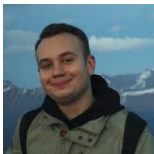


Generative Flow Networks as Entropy-Regularized RL

Daniil Tiapkin^{1,3}, Nikita Morozov¹, Alexey Naumov¹, Dmitry Vertov¹

¹HSE University ²AIRI ³École Polytechnique

Based on <https://arxiv.org/abs/2310.12934>



Problem

Suppose we have a finite discrete space of objects \mathcal{X} with a black-box non-negative reward function $\mathcal{R}(x)$ given for every $x \in \mathcal{X}$.

Goal: Construct and train a model that will sample objects from \mathcal{X} from the probability distribution $\mathcal{R}(x)/Z$ proportional to the reward, where Z is an unknown normalizing constant.

Standard approach: MCMC methods ...

- What can be done if drawing i.i.d. samples from $\mathcal{R}(\cdot)/Z$ is not an option?
- If we run the (ergodic) Markov chain $(Z_k)_{k \geq 0}$ for a long time (started from anywhere), then for large N the distribution of Z_N will be approximately invariant: $\text{Law}(Z_N) \approx \mathcal{R}(\cdot)/Z$. We can then set $X_1 = Z_N$, and then restart and rerun the Markov chain to obtain X_2, X_3, \dots ,

Important question

How to construct $P(x, A)$ such that the distribution of X_n converges to invariant distribution $\mathcal{R}(\cdot)/Z$ as quickly as possible for arbitrary initial distribution ξ ?

Example: Metropolis-Hastings algorithm

Let $Q(x, A) = \int_A q(x, y) dy$ be some MK (e.g. Gaussian)

1. Choose X_0 .
2. Given X_k , a candidate move Y_{k+1} is sampled from $Q(X_k, \cdot)$
3. $X_{k+1} = Y_{k+1}$ with probability $\alpha(X_k, Y_{k+1})$, otherwise $X_{k+1} = X_k$, where acceptance ratio

$$\alpha(x, y) = \min \left\{ 1, \frac{\mathcal{R}(y)q(y, x)}{\mathcal{R}(x)q(x, y)} \right\}$$

Example: Random walk MH

Take $q(x, y) = \bar{q}(y - x)$, where $\bar{q}(x) = \bar{q}(-x)$. Then

$$Y_{k+1} = X_k + Z_{k+1}, \quad Z_{k+1} \sim \bar{q}$$

In this case

$$\alpha(x, y) = \min \left\{ 1, \frac{\mathcal{R}(y)}{\mathcal{R}(x)} \right\}$$

Example: Langevin Dynamics

Langevin Dynamics Itô SDE:

$$dX_t = -\nabla U(X_t) dt + \sqrt{2}dW_t,$$

Invariant measure: $\mathcal{R}(x) = e^{-U(x)}$ and $\text{Law}(X_t) \rightarrow \mathcal{R}(\cdot)/Z$ as $t \rightarrow \infty$.

1. Take $\mathcal{R}(x)/Z = (2\pi)^{-1/2}e^{-x^2/2}$.
2. SDE: $dX_t = -X_t dt + \sqrt{2}dW_t$, X_0 is independent of W . This is Ornstein–Uhlenbeck process
3. Apply Ito's formula to obtain

$$X_t = X_0 e^{-t} + \sqrt{2} \int_0^t e^{-(t-s)} dW_s$$

4. Since the Itô integral of deterministic integrand is normally distributed, we readily have

$$\text{Law}(X_t) = \mathcal{N}(\theta_0 e^{-t}, 1 - e^{-2t}) \rightarrow \mathcal{N}(0, 1)$$

Example: Langevin Dynamics

Itô SDE:

$$dX_t = -\nabla U(X_t) dt + \sqrt{2} dW_t,$$

Invariant measure: $\mathcal{R}(x) = e^{-U(x)}$

1. First-order discretization (Unadjusted Langevin Algorithm, ULA):

$$Y_{k+1} = Y_k - \gamma \nabla U(Y_k) + \sqrt{2\gamma} Z_{k+1}, \quad i.i.d. \ Z_k \sim \mathcal{N}(0, I_d)$$

Equivalently, $Y_{k+1} \sim \mathcal{N}(Y_k - \gamma \nabla U(Y_k), 2\gamma I_d)$

2. Metropolis-adjusted Langevin Algorithm (MALA):
ULA + Metropolis-Hastings correction;
3. Demo: <https://chi-feng.github.io/mcmc-demo>
4. If we can't calculate ∇U replace it by its estimate over batch (SGLD, SGLD-FP, SAGA etc)
5. Mixture of global (importance sampling) and local steps with trainable proposals (Ex²MCMC), see [Samsonov et al., 2022]

Problem

Suppose we have a finite discrete space of objects \mathcal{X} with a black-box non-negative reward function $\mathcal{R}(x)$ given for every $x \in \mathcal{X}$.

