

Pessimistic Backward Policy for GFlowNets

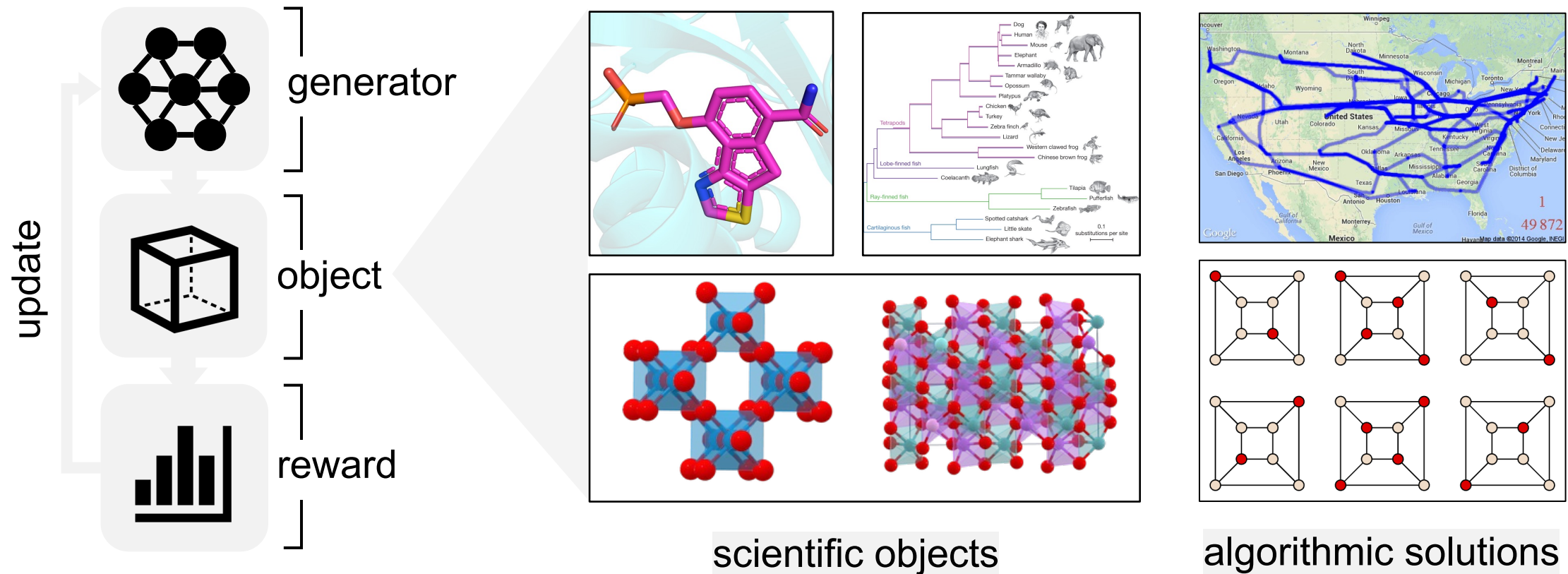
Hyosoon Jang¹ Yunhui Jang¹ Minsu Kim² Jinkyoo Park² Sungsoo Ahn¹

