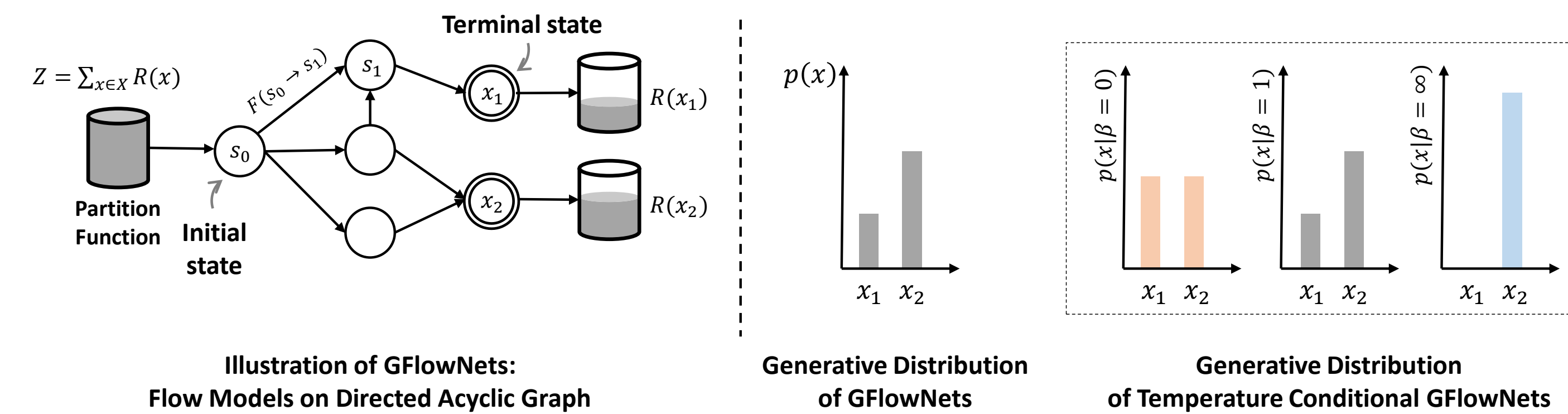


# Learning to Scale Logits for Temperature-Conditional GFlowNets

We propose a simple technique to adjust the SoftMax logit temperature for temperature-conditional GFlowNets and demonstrate its useful applications in drug discovery.

## 1. Preliminary



### 1.1 What are GFlowNets?

GFlowNets aim to convert a non-negative reward function  $R(x)$  into a generative policy:  $p(x) \propto R^\beta(x)$  with inverse temperature  $\beta$  which determine steepness of target distribution.

### 1.2 Temperature Conditional GFlowNets

Our aim is to learn the temperature-conditional GFlowNets (Temp-GFN):

$$p(x|\beta) \propto R^\beta(x)$$

The major benefit of temperature-conditional GFlowNets is the controllability of GFlowNets' exploration and exploitation through adjusting temperature.

## 2. Research Objective

### Research Scope 1: Training stability of Temp-GFN

- Temp-GFN training is challenging due to varying target distribution steepness.
- Representing various temperature distributions can be numerically difficult and unstable.
- Incorporating additional inductive bias into SoftMax temperature may stabilize Temp-GFN training.**

