky, Marcella 566 Iannuzzi, Michele Invernizzi, Kim E. Jelfs, Alexander Jussupow, Evgeny Kirilin, Alessandro Laio, 567 Vittorio Limongelli, Kresten Lindorff-Larsen, Thomas Löhr, Fabrizio Marinelli, Layla Martin-568 Samos, Matteo Masetti, Ralf Meyer, Angelos Michaelides, Carla Molteni, Tetsuya Morishita, 569 Marco Nava, and The PLUMED consortium. Promoting transparency and reproducibility in 570 enhanced molecular simulations. Nature Methods, 16(8):670–673, August 2019. ISSN 1548-7091. 571 doi: 10.1038/s41592-019-0506-8.
- Luca Bortolussi and Alberto Policriti. Hybrid Systems and Biology. In Marco Bernardo, Pierpaolo
  Degano, and Gianluigi Zavattaro, editors, *Formal Methods for Computational Systems Biology*,
  pages 424–448, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg. ISBN 978-3-540-68894-5.
- Giovanni Bussi and Alessandro Laio. Using metadynamics to explore complex free-energy landscapes.
   *Nature Reviews Physics*, 2:200–212, March 2020. doi: 10.1038/s42254-020-0153-0. URL
   https://ui.adsabs.harvard.edu/abs/2020NatRP...2..200B. ADS Bibcode: 2020NatRP...2..200B.
- Flaviu Cipcigan, Jonathan Booth, Rodrigo Neumann Barros Ferreira, Carine Ribeiro dos Santos, and Mathias Steiner. Discovery of Novel Reticular Materials for Carbon Dioxide Capture using GFlowNets, October 2023. URL http://arxiv.org/abs/2310.07671. arXiv:2310.07671 [cond-mat].
- Marco De Vivo, Matteo Masetti, Giovanni Bottegoni, and Andrea Cavalli. Role of Molecular
   Dynamics and Related Methods in Drug Discovery. *Journal of Medicinal Chemistry*, 59(9):
   4035–4061, May 2016. ISSN 0022-2623. doi: 10.1021/acs.jmedchem.5b01684. URL https:
   //doi.org/10.1021/acs.jmedchem.5b01684. Publisher: American Chemical Society.
- Tristan Deleu, António Góis, Chris Emezue, Mansi Rankawat, Simon Lacoste-Julien, Stefan Bauer, and Yoshua Bengio. Bayesian Structure Learning with Generative Flow Networks, June 2022. URL http://arxiv.org/abs/2202.13903. arXiv:2202.13903 [cs, stat].
- Peter Eastman, Jason Swails, John D. Chodera, Robert T. McGibbon, Yutong Zhao, Kyle A. Beauchamp, Lee-Ping Wang, Andrew C. Simmonett, Matthew P. Harrigan, Chaya D. Stern,

Rafal P. Wiewiora, Bernard R. Brooks, and Vijay S. Pande. OpenMM 7: Rapid development 595 of high performance algorithms for molecular dynamics. *PLoS computational biology*, 13(7): 596 e1005659, July 2017. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1005659. 597 Giacomo Fiorin, Michael L. Klein, and Jérôme Hénin. Using collective variables to 598 drive molecular dynamics simulations. Molecular Physics, 111(22-23):3345-3362, December 2013. ISSN 0026-8976. doi: 10.1080/00268976.2013.813594. URL https:// 600 doi.org/10.1080/00268976.2013.813594. Publisher: Taylor & Francis \_eprint: 601 https://doi.org/10.1080/00268976.2013.813594. 602 Paul C. D. Hawkins. Conformation Generation: The State of the Art. Journal of Chemical Information 603 and Modeling, 57(8):1747–1756, August 2017. ISSN 1549-9596. doi: 10.1021/acs.jcim.7b00221. 604 URL https://doi.org/10.1021/acs.jcim.7b00221. Publisher: American Chemical 605 Society. 606 607 Jan Hermans. The amino acid dipeptide: Small but still influential after 50 years. Proceedings of the 608 National Academy of Sciences, 108(8):3095–3096, February 2011. doi: 10.1073/pnas.1019470108. 609 URL https://www.pnas.org/doi/full/10.1073/pnas.1019470108. Publisher: Proceedings of the National Academy of Sciences. 610 611 Edward J. Hu, Nikolay Malkin, Moksh Jain, Katie Everett, Alexandros Graikos, and Yoshua Bengio. 612 GFlowNet-EM for learning compositional latent variable models, June 2023. URL http:// 613 arxiv.org/abs/2302.06576. arXiv:2302.06576 [cs, stat]. 614 Michele Invernizzi. OPES: On-the-fly Probability Enhanced Sampling Method. Il Nuovo Cimento 615 C, 44(405):1–4, September 2021. ISSN 03905551, 03905551. doi: 10.1393/ncc/i2021-21112-8. 616 URL http://arxiv.org/abs/2101.06991. arXiv:2101.06991 [physics]. 617 618 Moksh Jain, Emmanuel Bengio, Alex Hernandez-Garcia, Jarrid Rector-Brooks, Bonaventure F. P. 619 Dossou, Chanakya Ajit Ekbote, Jie Fu, Tianyu Zhang, Michael Kilgour, Dinghuai Zhang, Lena 620 Simine, Payel Das, and Yoshua Bengio. Biological Sequence Design with GFlowNets. In 621 Proceedings of the 39th International Conference on Machine Learning, pages 9786–9801. PMLR, 622 June 2022. URL https://proceedings.mlr.press/v162/jain22a.html. ISSN: 2640-3498. 623 624 Moksh Jain, Sharath Chandra Raparthy, Alex Hernandez-Garcia, Jarrid Rector-Brooks, Yoshua 625 Bengio, Santiago Miret, and Emmanuel Bengio. Multi-Objective GFlowNets, July 2023. URL 626 http://arxiv.org/abs/2210.12765. arXiv:2210.12765 [cs, stat]. 627 William L. Jorgensen, Jayaraman Chandrasekhar, Jeffry D. Madura, Roger W. Impey, and Michael L. 628 Klein. Comparison of simple potential functions for simulating liquid water. The Journal of 629 Chemical Physics, 79(2):926–935, July 1983. ISSN 0021-9606. doi: 10.1063/1.445869. URL 630 https://doi.org/10.1063/1.445869. 631 632 Minsu Kim, Taeyoung Yun, Emmanuel Bengio, Dinghuai Zhang, Yoshua Bengio, Sungsoo Ahn, and Jinkyoo Park. Local Search GFlowNets, March 2024. URL http://arxiv.org/abs/2310. 633 02710. arXiv:2310.02710 [cs, stat]. 634 635 Salem Lahlou, Tristan Deleu, Pablo Lemos, Dinghuai Zhang, Alexandra Volokhova, Alex Hernández-636 García, Léna Néhale Ezzine, Yoshua Bengio, and Nikolay Malkin. A theory of continuous 637 generative flow networks. Proceedings of the 40th International Conference on Machine Learning, 638 2023. doi: 10.48550/ARXIV.2301.12594. URL https://arxiv.org/abs/2301.12594. 639 Publisher: arXiv Version Number: 2. 640 Alessandro Laio and Francesco L. Gervasio. Metadynamics: a method to simulate rare events and 641 reconstruct the free energy in biophysics, chemistry and material science. Reports on Progress 642 in Physics, 71(12):126601, November 2008. ISSN 0034-4885. doi: 10.1088/0034-4885/71/12/ 643 126601. URL https://dx.doi.org/10.1088/0034-4885/71/12/126601. 644 Alessandro Laio and Michele Parrinello. Escaping free-energy minima. Proceedings of the Na-645 tional Academy of Sciences, 99(20):12562–12566, October 2002. doi: 10.1073/pnas.202427399.