¹POSTECH ²KAIST



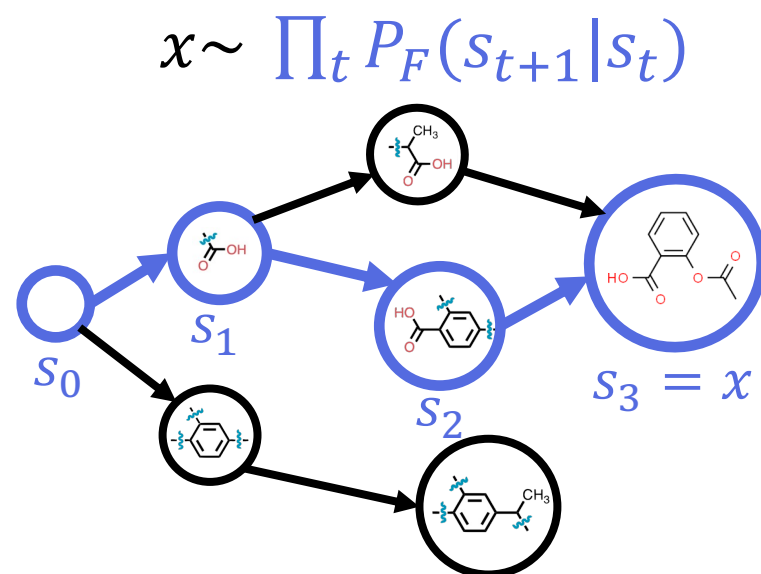
Function-driven Learning

- Given a reward function, we aim to obtain diverse high-reward objects



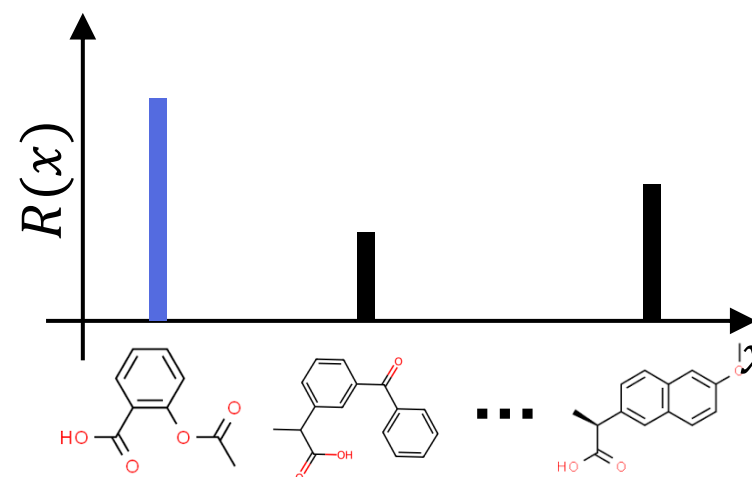
Generative Flow Networks (GFlowNets)

- GFlowNets construct objects with a **forward policy** P_F :



- compositionality in the generation:**
inducing an object x with a trajectory of states $\tau = (s_0, \dots, s_T = x)$

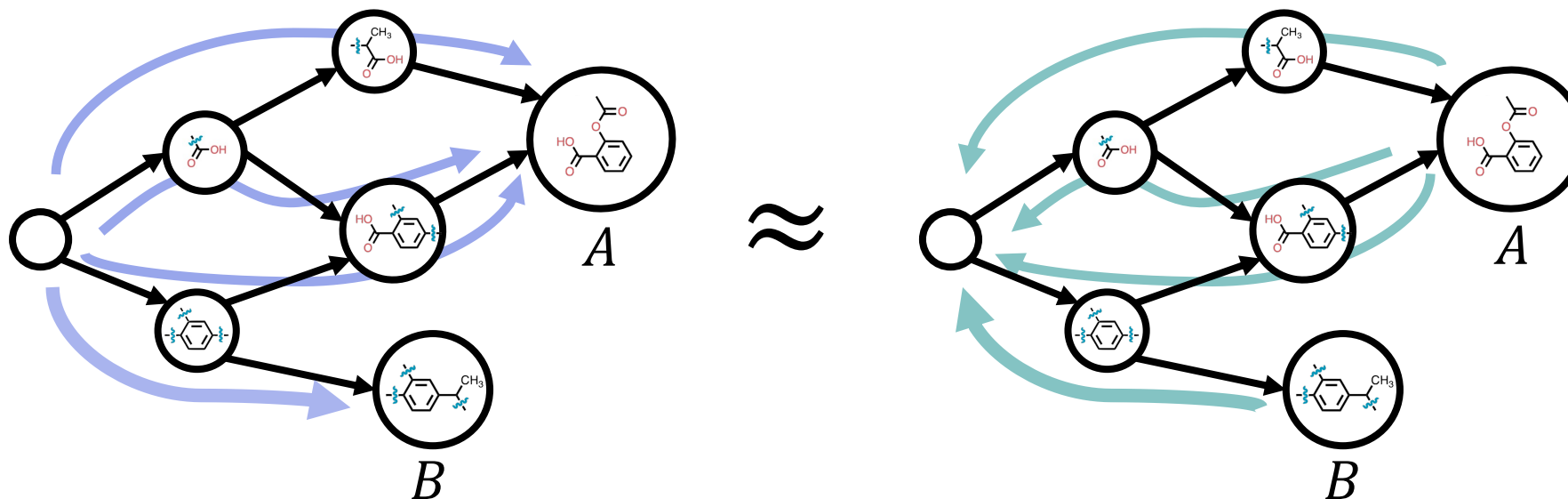
$P_F(x) \propto R(x)$



- diversity in the generation:**
sampling with a probability proportional to the reward $R(x)$