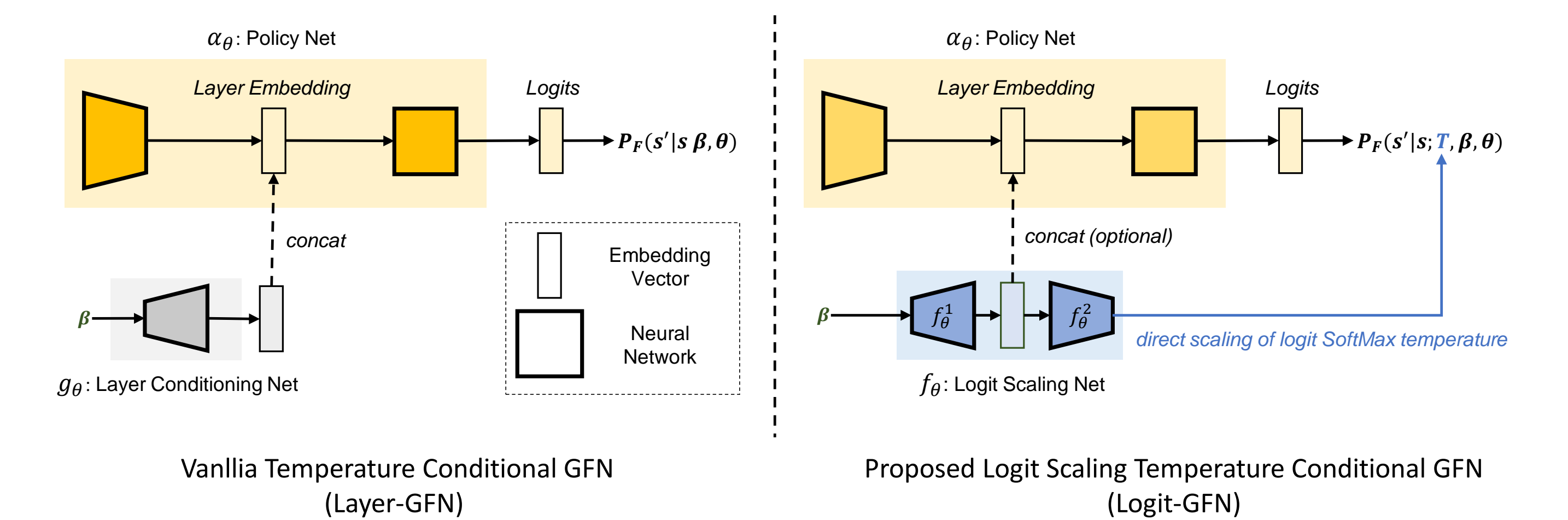
### Research Scope 2: Usefulness of Temp-GFN

- Training GFN at low temperatures (steep distributions) is challenging.
- Generalizing to steep, high-reward regions is difficult with limited observations.
- Training at various high temperatures (low  $\beta$ ) and querying at low temperatures (high  $\beta$ ) may generalize better than fixed low temperature training.**

## 3. Methodology: Logit-GFN

### 3.1 Overall Idea

Introducing the Logit-Scaling Net: a novel auxiliary network that optimally scales the SoftMax temperature of policy network logits, enhancing Temp-GFN training and enabling effective generalization across various temperature distributions.



### 3.2 Logit Scaling

**Objective:** Facilitate training of temperature-conditional GFlowNets,  $p(x|\beta) \propto R(x)^\beta$ , over varying inverse temperatures  $\beta$ .

**Logit-Scaling Trick:** Uses a skip connection to adjust the SoftMax temperature  $T$  of  $P_F$  logits based on  $\beta$ . Defined as:

$$P_F(s'|s; \beta, \theta) := \frac{\exp(\alpha_\theta(s, s')/f_\theta(\beta))}{\sum_{s'' \in \text{Ch}(s)} \exp(\alpha_\theta(s, s'')/f_\theta(\beta))} \quad (1)$$

- $\alpha_\theta(s, s')$ : Neural net independent of  $\beta$
- $f_\theta(\beta)$ : Logit-scaling net transforming  $\beta$  into SoftMax temperature  $T$

**Key Point:** Logit-scaling is agnostic to the policy network's form, including layer-conditioning networks.

(Note that the logit-scaling can be applied to layer-conditioning networks (where  $\alpha$  is also dependent on  $\beta$ ), ensuring their **full expressive capacity** is maintained across different temperature regimes.)