tional Academy of Sciences, 99(20):12562–12566, October 2002. doi: 10.1073/pnas.202427399.
 URL https://www.pnas.org/doi/full/10.1073/pnas.202427399.
 Publisher: Proceedings of the National Academy of Sciences.

648 Pablo Lemos, Nikolay Malkin, Will Handley, Yoshua Bengio, Yashar Hezaveh, and Laurence 649 Perreault-Levasseur. Improving Gradient-guided Nested Sampling for Posterior Inference, Decem-650 ber 2023. URL http://arxiv.org/abs/2312.03911. arXiv:2312.03911 [cs, stat]. 651 Shuang Luo, Yinchuan Li, Shunyu Liu, Xu Zhang, Yunfeng Shao, and Chao Wu. Multi-agent 652 Continuous Control with Generative Flow Networks. Neural Networks, 174:106243, June 2024. 653 ISSN 0893-6080. doi: 10.1016/j.neunet.2024.106243. URL https://www.sciencedirect. 654 com/science/article/pii/S0893608024001679. 655 656 Kanika Madan, Jarrid Rector-Brooks, Maksym Korablyov, Emmanuel Bengio, Moksh Jain, Andrei Nica, Tom Bosc, Yoshua Bengio, and Nikolay Malkin. Learning GFlowNets from partial episodes 657 for improved convergence and stability. Proceedings of the 40th International Conference on 658 Machine Learning, 2022. doi: 10.48550/ARXIV.2209.12782. URL https://arxiv.org/ 659 abs/2209.12782. Version Number: 3. 660 661 Nikolay Malkin, Moksh Jain, Emmanuel Bengio, Chen Sun, and Y. Bengio. Trajectory Balance: 662 Improved Credit Assignment in GFlowNets. ArXiv, January 2022. 663 Nikolay Malkin, Salem Lahlou, Tristan Deleu, Xu Ji, Edward Hu, Katie Everett, Dinghuai Zhang, 664 and Yoshua Bengio. GFlowNets and variational inference, March 2023. URL http://arxiv. 665 org/abs/2210.00580. arXiv:2210.00580 [cs, stat]. 666 667 Andreas Mardt, Luca Pasquali, Hao Wu, and Frank Noé. VAMPnets for deep learning of molecular kinetics. Nature Communications, 9(1):5, January 2018. ISSN 2041-668 1723. doi: 10.1038/s41467-017-02388-1. URL https://www.nature.com/articles/ 669 s41467-017-02388-1. Publisher: Nature Publishing Group. 670 671 L. Molgedey and H. G. Schuster. Separation of a mixture of independent signals using time 672 delayed correlations. Physical Review Letters, 72(23):3634–3637, June 1994. ISSN 0031-673 9007. doi: 10.1103/PhysRevLett.72.3634. URL https://link.aps.org/doi/10.1103/ 674 PhysRevLett.72.3634. 675 Michael Neunert, Abbas Abdolmaleki, Markus Wulfmeier, Thomas Lampe, Tobias Springen-676 berg, Roland Hafner, Francesco Romano, Jonas Buchli, Nicolas Heess, and Martin Ried-677 miller. Continuous-Discrete Reinforcement Learning for Hybrid Control in Robotics. In 678 Proceedings of the Conference on Robot Learning, pages 735–751. PMLR, May 2020. URL 679 https://proceedings.mlr.press/v100/neunert20a.html. ISSN: 2640-3498. 680 Esa Nummelin. General Irreducible Markov Chains and Non-Negative Operators. Cambridge 681 Tracts in Mathematics. Cambridge University Press, Cambridge, 1984. ISBN 978-0-521-60494-9. 682 doi: 10.1017/CBO9780511526237. URL https://www.cambridge.org/core/books/ 683 general-irreducible-markov-chains-and-nonnegative-operators/ 684 0557D49C011AA90B761FC854D5C14983. 685 686 Ling Pan, Dinghuai Zhang, Aaron Courville, Longbo Huang, and Yoshua Bengio. Generative Augmented Flow Networks, October 2022. URL http://arxiv.org/abs/2210.03308. 687 arXiv:2210.03308 [cs]. 688 689 Ling Pan, Nikolay Malkin, Dinghuai Zhang, and Yoshua Bengio. Better Training of GFlowNets with 690 Local Credit and Incomplete Trajectories. In Proceedings of the 40th International Conference 691 on Machine Learning, pages 26878–26890. PMLR, July 2023. URL https://proceedings. 692 mlr.press/v202/pan23c.html. ISSN: 2640-3498. 693 Grigorios A. Pavliotis. Stochastic Processes and Applications: Diffusion Processes, the Fokker-694 Planck and Langevin Equations, volume 60 of Texts in Applied Mathematics. Springer, New York, NY, 2014. ISBN 978-1-4939-1322-0. doi: 10.1007/978-1-4939-1323-7. URL https: 696 //link.springer.com/10.1007/978-1-4939-1323-7. 697 Venkata K. Ramaswamy, Samuel C. Musson, Chris G. Willcocks, and Matteo T. Degiacomi. Deep Learning Protein Conformational Space with Convolutions and Latent Interpolations. Physical 699 *Review X*, 11(1):011052, March 2021. doi: 10.1103/PhysRevX.11.011052. URL https:// 700 link.aps.org/doi/10.1103/PhysRevX.11.011052. Publisher: American Physical Society.