Training of GFlowNets

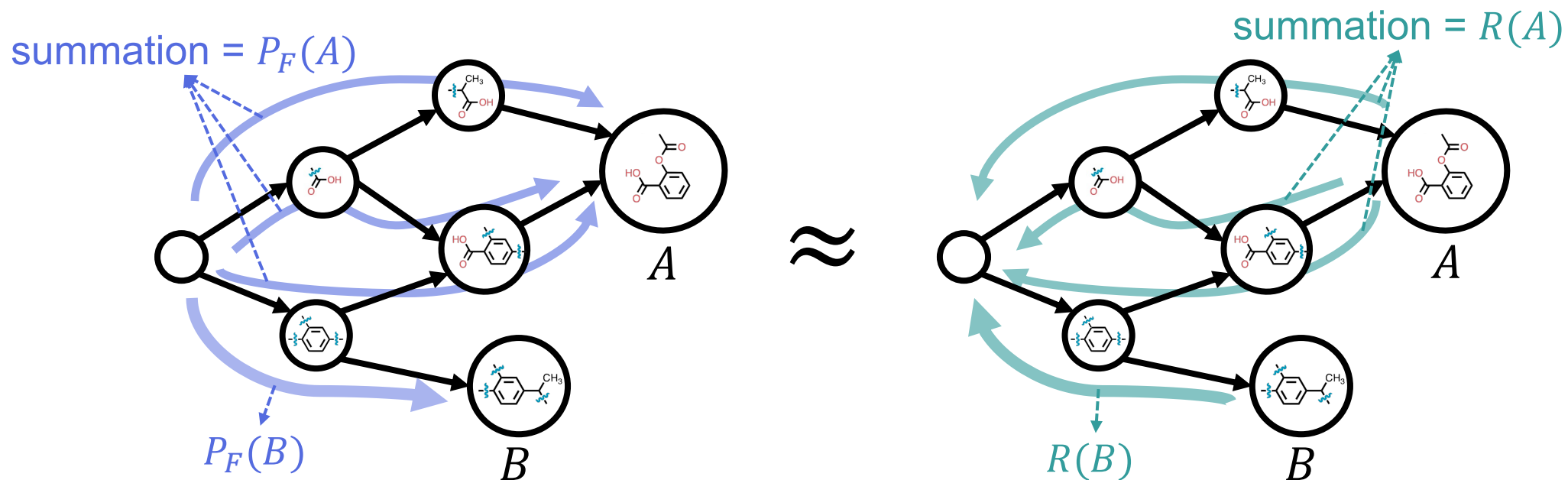
- GFlowNets align the **forward policy** with a **backward policy** over trajectories
 - similar to aligning diffusion with reverse diffusion or encoder with decoder



$$\Rightarrow P_F(x) \propto R(x)$$

Training of GFlowNets

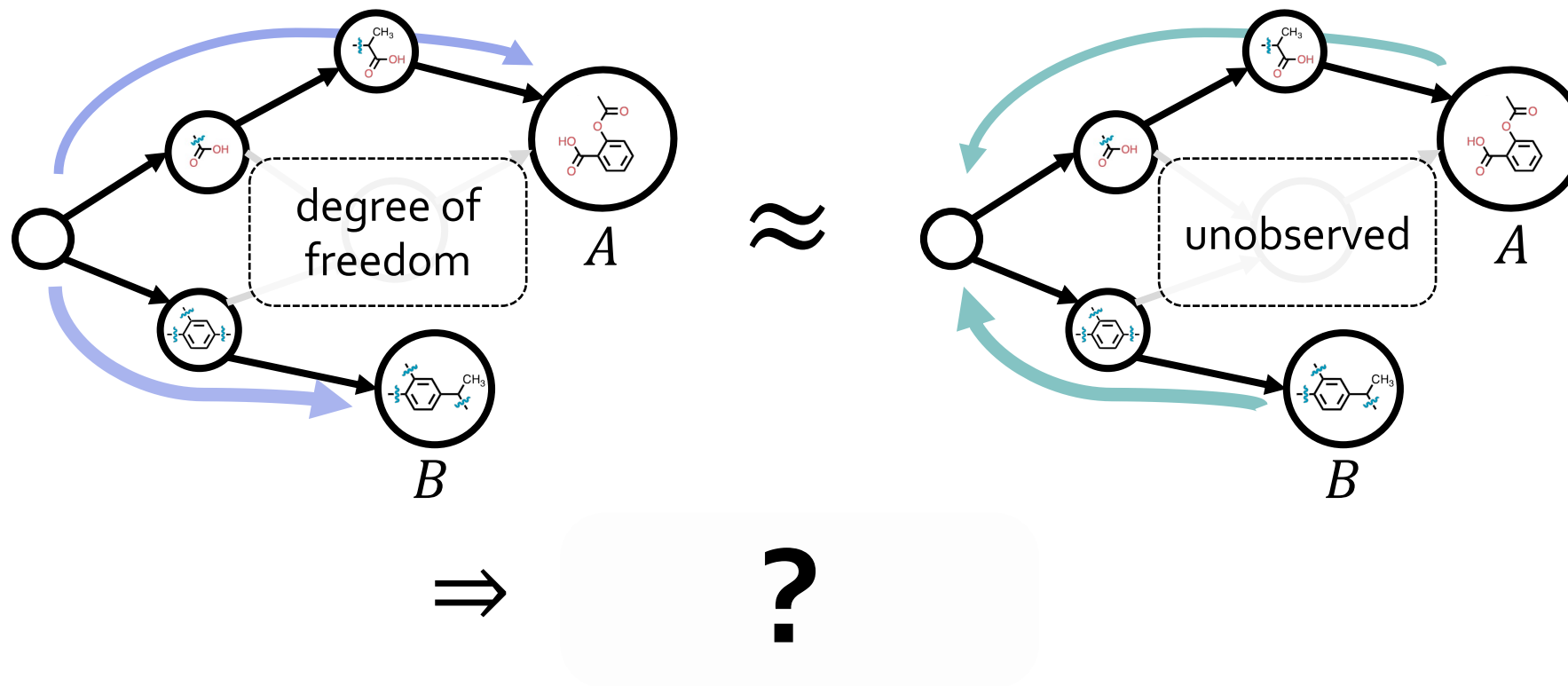
- The **forward policy** captures the **target densities** over trajectories that decompose the **object's reward** through the **backward policy**