Goal: Construct and train a model that will sample objects from \mathcal{X} from the probability distribution $\mathcal{R}(x)/Z$ proportional to the reward, where Z is an unknown normalizing constant.

Standard approach: MCMC methods usually suffer from the mixing problem.

GFlowNet as an alternative: it learns a stochastic policy that starts with “empty” objects and modifies them until completion. GFlowNet generates each sample independently and, like RL methods, relies on exploration and the generalization power of ML to guess and discover new modes of the reward function (or regions of low energy).

References: [Bengio et al., 2021], [Bengio et al., 2023] etc.

Flows Over Directed Acyclic Graphs

DAG: $\mathcal{G} = (\mathcal{S}, \mathcal{E})$: \mathcal{S} is a state space, $\mathcal{E} \subseteq \mathcal{S} \times \mathcal{S}$ is a set of edges. For $s \in \mathcal{S}$, \mathcal{A}_s is the set of actions, i.e. possible next states: $s' \in \mathcal{A}_s \Leftrightarrow s \rightarrow s' \in \mathcal{E}$; $s_0 \in \mathcal{S}$ is an initial state ("empty" object) and the set of terminal states (states with no outgoing edges) directly corresponds to \mathcal{X}

Denote \mathcal{T} as a set of all complete trajectories in the graph

$\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_{n_\tau})$, where n_τ is the trajectory length,

$s_i \rightarrow s_{i+1} \in \mathcal{E}$, s_0 is always the initial state, and s_{n_τ} is always a terminal state, thus $s_{n_\tau} \in \mathcal{X}$

A **trajectory flow** is any nonnegative function $\mathcal{F}: \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}$.

$$\mathcal{F}(s) \triangleq \sum_{\tau \ni s} \mathcal{F}(\tau),$$

$$\mathcal{F}(s \rightarrow s') \triangleq \sum_{\tau \ni (s \rightarrow s')} \mathcal{F}(\tau).$$

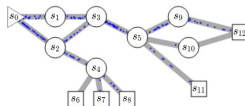


Figure: Directed Acyclic Graphs

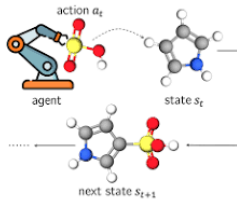


Figure: Molecule synthesis experiment

Reward Matching Condition

If \mathcal{F} is not identically zero, it can be used to introduce a probability distribution over complete trajectories:

$$\mathcal{P}(\tau) \triangleq \frac{1}{Z_{\mathcal{F}}} \mathcal{F}(\tau), \quad Z_{\mathcal{F}} \triangleq \sum_{\tau \in \mathcal{T}} \mathcal{F}(\tau).$$

Then, $\mathcal{F}(s)/Z_{\mathcal{F}}$ corresponds to the probability that a random trajectory contains the state s , and $\mathcal{F}(s \rightarrow s')/Z_{\mathcal{F}}$ corresponds to the probability that a random trajectory contains the edge $s \rightarrow s'$. These definitions also imply $Z_{\mathcal{F}} = \mathcal{F}(s_0) = \sum_{x \in \mathcal{X}} \mathcal{F}(x)$.