702 703 704 705	Jarrid Rector-Brooks, Kanika Madan, Moksh Jain, Maksym Korablyov, Cheng-Hao Liu, Sarath Chandar, Nikolay Malkin, and Yoshua Bengio. Thompson sampling for improved exploration in GFlowNets, June 2023. URL http://arxiv.org/abs/2306.17693. arXiv:2306.17693 [cs].
707 708 709	Desik Rengarajan, Gargi Vaidya, Akshay Sarvesh, Dileep Kalathil, and Srinivas Shakkottai. Reinforcement Learning with Sparse Rewards using Guidance from Offline Demonstration, February 2022. URL http://arxiv.org/abs/2202.04628. arXiv:2202.04628.
710 711 712 713	Romelia Salomon-Ferrer, David A. Case, and Ross C. Walker. An overview of the Amber biomolecular simulation package. <i>WIREs Computational Molecular Science</i> , 3(2):198–210, 2013. ISSN 1759-0884. doi: 10.1002/wcms.1121. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/wcms.1121eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/wcms.1121.
714 715 716 717 718	Marcin Sendera, Minsu Kim, Sarthak Mittal, Pablo Lemos, Luca Scimeca, Jarrid Rector-Brooks, Alexandre Adam, Yoshua Bengio, and Nikolay Malkin. Improved off-policy training of diffusion samplers, May 2024. URL http://arxiv.org/abs/2402.05098. arXiv:2402.05098 [cs, stat].
719 720 721 722	Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando De Freitas. Taking the Human Out of the Loop: A Review of Bayesian Optimization. <i>Proceedings of the IEEE</i> , 104(1): 148–175, January 2016. ISSN 0018-9219, 1558-2256. doi: 10.1109/JPROC.2015.2494218. URL https://ieeexplore.ieee.org/document/7352306/.
723 724 725 726 727	Max W. Shen, Emmanuel Bengio, Ehsan Hajiramezanali, Andreas Loukas, Kyunghyun Cho, and Tommaso Biancalani. Towards Understanding and Improving GFlowNet Training. <i>Proceedings of</i> <i>the 40th International Conference on Machine Learning</i> , 2023a. doi: 10.48550/ARXIV.2305.07170. URL https://arxiv.org/abs/2305.07170. Publisher: arXiv Version Number: 1.
728 729 730	Tony Shen, Mohit Pandey, Jason Smith, Artem Cherkasov, and Martin Ester. TacoGFN: Target Conditioned GFlowNet for Structure-Based Drug Design, December 2023b. URL http://arxiv.org/abs/2310.03223. arXiv:2310.03223 [cs].
731 732 733 734 735	Hythem Sidky, Wei Chen, and Andrew L. Ferguson. Machine learning for collective variable discovery and enhanced sampling in biomolecular simulation. <i>Molecular Physics</i> , 118(5), March 2020. ISSN 0026-8976. doi: 10.1080/00268976.2020.1737742. URL https://www.osti.gov/pages/biblio/1651117. Institution: Argonne National Laboratory (ANL), Argonne, IL (United States) Publisher: Taylor & Francis.
736 737 738 739 740	Laura Painton Swiler, Patricia Diane Hough, Peter Qian, Xu Xu, Curtis B. Storlie, and Herbert K. H. Lee. Surrogate models for mixed discrete-continuous variables. Technical Report SAND2012-0491, Sandia National Laboratories (SNL), Albuquerque, NM, and Livermore, CA (United States), August 2012. URL https://www.osti.gov/biblio/1055621.
741 742 743 744	Alexandra Volokhova, Michał Koziarski, Alex Hernández-García, Cheng-Hao Liu, Santiago Miret, Pablo Lemos, Luca Thiede, Zichao Yan, Alán Aspuru-Guzik, and Yoshua Bengio. Towards equilibrium molecular conformation generation with GFlowNets, October 2023. URL http: //arxiv.org/abs/2310.14782. arXiv:2310.14782 [cs].
745 746 747 748	Yi Isaac Yang, Qiang Shao, Jun Zhang, Lijiang Yang, and Yi Qin Gao. Enhanced sampling in molecular dynamics. <i>The Journal of Chemical Physics</i> , 151(7):070902, August 2019. ISSN 0021-9606. doi: 10.1063/1.5109531. URL https://doi.org/10.1063/1.5109531.
749 750 751	David W. Zhang, Corrado Rainone, Markus Peschl, and Roberto Bondesan. Robust Schedul- ing with GFlowNets, February 2023. URL http://arxiv.org/abs/2302.05446. arXiv:2302.05446 [cs].
752 753 754 755	Dinghuai Zhang, Nikolay Malkin, Zhen Liu, Alexandra Volokhova, Aaron Courville, and Yoshua Bengio. Generative Flow Networks for Discrete Probabilistic Modeling. In <i>Proceedings of the 39th International Conference on Machine Learning</i> , pages 26412–26428. PMLR, June 2022. URL https://proceedings.mlr.press/v162/zhang22v.html. ISSN: 2640-3498.