$$\Rightarrow R(A) > R(B) \text{ and } P_F(A) > P_F(B)$$

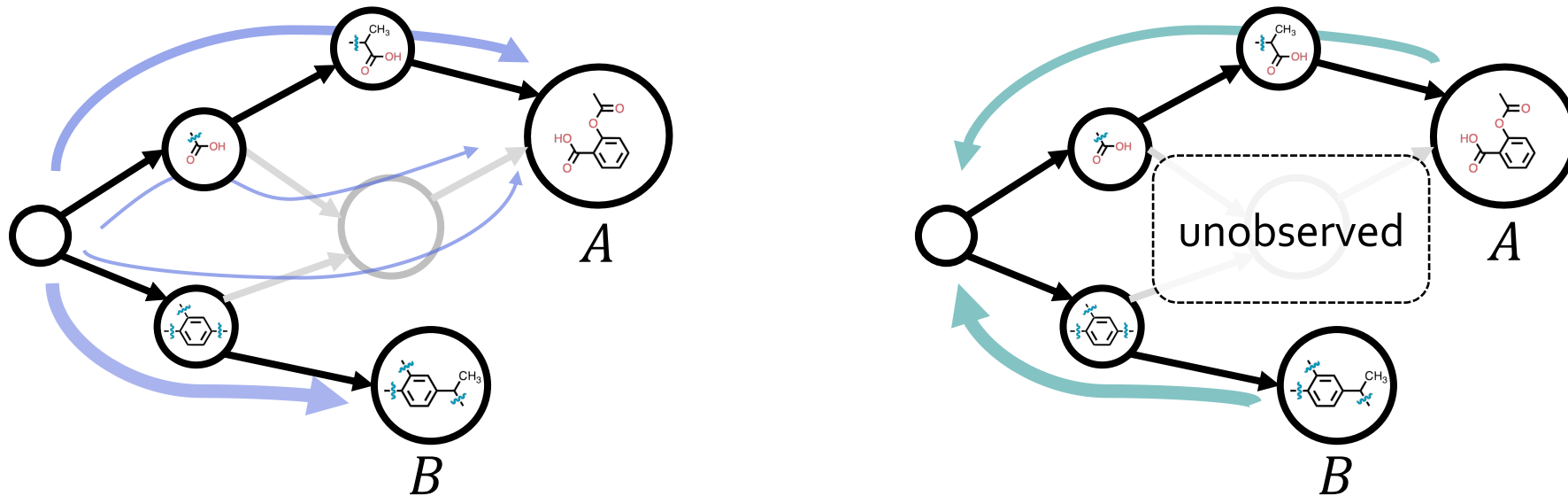
Motivation: Under-determined Probability

- Unobserved trajectories under-determine the sampling probability
 - a degree of freedom in the **forward probability** over unobserved trajectories



Motivation: Under-determined Probability

- This can cause the **forward policy** to favor low-reward objects
 - due to the degree of freedom lower than unobserved target densities

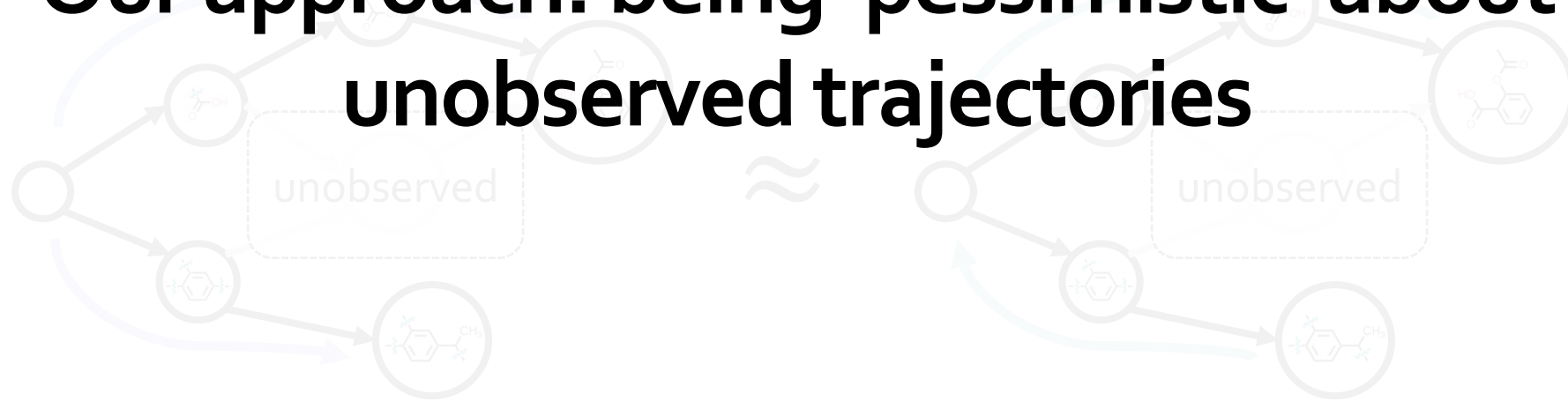


$$\Rightarrow R(A) > R(B) \text{ but } P_F(A) < P_F(B)$$

Motivation: Under-determined Probability

- Unobserved trajectories **under-determine** the sampling probability
 - even make GFlowNets to favor low-reward objects