Denote $\mathcal{P}(x)$ as the probability that x is the terminal state of a trajectory sampled from $\mathcal{P}(\tau)$. If the *reward matching constraint* is satisfied for all terminal states $x \in \mathcal{X}$:

$$\mathcal{F}(x) = \mathcal{R}(x), \tag{1}$$

$\mathcal{P}(x)$ will be equal to $\mathcal{R}(x)/Z$. **Thus by sampling trajectories we can sample the objects from the distribution of interest**

Markovian flow

The trajectory flow \mathcal{F} is said to be Markovian if $\forall s \in \mathcal{S} \setminus \mathcal{X}$ there exist distributions $\mathcal{P}_F(-|s)$ over the children of s such that \mathcal{P} can be factorized as

$$\mathcal{P}(\tau = (s_0 \rightarrow \dots \rightarrow s_{n_\tau})) = \prod_{t=1}^{n_\tau} \mathcal{P}_F(s_t | s_{t-1}).$$

Equivalently, \mathcal{F} is Markovian if $\forall s \in \mathcal{S} \setminus \{s_0\}$ there exist distributions $\mathcal{P}_B(-|s)$ over the parents of s such that $\forall x \in \mathcal{X}$

$$\mathcal{P}(\tau = (s_0 \rightarrow \dots \rightarrow s_{n_\tau}) | s_{n_\tau} = x) = \prod_{t=1}^{n_\tau} \mathcal{P}_B(s_{t-1} | s_t),$$

If \mathcal{F} is Markovian, then \mathcal{P}_F and \mathcal{P}_B can be uniquely defined as

$$\mathcal{P}_F(s' | s) = \frac{\mathcal{F}(s \rightarrow s')}{\mathcal{F}(s)}, \quad \mathcal{P}_B(s | s') = \frac{\mathcal{F}(s \rightarrow s')}{\mathcal{F}(s')}. \quad (2)$$

We call \mathcal{P}_F and \mathcal{P}_B the **forward policy** and the **backward policy**

Trajectory and detailed balance

The detailed balance constraint is a relation between them:

$$\mathcal{F}(s)\mathcal{P}_F(s'|s) = \mathcal{F}(s')\mathcal{P}_B(s|s'). \quad (3)$$

The trajectory balance constraint states that for any complete trajectory

$$\prod_{t=1}^{n_\tau} \mathcal{P}_F(s_t|s_{t-1}) = \frac{\mathcal{F}(s_{n_\tau})}{Z_{\mathcal{F}}} \prod_{t=1}^{n_\tau} \mathcal{P}_B(s_{t-1}|s_t). \quad (4)$$

Thus any Markovian flow can be uniquely determined from $Z_{\mathcal{F}}$ and $\mathcal{P}_F(-|s) \forall s \in \mathcal{S} \setminus \mathcal{X}$ or from $\mathcal{F}(x) \forall x \in \mathcal{X}$ and $\mathcal{P}_B(-|s) \forall s \in \mathcal{S} \setminus \{s_0\}$.

Training GFlowNet

To train GFlowNet we parameterize a Markovian flow by some model and train it to minimize some objective. If the objective is globally minimized, we should have $\mathcal{F}(x) = \mathcal{R}(x)$.

Trajectory Balance (TB): for each complete trajectory $\tau = (s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_{n_\tau} = x)$:

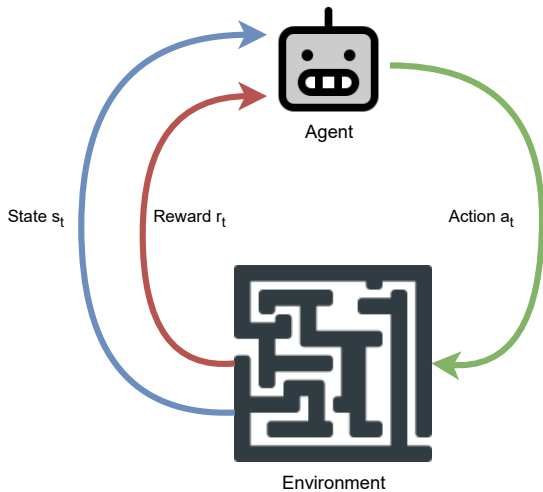
$$\mathcal{L}_{\text{TB}}(\tau) = \left(\log \frac{Z_\theta \prod_{t=1}^{n_\tau} \mathcal{P}_F(s_t | s_{t-1}, \theta)}{\mathcal{R}(x) \prod_{t=1}^{n_\tau} \mathcal{P}_B(s_{t-1} | s_t, \theta)} \right)^2. \quad (5)$$

Detailed Balance (DB): for each edge $s \rightarrow s' \in \mathcal{E}$:

$$\mathcal{L}_{\text{DB}}(s, s') = \left(\log \frac{\mathcal{F}_\theta(s) \mathcal{P}_F(s' | s, \theta)}{\mathcal{F}_\theta(s') \mathcal{P}_B(s | s', \theta)} \right)^2. \quad (6)$$

For terminal states $\mathcal{F}_\theta(x)$ is substituted by $\mathcal{R}(x)$ in the objective.