## A LOSS FUNCTIONS

For a complete trajectory  $\tau$ , the *detailed balanced loss* (DB) is

759 760 761

762

763 764

765 766

756

757 758

$$L_{DB}(\tau) = \sum_{t=0}^{n-1} \left( \log \frac{\hat{f}(s_t; \theta) \hat{p}_F(s_t, s_{t+1}; \theta)}{\hat{f}(s_{t+1}; \theta) \hat{p}_B(s_{t+1}, s_t; \theta)} \right)^2$$

where  $\hat{f}(s_{t+1}; \theta)$  is replaced with  $r(s_n)$  if  $s_n$  is terminal.

The subtrajectory balance loss (STB) is

$$L_{STB}(\tau) = \frac{\sum_{0 \le i < j \le n} \lambda^{j-1} \mathcal{L}_{TB}(\tau_{i:j})}{\sum_{0 \le i < j \le n} \lambda^{j-i}},$$

767 768 769

770

775

776 777

778

779 780

$$\mathcal{L}_{STB}(\tau_{i:j}) \coloneqq \left( \log \frac{\hat{f}(s_i;\theta) \prod_{t=i}^{j-1} \hat{p}_F(s_{t+1}|s_t;\theta)}{\hat{f}(s_j;\theta) \prod_{t=i+1}^{j} \hat{p}_B(s_{t-1}|s_t;\theta)} \right)^2$$

where  $\hat{f}(s_j; \theta)$  is replaced with  $r(s_j)$  if  $s_j$  is terminal. In the above,  $\lambda < 0$  is a hyperparameter. The limit  $\lambda \to 0^+$  leads to average detailed balance. The  $\lambda \to \infty$  limit gives the trajectory balance objective. We use  $\lambda = 0.9$  in our experiments.

## **B** LANGEVIN DYNAMICS

Langevin dynamics (LD), is defined through the Stochastic Differential Equation (SDE):

$$\mathrm{d}x = M^{-1}p\mathrm{d}t \tag{3}$$

$$dp = F(x)dt - \gamma pdt + \sqrt{2\gamma\beta^{-1}}M^{1/2}dW.$$
(4)

<sup>781</sup> In the above,  $x, p \in \mathbb{R}^D$  are vectors of instantaneous position and momenta respectively,  $F : \mathbb{R}^D \to \mathbb{R}^D$  is a force function, W(t) is a vector of D independent Wiener processes,  $M \in \mathbb{R}^D \times \mathbb{R}^D$  is a constant diagonal mass matrix, and  $\gamma, \beta > 0$  are constant scalars which can be interpreted as a friction coefficient and inverse temperature respectively. In conventional Langevin dynamics, the force function is given by the gradient of the potential energy function,  $F(x) = -\nabla V(x)$ , where  $V : \mathbb{R}^D \to \mathbb{R}$  and the dynamics are ergodic with respect to the Gibbs-Boltzmann density  $\rho_\beta(x, p) \propto e^{-\beta H(x, p)}$ ,

where  $H(x, p) = p^T M^{-1} p/2 + V(x)$  is the Hamiltonian. Since the Hamiltonian is separable in position and momenta terms, the marginal Gibbs-Boltzmann density is position space is simply  $\rho_{\beta}(x) \propto e^{-\beta V(x)}$ , often referred to as the *Gibbs density*.

We present the Euler-Maruyama numerican scheme for integrating equations 3 below. This is called
in line 14 of Adapted Metadynamics (Algorithm 1).

795	Algorithm 3: Euler-Maruyama Langevin Dynamics Step		
796	Input	:Current state $(x_t, p_t)$ . Force F.	
797	Paramet	ters: Friction coefficient $\gamma$ . Thermodynamic beta $\beta$ . Timestep $\Delta t$ .	
798	Output	:State $(x_{t+\Delta t}, p_{t+\Delta t})$ at the next timestep.	
799	1 Sample a	a random vector R, with the same dimension as $x_t$ , where each element is an	
800	indepen	ident sample from a standard normal.	
801	2 $x_{t+\Delta t} =$	$x_t + p_t \Delta t$	
802	$p_{t+\Delta t} =$	$p_t + F\Delta t - \gamma p_t \Delta t + \sqrt{2\gamma\Delta t/\beta} \cdot R$	
803	4 return $(x_{t+\Delta t}, p_{t+\Delta t})$		
804			
805		0.0.5%	
806	C PR	UUFS	
807	Ŧ		
	Lemma	<b>U.I.</b> Let $(f_n(x))$ and $(q_n(x))$ be sequences of real functions where $\lim_{n\to\infty} f_n(x) = \infty$	

Lemma C.1. Let  $(f_n(x))$  and  $(g_n(x))$  be sequences of real functions where  $\lim_{n\to\infty} f_n(x) = \infty$ ,  $\lim_{n\to\infty} g_n(x) = \infty$  and  $\lim_{n\to\infty} \frac{f_n(x)}{g_n(x)} = h(x)$ . Then, for all  $\epsilon > 0$ , we have  $\lim_{n\to\infty} \frac{f_n(x)}{g_n(x)+\epsilon} = h(x)$ . Proof. 

$$\lim_{n \to \infty} \frac{f_n(x)}{g_n(x) + \epsilon} = \lim_{n \to \infty} \frac{f_n(x)}{g_n(x)} \frac{1}{1 + \epsilon/g_n(x)}$$

and the right-hand side is the product of two functions whose limit exists so, by the product rule of limits

$$\lim_{n \to \infty} \frac{f_n(x)}{g_n(x)} \lim_{n \to \infty} \frac{1}{1 + \epsilon/g_n(x)} = h(x) \cdot 1 = h(x),$$

so done.

**Lemma C.2.** Let  $(X_i)$  be a sequence of continuous random variables that take values on a bounded domain  $D \subset \mathbb{R}^d$  that asymptotically approaches the uniform random variable U on D, i.e.  $X_i \to U$ uniformly. Further, suppose

$$h(x) := \frac{\sum_{i=1}^{\infty} f(x_i)g(x, x_i)}{\sum_{i=1}^{\infty} g(x, x_i)}$$

exists, where  $x_i \in D$  is a sample from  $X_i$  and f(x) and g(x, x') are analytic functions on D and  $D \times D$  respectively. Then,

$$h(x) = \frac{\sum_{i=1}^{\infty} f(u_i)g(x, u_i)}{\sum_{i=1}^{\infty} g(x, u_i)},$$

where the  $u_i$  are samples from U. We make no assumption of independence of samples.