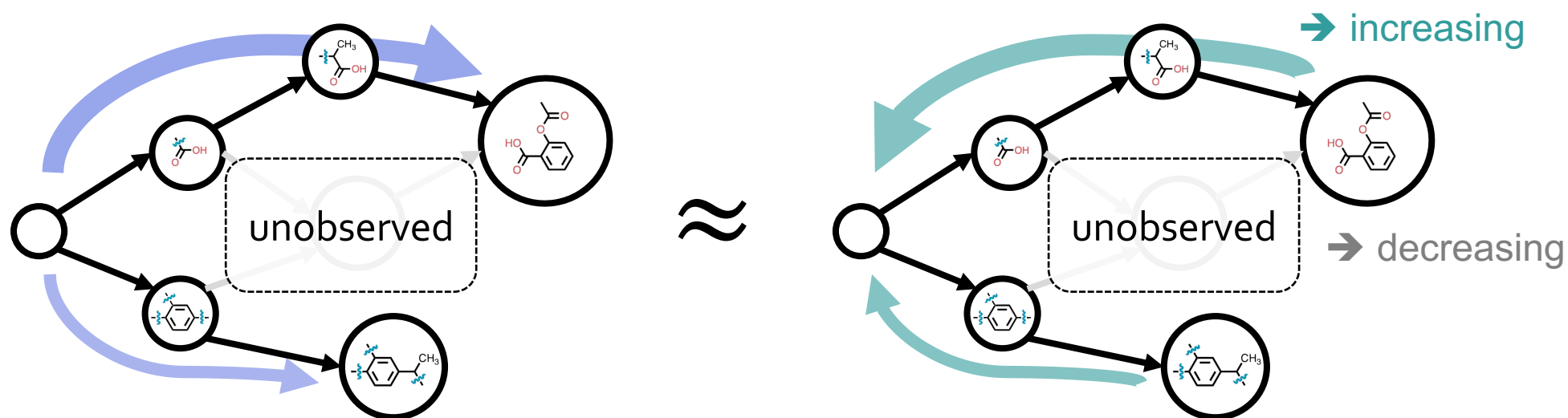
Our approach: being 'pessimistic' about unobserved trajectories



$$\Rightarrow P(\alpha) \stackrel{?}{\propto} R(\alpha)$$

Our: Pessimistic Backward Policy

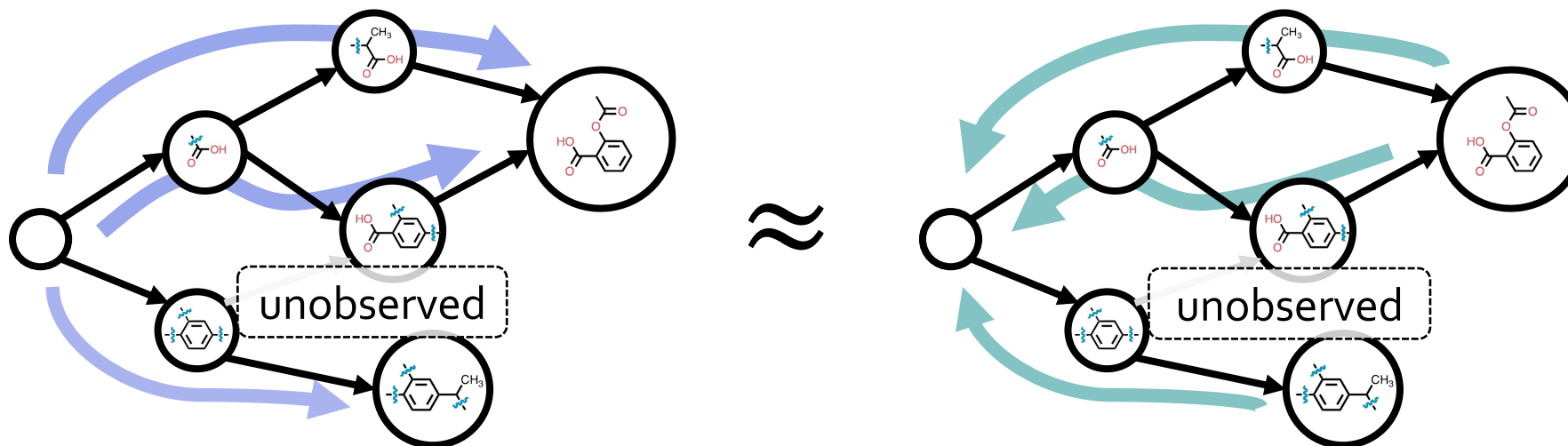
- For an observed object, **our backward policy** minimizes and increases the unobserved and **observed** target densities, respectively.