Reinforcement Learning



Bridging tabular and Deep RL

Tabular setting: $S \approx 100$;

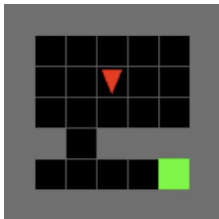


Figure: Minigrid environment

Deep RL setting: $S \approx 10^{100}$;



Figure: Atari Breakout

Is there an algorithm that at the same time

- provably optimal in tabular setting?
- practically good in Deep RL setting?

Soft Reinforcement learning

Markov Decision Process (MDP): $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma, s_0)$, where \mathcal{S} and \mathcal{A} are finite state and action space, P is a Markovian transition kernel, r is a bounded deterministic reward function, $\gamma \in [0, 1]$ is a discount factor, and s_0 is a fixed initial state. At step t

- state $s_t \in \mathcal{S}$;
- action $a_t \in \mathcal{A}$;
- next state $s_{t+1} \sim P(\cdot | s_t, a_t)$;
- reward $r(s_t, a_t)$ - known.

An RL agent interacts with MDP by a policy π that maps each state s to a probability distribution over possible actions.

Entropy-Regularized RL

Performance measure:

$$V_{\lambda}^{\pi}(s) \triangleq \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda \mathcal{H}(\pi(s_t))) \mid s_0 = s \right], \quad (7)$$

where λ is a regularization coefficient, \mathcal{H} is a Shannon entropy function, $a_t \sim \pi(s_t)$, $s_{t+1} \sim P(s_t, a_t)$ for all $t \geq 0$. In a similar manner we can define regularized Q-values $Q_{\lambda}^{\pi}(s, a)$ as an expected (discounted) sum of rewards augmented by Shannon entropy given fixed initial state $s_0 = s$ and action $a_0 = a$.

Goal: find a **regularized optimal policy** π_{λ}^{\star} that maximizes $V_{\lambda}^{\pi}(s)$ for any initial state s . In the case $\lambda = 1$ we will drop a subscript λ for brevity.

Soft Bellman Equations

Let V_λ^* and Q_λ^* be a value and Q-value of an optimal policy π_λ^* correspondingly. Then for any $s \in \mathcal{S}, a \in \mathcal{A}_s$

$$\begin{aligned} Q_\lambda^*(s, a) &\triangleq r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} [V_\lambda^*(s')] \\ V_\lambda^*(s) &\triangleq \text{LogSumExp}_\lambda(Q_\lambda^*(s, \cdot)), \end{aligned} \tag{8}$$

where $\text{LogSumExp}_\lambda \triangleq \lambda \log(\sum_{i=1}^d e^{x_i/\lambda})$, and an optimal policy is defined as a λ -softmax policy with respect to optimal Q-values $\pi_\lambda^*(s) = \text{Softmax}_\lambda(Q_\lambda^*(s, \cdot))$, where $[\text{Softmax}_\lambda(x)]_i \triangleq e^{x_i/\lambda} / (\sum_{j=1}^d e^{x_j/\lambda})$, or, alternatively

$$\pi_\lambda^*(a|s) = \exp\left(\frac{1}{\lambda}(Q_\lambda^*(s, a) - V_\lambda^*(s))\right). \tag{9}$$

Notice that as $\lambda \rightarrow 0$ we recover the well-known Bellman equations.

Soft Deep Q-Network

The goal of SoftDQN: approximate a solution to soft Bellman equations (8) using a Q-value parameterized by a neural network.

Let Q_θ be an online Q-network, it represents the current approximation of optimal regularized Q-value, and let $\pi_\theta = \text{Softmax}_\lambda(Q_\theta)$. Then, akin to DQN, agent interacts with MDP using policy π_θ (with an additional ε -greedy exploration or even without it), collect transitions (s_t, a_t, r_t, s_{t+1}) for $r_t = r(s_t, a_t)$ in a replay buffer \mathcal{B} .