Ì

*Proof.* Fix a probability space  $(\Omega, \mathcal{F}, P)$  on which  $(X_i)$  and U are defined. Recall that a continuous random variable X that takes values on  $D \subset \mathbb{R}^d$  is a measurable function  $X : \Omega \to D$  where  $(D, \mathcal{B})$ is a measure space and  $\mathcal{B}$  is the Borel  $\sigma$ -algebra on D. Let  $\omega \subset \Omega$  denote an arbitrary element of the sample space. The requirement that  $X_i \rightarrow U$  uniformly can be written formally as: 

$$\forall \epsilon > 0, \exists N(\epsilon) \text{ s.t. } \forall i > N(\epsilon), \forall \omega \subset \Omega, |X_i(\omega) - U(\omega)| < \epsilon.$$

We prove the Lemma by showing that equality holds for all possible sequences of outcomes  $\omega_1, \omega_2, \ldots$  That is, we prove:

$$\frac{\sum_{i=1}^{\infty} f(X_i(\omega_i))g(x, X_i(\omega_i))}{\sum_{i=1}^{\infty} g(x, X_i(\omega_i))} = \frac{\sum_{i=1}^{\infty} f(U(\omega_i))g(x, U(\omega_i)))}{\sum_{i=1}^{\infty} g(x, U(\omega_i)))}.$$
(5)

Since these are ratios of infinite series, to prove their equality it is sufficient to show that the numerator of the LHS is asymptotically equivalent to the numerator of the RHS and that the denominator of the LHS is asymptotically equivalent to the denominator of the RHS. Recall that two sequences of real functions  $(a_n)$  and  $(b_n)$  are asymptotically equivalent if  $\lim_{n\to\infty} \frac{a_n(x)}{b_n(x)} = c$  where c is a constant. First, we prove that this holds with 

$$a_n := \sum_{i=1}^n f(X_i(\omega_i))g(x, X_i(\omega_i))$$
(6)

$$b_n := \sum_{i=1}^n f(U(\omega_i))g(x, U(\omega_i))).$$
(7)

We write

and

$$\frac{a_n}{b_n} = \frac{\sum_{i=1}^{N(\epsilon)} f(X_i(\omega_i))g(x, X_i(\omega_i)) + \sum_{i=N(\epsilon)+1}^n f(X_i(\omega_i))g(x, X_i(\omega_i))}{\sum_{i=1}^{N(\epsilon)} f(U(\omega_i))g(x, U(\omega_i))) + \sum_{i=N(\epsilon)+1}^n f(U(\omega_i))g(x, U(\omega_i)))}$$

Dividing by  $\sum_{i=N(\epsilon)+1}^{n} f(U(\omega_i))g(x, U(\omega_i)))$  and taking the limit  $n \to \infty$  we have

$$\lim_{n \to \infty} \frac{a_n}{b_n} = \lim_{n \to \infty} \frac{\sum_{i=N(\epsilon)+1}^n f(X_i(\omega_i))g(x, X_i(\omega_i))}{\sum_{i=N(\epsilon)+1}^n f(U(\omega_i))g(x, U(\omega_i)))}$$

Since f and g are analytic and  $i > N(\epsilon)$  for all terms in the sums we have, by Taylor expansion,  $f(X_i(\omega_i)) = f(U(\omega_i)) + O(\epsilon)$  and  $g(x, X_i(\omega_i)) = g(x, U(\omega_i)) + O(\epsilon)$ , hence

$$\lim_{n \to \infty} \frac{a_n}{b_n} = \lim_{n \to \infty} \left( 1 + \frac{nO(\epsilon)}{\sum_{i=N(\epsilon)+1}^n f(U(\omega_i))g(x, U(\omega_i)))} \right) = 1 + \lim_{n \to \infty} \frac{nO(\epsilon)}{O(n)} = 1 + O(\epsilon).$$

Finally, since  $\epsilon$  can be made arbitrarily small by partitioning the sum at a  $N(\epsilon)$  that is sufficiently large, we conclude that  $\lim_{n\to\infty} \frac{a_n}{b_n} = 1$ , hence  $(a_n)$  and  $(b_n)$  as defined in equation 6 and equation 7 are asymptotically equivalent. By a similar argument, it can be shown that

$$c_n := \sum_{i=1}^n g(x, X_i(\omega_i))$$

871 and

$$d_n := \sum_{i=1}^n g(x, U(\omega_i)))$$

are also asymptotically equivalent. This proves equation 5.

Below, we present the proof of Theorem 3.1 that appears in the main text.

*Proof.* For concreteness, throughout this proof we assume that the kernel function is a Gaussian. We explain at the appropriate stage in the proof, indicated by (\*), how this assumption can be relaxed.

First we take the  $n \to \infty$  limit. Recall the notation from Section 3, i.e. assume collective variables  $z(x) = (z_1(x), \ldots, z_d(x))$  where  $z : \mathcal{X} \to \mathcal{Z}$ , and  $\mathcal{Z}$  is *d*-dimensional. Since the log function is continuous, the limit and log can be interchanged and we have

$$\lim_{n \to \infty} \hat{V}(z, t_n) = -\frac{1}{\beta'} \log \left( \lim_{n \to \infty} \left( \frac{\hat{R}(z, t_n)}{\hat{N}(z, t_n) + \epsilon} \right) + \epsilon \right)$$

where, from Algorithm 1:

$$\hat{R}(z,t_n) = \sum_{i=1}^n r(x_{t_i}) \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(x_{t_i}))^2}{2\sigma_j^2}\right),\tag{8}$$

 $\hat{N}(z,t_n) = \sum_{i=1}^n \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(x_{t_i}))^2}{2\sigma_j^2}\right).$ (9)