### 3.3 Training Objective

**Training Procedure.** We minimize the trajectory balance (TB) loss with a replay buffer  $\mathcal{D}$ , but train GFlowNets over multiple  $\beta \sim P_{\text{train}}(\beta)$ :

$$\mathcal{L}(\theta; \mathcal{D}) = \mathbb{E}_{P_{\text{train}}(\beta)} \mathbb{E}_{P_{\mathcal{D}}(\tau)} \left[ \left( \log \frac{Z_\theta(\beta) \prod_{t=1}^n P_F(\cdot)}{R(x)^\beta \prod_{t=1}^n P_B(\cdot)} \right)^2 \right], \quad (2)$$

where  $P_F(\cdot) = P_F(s_t|s_{t-1}; \beta, \theta)$  and  $P_B(\cdot) = P_B(s_{t-1}|s_t; \beta, \theta)$ .

**Implementation.** We condition the partition function  $Z_\theta(\beta)$  on  $\beta$  and use DNNs  $f_\theta, Z_\theta$  to map scalars, minimizing parameter overhead.

### 3.4 Online discovery algorithm with Logit-GFN

**Objective.** Discover diverse, high-reward candidates  $x \in \mathcal{X}$  (e.g., molecules with high binding affinity) focusing on top rewards, diversity, and modes.

**Method.** Enhance GFlowNets exploration by sampling multiple  $\beta$  values, forming a diverse replay buffer  $\mathcal{D}$ :

$$\begin{aligned} \mathcal{D} &\leftarrow \mathcal{D} \cup \{\tau_1, \dots, \tau_M\} \\ \tau_1, \dots, \tau_M &\sim \int_{\beta} P_F(\tau|\beta) dP_{\text{exp}}(\beta). \end{aligned} \quad (3)$$

Dynamic control policy  $P_{\text{exp}}(\beta)$  varies exploration range.

**Algorithm 1** Scientific Discovery with Temperature-Conditional GFlowNets

```

1: Set  $\mathcal{D} \leftarrow \emptyset$  ▷ Initialize dataset.
2: for  $t = 1, \dots, T$  do ▷ Training  $T$  rounds
3:    $\beta_1, \dots, \beta_M \sim P_{\text{exp}}(\beta)$  ▷ Sample temperatures from exploration query prior
4:   for  $m = 1, \dots, M$  do
5:      $\tau_m \sim P_F(\tau|\beta = \beta_m; \theta)$  ▷ Sample trajectories from Logit-GFN.
6:    $\mathcal{D} \leftarrow \mathcal{D} \cup \{\tau_m\}$ 
7: end for
8: for  $k = 1, \dots, K$  do ▷ Training  $K$  epochs per each training rounds
9:   Use ADAM for gradually minimizing  $\mathcal{L}(\theta; \mathcal{D})$ .
10: end for
11: end for
12: Output:  $\mathcal{D}$ 