Then SoftDQN performs stochastic gradient descent on the regression loss $\mathcal{L}_{\text{SoftDQN}}(\mathcal{B}) = \mathbb{E}_{\mathcal{B}}[(Q_\theta(s_t, a_t) - y_{\text{SoftDQN}}(r_t, s_{t+1}))^2]$ with a target Q-value

$$y_{\text{SoftDQN}}(r_t, s_{t+1}) = r_t + \gamma \text{LogSumExp}_\lambda(Q_{\bar{\theta}}(s_{t+1}, \cdot)), \quad (10)$$

where $Q_{\bar{\theta}}$ is a target Q-network with parameters $\bar{\theta}$ that are periodically copied from parameters of online Q-network θ .

GFlowNet Training as Soft RL Problem

Our main result is the reduction of GFlowNet learning problem to Soft RL:

- Construct a deterministic MDP based on GFlowNet DAG. Add an absorbing state s_f and add edges from all terminal states to it.
- For any s we select $\mathcal{A}_s \subseteq \mathcal{S}$ as a set of actions in state s in $\mathcal{S} \setminus \mathcal{X}$, and define $\mathcal{A}_x \triangleq \{s_f\}$ for any $x \in \mathcal{X} \cup \{s_f\}$.
- Deterministic Markov kernel is simple: $P(s'|s, a) = \mathbb{1}\{s' = a\}$ since a corresponds to a next state \mathcal{S} given a corresponding action. In the sequel we will write s' instead of a to emphasize this connection.
- Reward function

$$r(s, s') = \begin{cases} \log \mathcal{P}_B(s|s') & s \notin \mathcal{X} \cup \{s_f\}, \\ \log \mathcal{R}(s) & s \in \mathcal{X}, \\ 0 & s = s_f \end{cases} \quad (11)$$

- Set $\gamma = 1$.

GFlowNet Training as Soft RL Problem

Theorem

The optimal policy $\pi^(s'|s)$ for a regularized MDP with a coefficient $\lambda = 1$ is equal to $\mathcal{P}_F(s'|s)$ for all $s \in \mathcal{S} \setminus \mathcal{X}, s' \in \mathcal{A}_s$;*

Moreover, $\forall s \neq s_f, s' \neq s_f$ regularized optimal value $V^(s)$ and Q-value $Q^*(s, s')$ coincide with $\log \mathcal{F}(s)$ and $\log \mathcal{F}(s \rightarrow s')$ respectively, where \mathcal{F} is the Markovian flow defined by \mathcal{P}_B and the GFlowNet reward \mathcal{R} .*

Soft DQN in the realm of GFlowNets

During the interaction, the agent collects transitions (s_t, a_t, r_t, s_{t+1}) into the replay buffer \mathcal{B} and optimizes a regression loss $\mathcal{L}_{\text{SoftDQN}}$ using targets (10) with rewards $r_t = \log \mathcal{P}_{\mathcal{B}}(s_t | s_{t+1})$ for any $s_t \notin \mathcal{X}$. For terminal transitions with $s_t \in \mathcal{X}$ and $s_{t+1} = s_f$ the target is equal to $\log \mathcal{R}(s_t)$.

Using an interpretation $Q_\theta(s, s') = \log \mathcal{F}_\theta(s \rightarrow s')$ we obtain the following loss

$$\mathcal{L}_{\text{SoftDQN}}(\mathcal{B}) = \mathbb{E}_{\mathcal{B}} \left[\left(\log \frac{\mathcal{F}_\theta(s_t \rightarrow s_{t+1})}{\mathcal{P}_{\mathcal{B}}(s_t | s_{t+1}) \mathcal{F}_{\bar{\theta}}(s_{t+1})} \right)^2 \right], \quad (12)$$

where $\mathcal{F}_{\bar{\theta}}(x) = \mathcal{R}(x) \forall x \in \mathcal{X}$, and otherwise $\mathcal{F}_{\bar{\theta}}(s) \triangleq \sum_{s'} \mathcal{F}_{\bar{\theta}}(s \rightarrow s')$ and $\bar{\theta}$ are parameters of a target network that is updated with weights θ from time to time.