Since the domain is bounded, we know that for fixed z, both equation 8 and equation 9 have limit at infinity, i.e.  $\lim_{n\to\infty} \hat{R}(z,t_n) = \infty$  and  $\lim_{n\to\infty} \hat{N}(z,t_n) = \infty$ . Hence, by Lemma C.1, we have

$$\lim_{n \to \infty} \frac{\hat{R}(z, t_n)}{\hat{N}(z, t_n) + \epsilon} = \lim_{n \to \infty} \frac{\hat{R}(z, t_n)}{\hat{N}(z, t_n)},$$

901 provided the limit on the RHS exists. Next, we show that this limit exists by computing it explicitly. 902 The limit can be written 

$$\lim_{n \to \infty} \frac{\hat{R}(z, t_n)}{\hat{N}(z, t_n)} = \lim_{n \to \infty} \frac{\sum_{i=1}^n r(x_{t_i}) \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(x_{t_i}))^2}{2\sigma_j^2}\right)}{\sum_{i=1}^n \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(x_{t_i}))^2}{2\sigma_j^2}\right)}$$

907 Recall that metadynamics eventually leads to uniform sampling over the domain, independent of the 908 potential. Hence, since R and z(x) are analytic, by Lemma C.2 we may replace the metadynamics 909 samples  $x_{t_i}$  with samples  $u_i$  from a uniform distribution over  $\mathcal{X}$ :

$$\lim_{n \to \infty} \frac{\hat{R}(z, t_n)}{\hat{N}(z, t_n)} = \lim_{n \to \infty} \frac{\sum_{i=1}^n r(u_i) \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(u_i))^2}{2\sigma_j^2}\right)}{\sum_{i=1}^n \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(u_i))^2}{2\sigma_j^2}\right)}$$

In the limit, the ratio of sums with uniform sampling becomes a ratio of integrals:

915  
916  
917 
$$\lim_{n \to \infty} \frac{\hat{R}(z, t_n)}{\hat{N}(z, t_n)} = \frac{\int_{\mathcal{X}} r(x') \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(x'))^2}{2\sigma_j^2}\right) \mathrm{d}x'}{\int_{\mathcal{X}} \exp\left(-\sum_{j=1}^d \frac{(z_j - z_j(x'))^2}{2\sigma_j^2}\right) \mathrm{d}x'}$$

The limit is therefore a (scaled) convolution of the reward function with a Gaussian in the collective variable space with width vector  $\sigma$ . Taking the limit  $\sigma_j \to 0$  for all  $j \in \{1, 2, ..., d\}$ , the Gaussian convergences to a delta distribution in the collective variable space and we have

921 922 923

924

925

926

927

$$\lim_{\sigma \to 0} \lim_{n \to \infty} \frac{\dot{R}(z, t_n)}{\dot{N}(z, t_n)} = \int_{\mathcal{X}} r(x')\delta(z - z(x')) \mathrm{d}x'.$$

(\*) This step also holds for any kernel that becomes distributionally equivalent to a Dirac delta function in the limit that its variance parameter goes to zero. In particular, it also holds for the von Mises distribution that we use in our alanine dipeptide experiment in  $\mathbb{T}^2$ .

Finally, we take the limit  $\epsilon \to 0$  to obtain

$$\lim_{\epsilon \to 0} \lim_{\sigma \to 0} \lim_{n \to \infty} \hat{V}(z, t_n) = -\frac{1}{\beta'} \lim_{\epsilon \to 0} \log\left(\int_{\mathcal{X}} r(x')\delta(z - z(x'))dx' + \epsilon\right)$$
(10)

$$= -\frac{1}{\beta'} \int_{\mathcal{X}} \log\left(r(x')\right) \delta(z - z(x')) \mathrm{d}x' \tag{11}$$

$$= \int_{\mathcal{X}} V(x')\delta(z - z(x'))\mathrm{d}x' := V(z), \tag{12}$$

where we have used the definition  $V(x') = -\frac{1}{\beta'} \log(r(x'))$  and in the last step we used the definition of the marginal potential energy in the collective variable space. If z(x) is invertible, then the delta function simplifies to a delta function in the original space and we obtain the original potential instead of the marginal potential.

=

## 940 941 942

943

944

964

965

966

935 936

937

938

939

## D EXPERIMENT DETAILS

### D.1 COMMON EXPERIMENTAL PARAMETERS

**Training parameters**: For all environments we train for  $B = 10^5$  batches use a learning rate with a linear schedule, starting at  $10^{-3}$  and finishing at 0. For the TB loss, we train the  $\log Z_{\theta}$  with a higher initial learning rate of  $10^{-1}$  (also linearly scheduled). For the STB loss, we use  $\lambda = 0.9$ , a value that has worked well in the discrete setting (Madan et al., 2022). We use the Adam optimiser with gradient clipping. For all loss functions, we clip the minimum log-reward signal at -10. This enables the model to learn despite regions of near-zero reward between the modes of r(x).

**Replay buffer**: We use a replay buffer with capacity for  $10^4$  trajectories. Trajectories are stored in the replay buffer only if the terminal state's reward exceeds  $10^{-3}$ . When drawing a replay buffer batch, trajectories are bias-sampled: 50% randomly drawn from the upper 30% of trajectories with the highest rewards, and the remaining 50% randomly drawn from the lower 70%.

**Noisy exploration**: An additional constant,  $\bar{\sigma}$ , is added to 956 the standard deviations of the Gaussian distributions of the 957 forward and backward policies. Specifically, the forward 958 policy becomes  $\hat{p}_F(s_t, s_{t-1}; \theta) = \sum_{i=1}^k w_i \mathcal{N}(\mu_i, (\sigma_i + \omega_i)^2)$ 959  $(\bar{\sigma})^2$ ), and similarly for the backward policy, where k is 960 the number of Gaussians in the mixture. We schedule the 961 value of  $\bar{\sigma}$  so that it decreases during training according to 962 an exponential-flat schedule: 963

$$\bar{\sigma} = \begin{cases} \bar{\sigma}_0 \left( e^{-2je/(B/2)} - e^{-2e} \right) & j < B/2\\ 0 & j \ge B/2, \end{cases}$$
(13)



where  $j \in (1, ..., B)$  is the batch number and  $\bar{\sigma}_0 = 2$ is the initial noise, plotted in Figure 7. Note that, for the figure 7: Exponential noise schedule. alanine dipeptide environment, the policy is a mixture of bivariate von Mises distributions and the noise  $\bar{\sigma}$  is added to the concentration parameter  $\kappa$ , where concentration is related to standard deviation through  $\sigma = \frac{1}{\kappa^2}$ .