```

## 4. Experiments

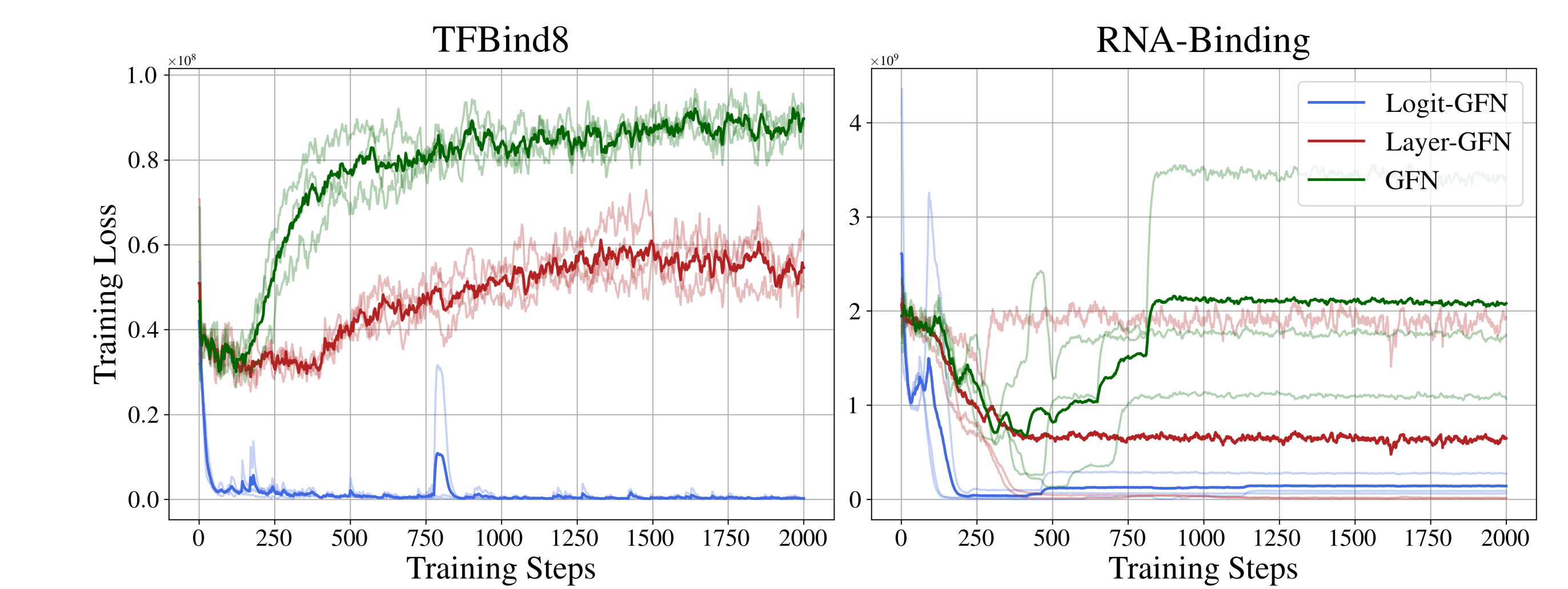
### 4.1 Tasks

We evaluate temperature-conditional GFlowNets on a toy grid world and four biochemical tasks: QM9, TFBind8, sEH, and RNA-binding.

- Grid world:** The agent can move forward but not backward to reach a high-reward goal state. With multiple high-reward modes, the agent aims to cover all modes in the grid space.
- QM9:** Generation of molecular graphs by sequentially adding atom components. The reward function is based on the Homo-Lumo gap.
- TFBind8:** Generation of DNA sequences by bidirectional token addition. The reward function measures binding activity with human transcription factors.
- sEH:** Generation of molecular graphs by sequentially adding predefined fragment components. The reward function is specified by enzymatic activity metrics.
- RNA-binding:** Generation of RNA sequences by bidirectional token addition. The reward function is based on binding activity with human transcription factors.

### 4.2 Training Stability

We assess the training stability of temperature-conditional GFlowNets.