Interpretation of Value

For a policy π in a deterministic MDP $\mathcal{M}_{\mathcal{G}}$ with all trajectories of length $n_{\tau} \leq N$ we can define induced trajectory distribution as follows

$$q^{\pi}(\tau) = \prod_{i=1}^N \pi(s_i | s_{i-1}) = \prod_{i=1}^{n_{\tau}} \pi(s_i | s_{i-1}),$$

where $\tau = (s_0 \rightarrow \dots \rightarrow s_N)$ with $s_N = s_f$ and $s_{n_{\tau}} \in \mathcal{X}$. The last identity holds since $\pi(s_f | s_f) = \pi(s_f | x) = 1$ for any $x \in \mathcal{X}$. Additionally, we defined a marginal distribution of sampling by policy π as $d^{\pi}(x) = \mathbb{P}_{\pi}[s_{n_{\tau}} = x]$.

Additionally, for a fixed GFlowNet reward \mathcal{R} and a backward policy \mathcal{P}_B we define the corresponding distribution over trajectories with a slight abuse of notation

$$\mathcal{P}_B(\tau) = \prod_{i=1}^{n_{\tau}} \mathcal{P}_B(s_{i-1} | s_i) \cdot \frac{\mathcal{R}(s_{n_{\tau}})}{Z}.$$

Notice that by trajectory balance constraint (4) $\mathcal{P}_B(\tau)$ is equal to a $\mathcal{P}(\tau)$ and, by Theorem 1, to $q^{\pi^*}(\tau)$.

Interpretation of Value

Then for any policy π its regularized value V^π in the initial state s_0 is equal to

$$V^\pi(s_0) = \log Z - \text{KL}(q^\pi \| \mathcal{P}_B).$$

Takeaway: Maximization of the value in the prescribed MDP \mathcal{M} is equivalent to minimization the KL-divergence between trajectory distribution induced by the current policy and by a backward policy.

Moreover, it shows that all near-optimal policies in soft RL-sense introduce a flow with close trajectory distributions, and, moreover, close marginal distribution:

$$V^*(s_0) - V^\pi(s_0) = \text{KL}(q^\pi \| \mathcal{P}_B) \geq \text{KL}(d^\pi \| \mathcal{R}(x)/Z),$$

Hypergrid Environment

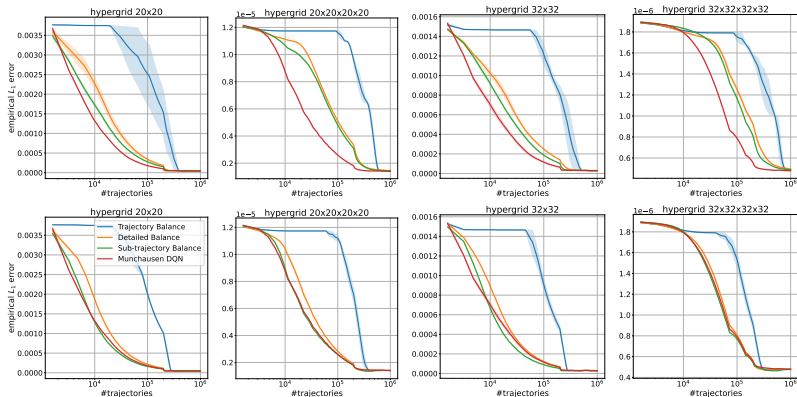


Figure: L_1 distance between target and empirical GFlowNet distributions over the course of training on hypergrid environment. *Top row:* \mathcal{P}_B is fixed to be uniform for all methods. *Bottom row:* \mathcal{P}_B is learnt for the baselines and fixed to be uniform for M-DQN. Mean and std values are computed over 3 runs.

Small Molecule Generation

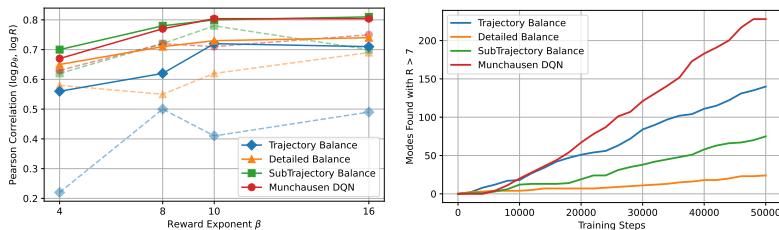


Figure: Small molecule generation results. *Left:* Pearson correlation between $\log \mathcal{R}$ and $\log \mathcal{P}_{\theta}$ on a test set for each method and varying $\beta \in \{4, 8, 10, 16\}$. Solid lines represent the best results over choices of learning rate, dashed lines — mean results. *Right:* Number of Tanimoto-separated modes with $\tilde{\mathcal{R}} > 7.0$ found over the course of training for $\beta = 10$.