**Thompson sampling**: We use 10 heads with the bootstrapping probability parameter set to p = 0.3.

**Local Search GFN**: Local search samples trajectories using a forward-backward reconstruction method. First, forward policy samples actions to generate batch size *b* on-policy trajectories. Then, the backward policy is used to re-sample the last *K* steps of each trajectory. Finally, the forward policy is used to reconstruct the last *K* steps, and trajectories in the batch with higher rewards from reconstruction replace the originals. We use K = 1 for the line environment, alanine dipeptide environment and two-dimensional grid. We use K = 2 and K = 3 for the three and four dimensional (hyper)grids, respectively.

979 Nested sampling: We use the Gradient Based Nested Sampling (GBNS) (Lemos et al., 2023) with a
 980 publicly available implementation: https://github.com/Pablo-Lemos/GGNS. For each
 981 environment, GBNS generates samples from the reward distribution. During training, we alternate
 982 between two types of backpropagation: one fully on-policy, and the other using backward-sampled
 983 trajectories from the generated samples, guided by the current backward policy.

984
 985
 Compute resources: Experiments are performed in PyTorch on a desktop with 32Gb of RAM and a 12-core 2.60GHz i7-10750H CPU.

986

987 D.2 LINE ENVIRONMENT DETAILS

989 **Parameterisation**: The forward and backward policies are a mixture of three Gaussians. We 990 parameterise  $\hat{p}_F$ ,  $\hat{p}_B$  and the flow  $\hat{f}$  through an MLP with 3 hidden layers, 256 hidden units per 991 layer. We use the GELU activation function and dropout probability 0.2 after each layer. This 992 defines the torso of the MLP. Connecting from this common torso, the MLP has three single-layer, fully-connected heads. The first two heads have output dimension 9 and parameterise the 3 means  $(\mu)$ , 993 standard deviations ( $\sigma$ ) and weights (w) of the mixture of Gaussians for the forward and backward 994 995 policies respectively. The third head has output dimension 1 and parameterises the flow function f. The mean and standard deviation outputs are passed through a sigmoid function and transformed so 996 that they map to the ranges  $\mu \in (-14, 14)$  and  $\sigma \in (0.1, 1)$ . The mixture weights are normalised with 997 the softmax function. The exception to this parameterisation is the backward transition to the source 998 state which is fixed to be the Dirac delta distribution centred on the source, i.e.  $\hat{p}_B(s_0|s_1;\theta) = \delta_{s_0}$ . 999 For the TB loss, we treat  $\log Z_{\theta}$  as a separate learnable parameter. We use batch size b = 64. It takes 1000 approximately 1 hour to train a continuous GFlowNet in this environment with  $B = 10^5$  batches. 1001

**MetaGFN**: We use  $\Delta t = 0.05$ , n = 2,  $\beta = 1$ ,  $\gamma = 2$ , w = 0.15,  $\sigma = 0.1$ ,  $\epsilon = 10^{-3}$ . The domain of Adapted Metadynamics is restricted to [-5, 23] and reflection conditions are imposed at the boundary. The bias and KDE potentials are stored on a uniform grid with grid spacing 0.01. Initial metadynamics samples are drawn from a Gaussian distribution, mean 0 and variance 1, and initial momenta from a Gaussian distribution, mean 0 and variance 0.5.

**Evaluation:** The L1 error between the known reward distribution,  $\rho(x) = r(x)/Z$ , and the empirical on-policy distribution, denoted  $\hat{\rho}(x)$ , estimated by sampling 10<sup>4</sup> on-policy trajectories and computing the empirical distribution over terminal states. Specifically, we compute error  $= \frac{1}{2} \int_{-5}^{23} \left| \hat{\rho}(x) - \frac{r(x)}{Z} \right| dx$ , where the integral is estimated by a discrete sum with grid spacing 0.01. Note that this error is normalised such that for all valid probability distributions  $\hat{\rho}(x)$ , we have  $0 \leq \operatorname{error} \leq 1$ .

**On-policy distributions**: Figure 9 shows the forward policy and replay buffer distributions (with bias sampling) after training for  $10^5$  iterations with TB loss. MetaGFN is the only method that can uniformly populate the replay buffer and consistently learn all three peaks.

1017Adapted Metadynamics: Figure 10a shows the L1 error between the density implied by the KDE1018potential and the true reward distribution during a typical training run. Figure 10b shows the resulting1019 $\hat{V}$  and  $V_{\text{bias}}$  at the end of the training. By (1), Adapted Metadynamics has fully explored the central1020peaks. At (2), the third peak is discovered, prompting rapid adjustment of the KDE potential. By1021 $2.5 \times 10^4$  iterations, a steady state is reached and the algorithm is sampling the domain uniformly.

**Sensitivity analysis of hyperparameters:** We compare MetaGFN performance for various freqRB and freqMD values, showing the mean and standard deviation of the final L1 error after  $10^5$ batches (Figure 8). Intuitively, we expect that if freqRB is set too large, then MetaGFN focuses predominantly on the on-policy dynamics and may fail to learn to sample distant modes. By contrast, setting freqRB = 1 means exclusively training on the replay buffer, without any on-policy training

1050

1051

1052

1053

1054

1056

1057

1058

1061

1062 1063 1064



Figure 8: Mean and standard deviation L1 loss after  $10^5$  training batches over 10 repeats in the line environment. Blue line: changing freqRB while freqMD = 10. Red line: changing freqMD while freqRB = 2.