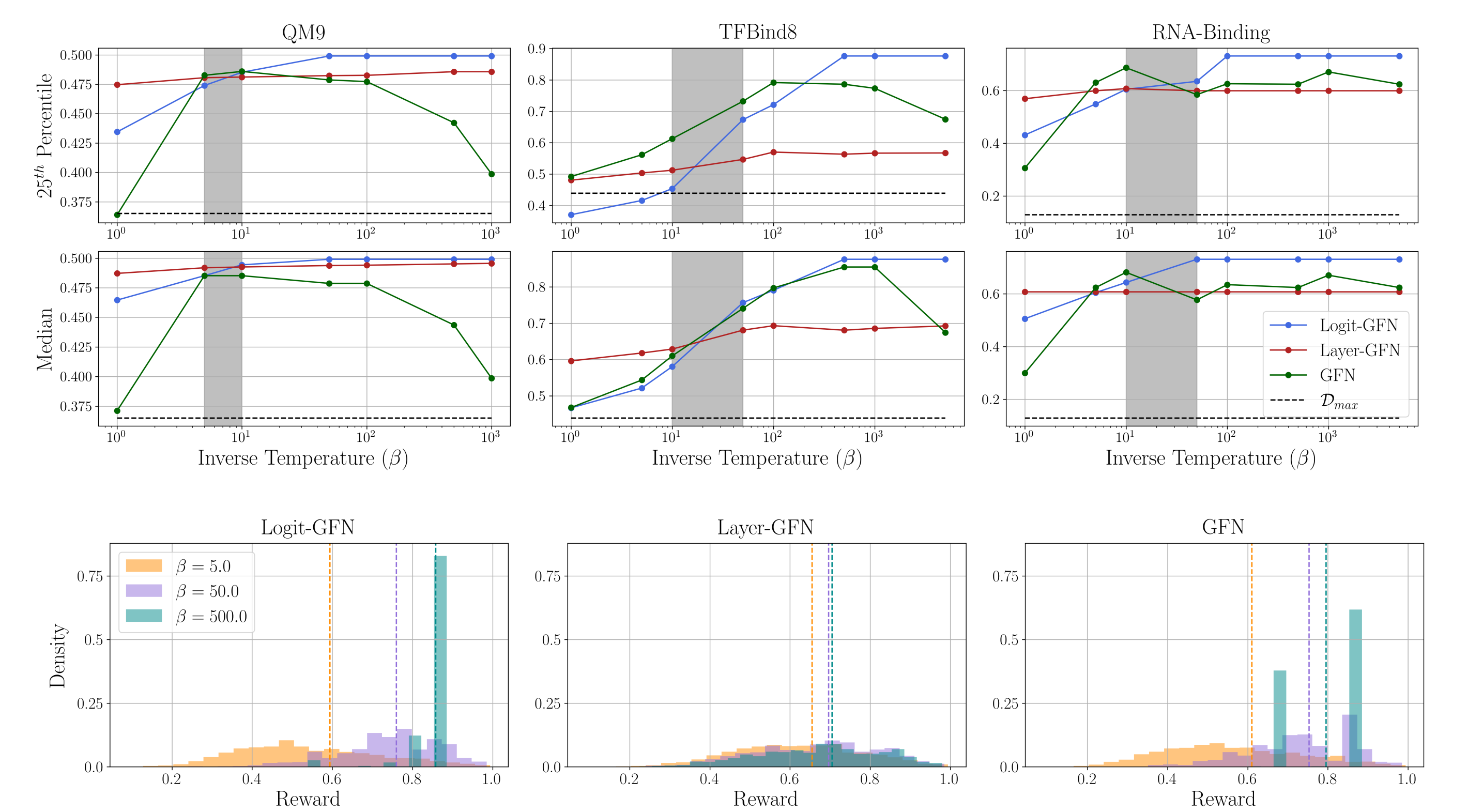


### 4.3 Offline Generalization

We examine the controllability of temperature-conditional GFlowNets, i.e.,  $p(x|\beta) \propto R(x)^\beta$ , in offline model-based optimization.

Minsu Kim<sup>1,2,\*</sup>, Joohwan Ko<sup>1,\*</sup>, Taeyoung Yun<sup>1,\*</sup>, Dinghuai Zhang<sup>2</sup>, Ling Pan<sup>3</sup>  
Woo Chang Kim<sup>1</sup>, Jinkyoo Park<sup>1</sup>, Emmanuel Bengio<sup>4</sup>, Yoshua Bengio<sup>2</sup>