$$\Rightarrow R(A) > R(B) \text{ and } P_F(A) > P_F(B)$$

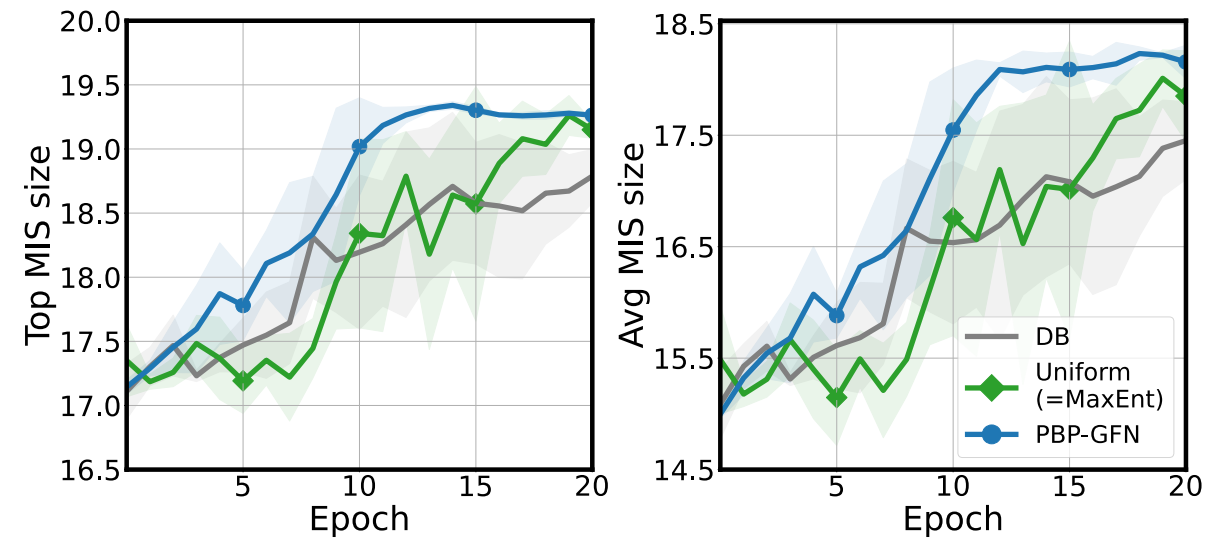
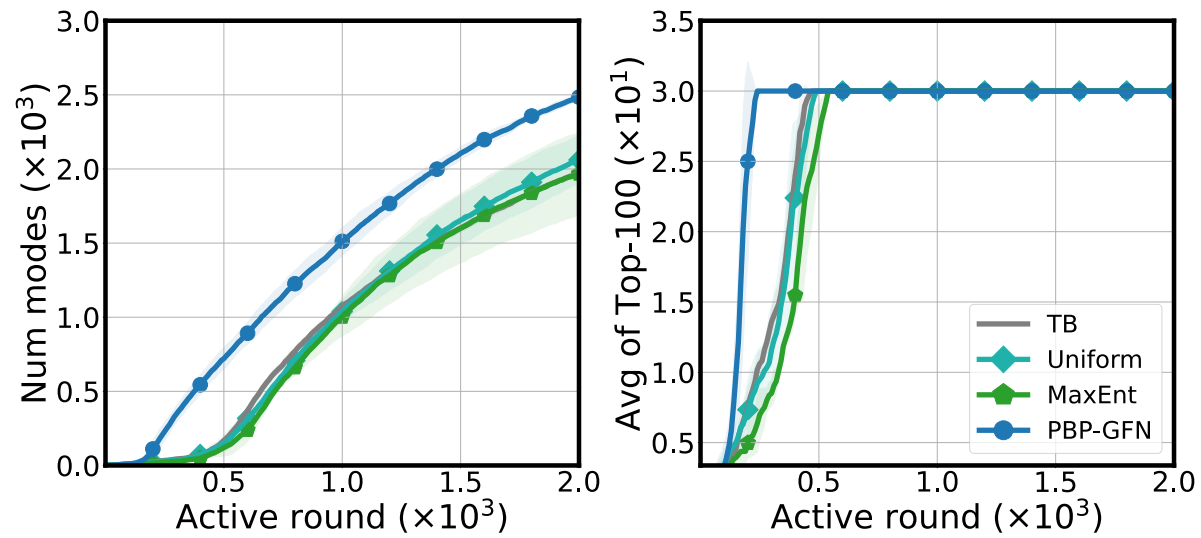
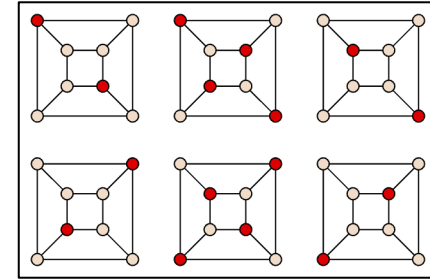
Our: Pessimistic Backward Policy