Bit sequence generation

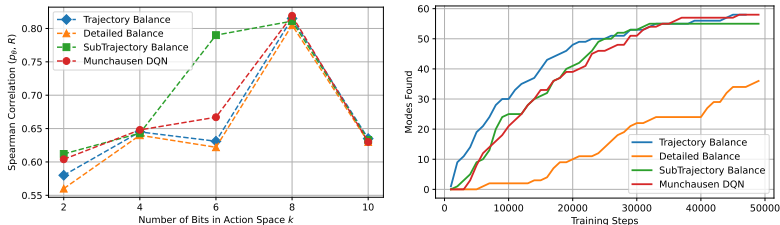


Figure: Bit sequence generation results. *Left:* Spearman correlation between \mathcal{R} and \mathcal{P}_θ on a test set for each method and varying $k \in \{2, 4, 6, 8, 10\}$. *Right:* The number of modes discovered over the course of training for $k = 8$.

Problem

We can interpret chain-of-thought reasoning, a paradigm of reasoning in language models, as a problem of intractable posterior inference. Given a question-answer pair (X, Y) , we are interested in finding latent chains of thought – token sequences Z that contribute the most to the conditional likelihood

$$p(Y|X) = \sum_Z p_{LM}(ZY|X) = \sum_Z p_{LM}(Y|XZ)p_{LM}(Z|X),$$

where p_{LM} denotes the likelihood assigned to a sequence by a language model and apposition of variables (e.g., XZY) denotes the concatenation of the token sequences. We are interested in sampling from

$$p(Z|X, Y) = \frac{p_{LM}(XZY)}{\sum_{Z'} p_{LM}(XZ'Y)}$$

Amortizing intractable inference in large language models

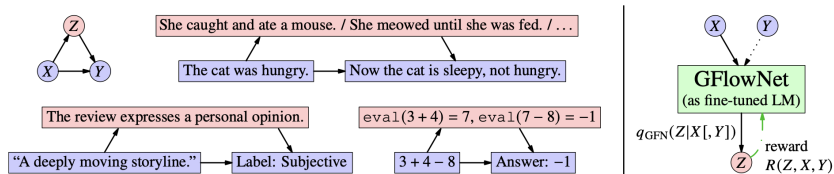


Figure 1: **Left:** Three problems of reasoning in language – sentence infilling, chain-of-thought reasoning, and problem-solving with external tool use – can all be seen as instances of the latent variable model at the top left, where an input (X) generates the output (Y) via a latent variable (Z). **Right:** We fine-tune an LLM to sample from the Bayesian posterior over Z , conditioning on X and optionally on Y . If conditioned on Y , the trained policy can be used to sample diverse latent sequences (e.g., for infilling, §4.2). If not conditioned on Y , the policy can sample Z , and thus predict Y , for inputs X not seen during training (e.g., for classification and multi-step reasoning, §4.3, 4.4). As shown in §4.4, modeling the full diversity of the posterior aids generalization.

Figure: See <https://arxiv.org/pdf/2310.04363.pdf>

Thank you!

Bibliography I



Bengio, E., Jain, M., Korablyov, M., Precup, D., and Bengio, Y. (2021). [Flow network based generative models for non-iterative diverse candidate generation](#). [Advances in Neural Information Processing Systems](#), 34:27381–27394.



Bengio, Y., Lahlou, S., Deleu, T., Hu, E. J., Tiwari, M., and Bengio, E. (2023). [Gflownet foundations](#). [Journal of Machine Learning Research](#), 24(210):1–55.



Samsonov, S., Lagutin, E., Gabri  , M., Durmus, A., Naumov, A., and Moulines, E. (2022). Local-global mcmc kernels: the best of both worlds. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A., editors, [Advances in Neural Information Processing Systems](#), volume 35, pages 5178–5193. Curran Associates, Inc.