Figure 9: Forward policy and replay buffer distributions after training for  $10^5$  iterations with TB loss. MetaGFN is the only method that can consistently learn all three peaks.

trajectories, which we would expect to lead to slow convergence. Therefore, an intermediate value of freqRB should be optimal. We found that setting freqRB = 2 (giving an equal weighting to both on-policy and replay buffer trajectories) showed the lowest loss. By contrast, performance is not too sensitive to freqMD. freqMD should be set so that Adapted Metadynamics explores the reward landscape in a comparable time it takes to train the GFN to convergence, hence only the order of magnitude of freqMD is important. For concreteness, we choose freqMD = 10 but the model performance is not too sensitive to this parameter.

Comparing MetaGFN variants: We consider three MetaGFN variants. The first variant, *always* backwards sample, regenerates the entire trajectory using the current backward policy when pulling from the replay buffer. The second variant, *reuse initial backwards sample*, generates the trajectory when first added to the replay buffer and reuses the entire trajectory if subsequently sampled. The third variant, *with noise*, is to always backwards sample with noisy exploration as per equation equation 13. We plot the L1 policy errors in Figure 11. We observe that *always backward sample* is better than *reuse initial backwards sample* for all loss functions. For DB and TB losses, there is no evidence for any benefit of adding noise, whereas noise improves training for STB loss, performing very similarly to TB loss without noise.



# 1134 D.3 Alanine dipeptide environment details

1136 **Computing the free energy surface**: To obtain a ground-truth 1137 free energy surface (FES) in  $\phi$ - $\psi$  space, we ran a 250ns NPT well-1138 tempered metadynamics MD simulation of alanine dipeptide at tem-1139 perature 300K ( $\beta = 0.4009$ ), pressure 1bar with the TIP3P explicit 1140 water model (Jorgensen et al., 1983). We used the PLUMMED 1141 plugin (Bonomi et al., 2019) for OpenMM (Eastman et al., 2017) 1142 with the AMBER14 force field (Salomon-Ferrer et al., 2013).

Parameterisation: The forward and backward policies are a mixture 1143 of three bivariate von Mises distributions. We parameterise  $\hat{p}_F$ ,  $\hat{p}_B$ 1144 and  $\hat{f}$  through three heads of an MLP with 3 hidden layers with 1145 512 hidden units per layer, with GeLU activations and dropout 1146 probability 0.2. The first two heads have output dimensions of 1147 15, parameterising the 6 means, 6 concentrations, and 3 weights 1148 of the mixture of von Mises policy. The third head has output 1149 dimension 1 and parameterises the flow function f. The means are 1150 mapped to the range  $(-\pi, \pi)$  through  $2 \arctan(\cdot)$ . Concentrations 1151



Figure 12: The potential KDE  $\hat{V}$  learnt using MetaGFN on the alanine dipeptide environment, concentration  $\kappa = 10$ .

are parameterised in log space and passed through a sigmoid to map to the range  $\ln(\kappa) \in (0, 5)$ . Mixture weights are normalised with the softmax function. We use batch size b = 64. It takes approximately 10 hours to train a continuous GFlowNet in this environment with  $B = 10^5$  batches.

**MetaGFN:** We use freqRB = 2, freqMD = 10,  $\Delta t = 0.01$ , n = 2,  $\beta = 0.4009$ ,  $\gamma = 0.1$ ,  $w = 10^{-5}$ ,  $\kappa = 10$ ,  $\epsilon = 10^{-6}$ . The bias and KDE potentials are stored on a uniform grid with grid spacing 0.1. Initial samples are drawn from a Gaussian distribution, mean centred  $P_{||}$ , variance  $\sigma^2 = (0.1, 0.1)$ , and initial momenta from a Gaussian, mean  $\mu = (0, 0)$ , variance  $\sigma^2 = (0.05, 0.05)$ . In Figure 12 we show the resulting learnt KDE potential with these parameters.

**Evaluation**: The L1 error of a histogram of on-policy samples,  $\hat{\rho}(\phi, \psi)$ , is computed via a twodimensional generalisation of equation D.2; error  $= \frac{1}{2} \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} \left| \hat{\rho}(\phi, \psi) - \frac{r(\phi, \psi)}{Z} \right| d\phi d\psi$ , estimated by a discrete sum with grid spacing 0.1.

- 1164
- D.4 GRID ENVIRONMENTS DETAILS

**Parameterisation**: For the two-dimensional grid, we use a mixture of 4 Gaussians. For the three and four dimensional (hyper)grids, we use a mixture of 2 Gaussians (but use longer trajectories, see Section 4). We parameterise  $\hat{p}_F$ ,  $\hat{p}_B$  and the flow  $\hat{f}$  through an MLP with 3 hidden layers, 512 hidden units per layer, GELU activation function and dropout probability 0.2. The means and standard deviations are mapped to the range  $\mu \in (-15, 15)$  and  $\sigma \in (0.1, 7)$  through sigmoid functions. Mixture weights are normalised with the softmax function. We use batch size b = 128. It takes approximately 10 hours to train a continuous GFlowNet in this environment with  $B = 10^5$  batches.

**MetaGFN:** We use freqRB = 2, freqMD = 10,  $\Delta t = 0.35$ , n = 3,  $\beta = 1$ ,  $\gamma = 2$ , w = 0.10,  $\sigma = 2$ . The bias and KDE potentials are stored on a uniform grid with grid spacing 0.075/0.3/0.6(2, 3, and 4 dimensions respectively). Initial samples are drawn from a Gaussian distribution, mean centred at the origin, isotropic variance  $\sigma^2 = 1$ . Momenta are initialised to zero.

**Evaluation**: The L1 error between the histogram of the terminal states of  $10^4$  on-policy samples and the exact distribution is computed, e.g. for the 2D grid, error  $= \frac{1}{2} \int_{-15}^{15} \int_{-15}^{15} \left| \hat{\rho}(x,y) - \frac{r(x,y)}{Z} \right| dxdy$ , estimated by a discrete sum with grid spacing 0.075.

- 1181
- 1182 1183
- 1184
- 1185

1186