- In each training round, we adapt our backward policy for the observed trajectories

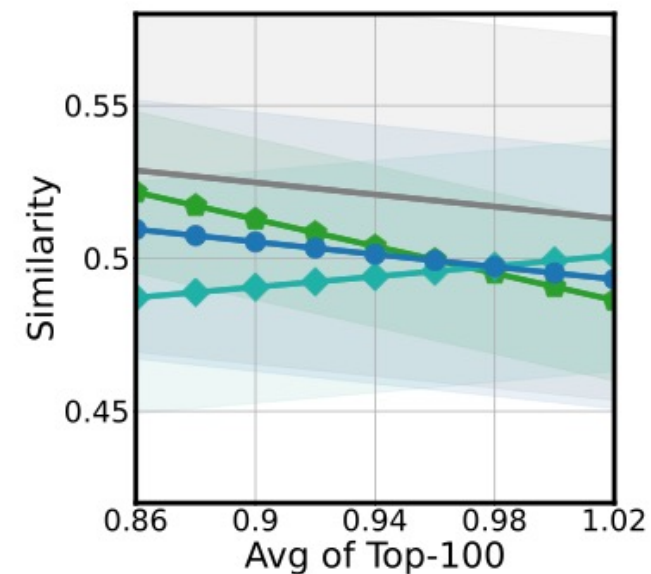
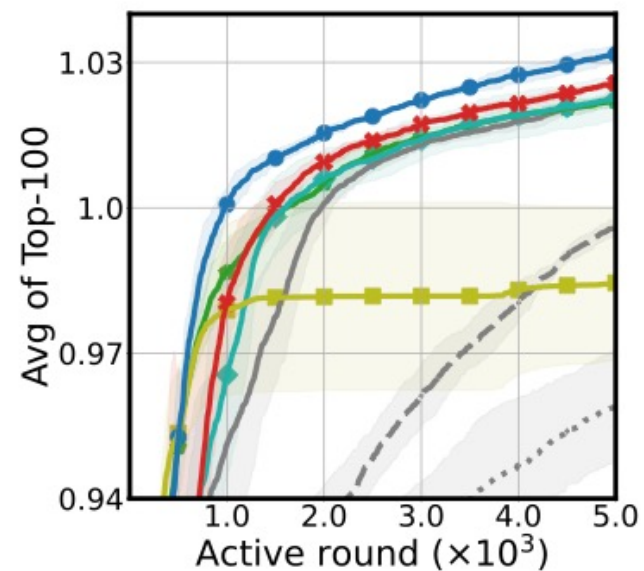
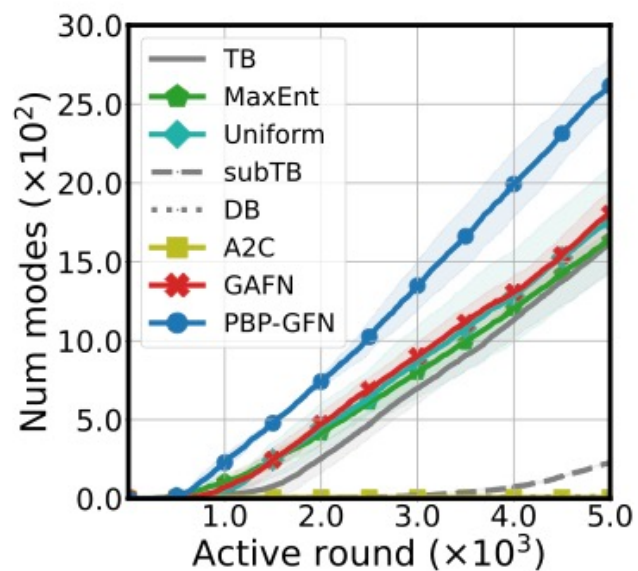
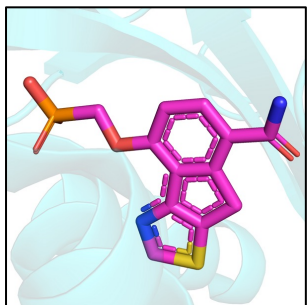