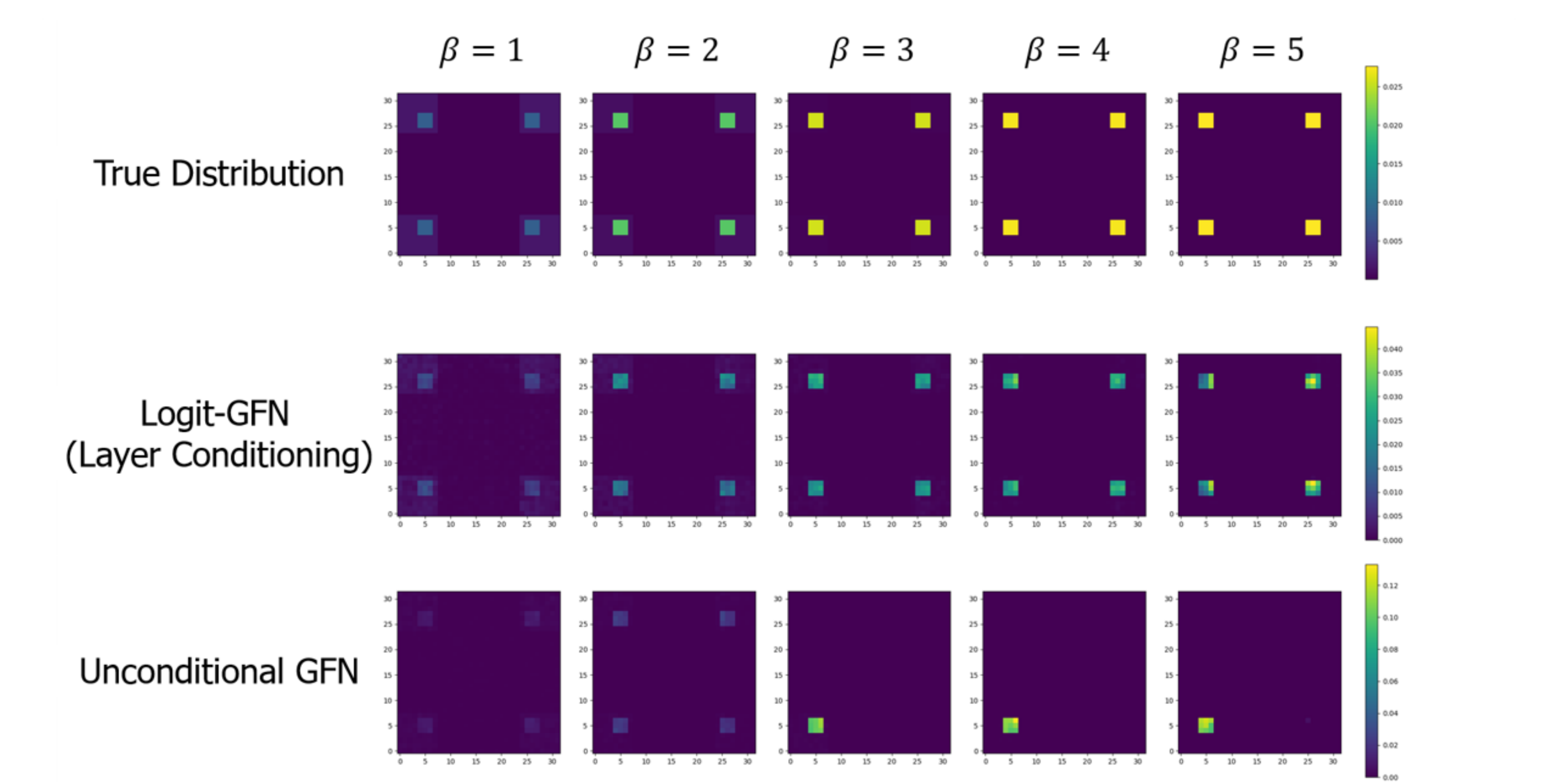
<sup>1</sup>KAIST <sup>2</sup>Mila <sup>3</sup>HKUST <sup>4</sup>Valance Labs \* Equal contribution



### 4.4 Online Mode-Seeking

We validate the effectiveness of temperature-conditional GFlowNets in solving online mode-seeking problems.

### Toy task: Grid world



### Biochemical discovery tasks