$$\Rightarrow R(A) > R(B) \text{ and } P_F(A) > P_F(B)$$

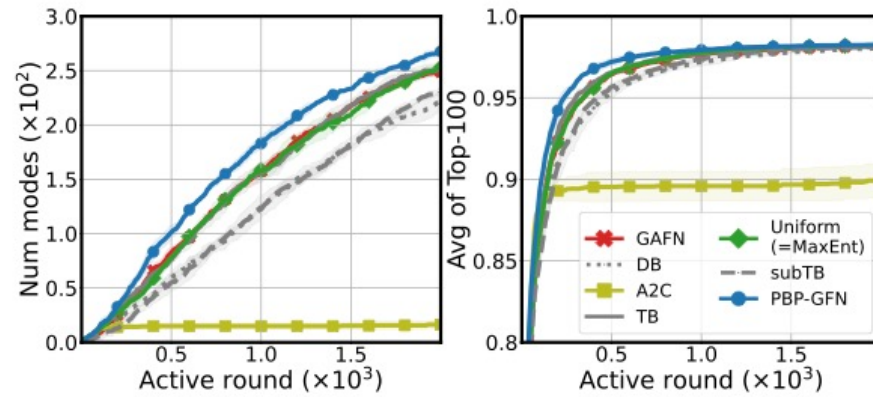
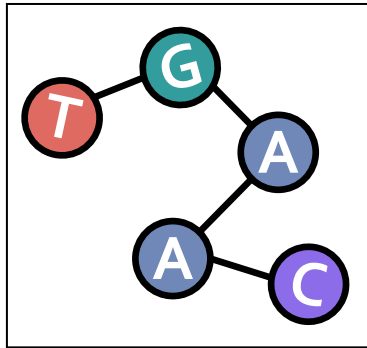
Experiments: Bag and Maximum Independent Set



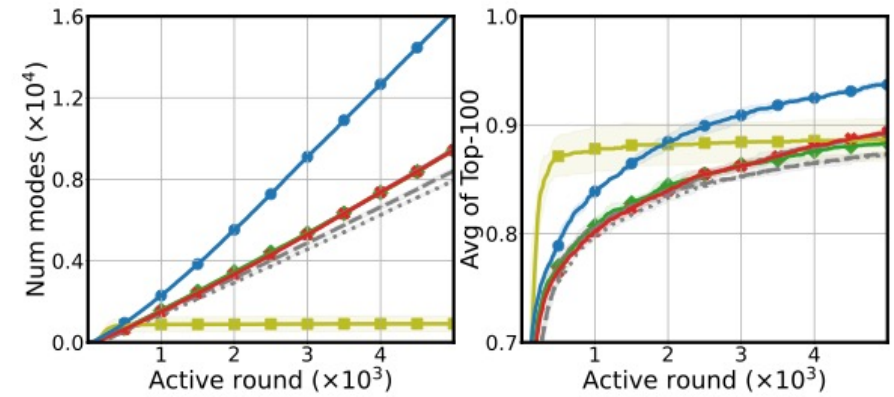
Experiments: Molecule Generation



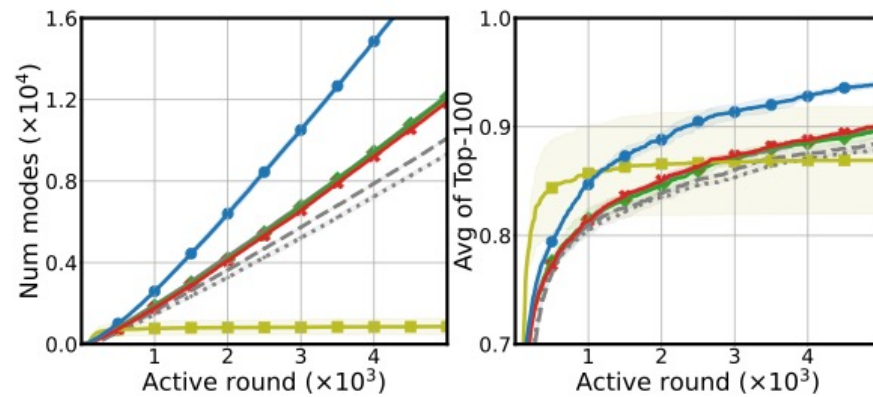
Experiments: RNA sequence



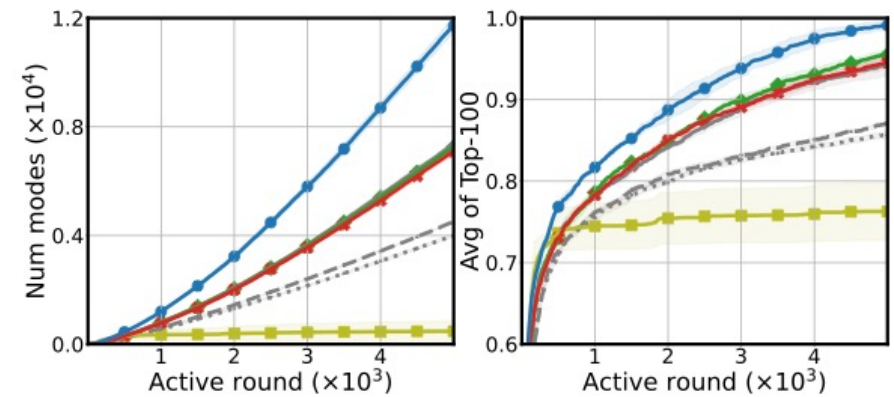
(a) TFBind8



(b) RNA-A



(c) RNA-B



(d) RNA-C