GENERATIVE FLOWS ON SYNTHETIC PATHWAY FOR DRUG DESIGN

Anonymous authors

Paper under double-blind review

ABSTRACT

Generative models in drug discovery have recently gained attention as efficient alternatives to brute-force virtual screening. However, most existing models do not account for synthesizability, limiting their practical use in real-world scenarios. In this paper, we propose RXNFLOW, which sequentially assembles molecules using predefined molecular building blocks and chemical reaction templates to constrain the synthetic chemical pathway. We then train on this sequential generating process with the objective of generative flow networks (GFlowNets) to generate both highly rewarded and diverse molecules. To mitigate the large action space of synthetic pathways in GFlowNets, we implement a novel action space subsampling method. This enables RXNFLOW to learn generative flows over extensive action spaces comprising combinations of 1.2 million building blocks and 71 reaction templates without significant computational overhead. Additionally, RXNFLOW can employ modified or expanded action spaces for generation without retraining, allowing for the introduction of additional objectives or the incorporation of newly discovered building blocks. We experimentally demonstrate that RXNFLOW outperforms existing reaction-based and fragment-based models in pocket-specific optimization across various target pockets. Furthermore, RXNFLOW achieves state-of-the-art performance on CrossDocked2020 for pocket-conditional generation, with an average Vina score of -8.85 kcal/mol and 34.8% synthesizability. Code is available at https://anonymous.4open.science/r/RxnFlow-B13E/.

029 030 031

032

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

1 INTRODUCTION

Structure-based drug discovery (SBDD) has emerged as a pivotal paradigm for early drug discovery (structure-based drug discovery (structure-based drug discovery) facilitated by the increasing accessibility of protein structure prediction tools (Jumper et al., 2021) and high-resolution crystallography (Liu et al., 2015). However, traditional brute-force virtual screening is computationally expensive (Graff et al., 2021), prompting the development of deep generative models that can bypass this inefficiency. In this context, various approaches such as deep reinforcement learning (Zhavoronkov et al., 2019), variational autoencoders (Zhung et al., 2024) generative adversarial network (Ragoza et al., 2022), and diffusion models (Guan et al., 2023a;b) have been proposed to directly sample candidate molecules against a given protein structure.

While generative models have shown success in molecular discovery with desirable biological properties, most overlook synthesizability which is a crucial factor for wet-lab validation (Gao & Coley, 2020). One line to improve synthesizability is multi-objective optimization using cheap functions to estimate the synthesizability (Ertl & Schuffenhauer, 2009), but this is too simplified to reflect complex synthetic principles (Cretu et al., 2024). Other efforts aim to project molecules from generative models into a synthesizable space (Gao et al., 2022b; Luo et al., 2024; Gao et al., 2024b), but chemical modifications in this process can degrade the optimized properties.

To address this issue, recent works formulate the generation of synthetic pathways as a Markov decision process (MDP) for molecular design (Gottipati et al., 2020). These approaches return a synthesizable molecule by assembling purchasable building blocks and reaction templates according to the generated synthetic pathway. Notably, the emergence of virtual libraries—created by combinatorially enumerating building blocks and reaction templates, such as Enamine REAL (Grygorenko et al., 2020)—allows the generated molecules to be readily synthesizable on demand with their synthetic pathways. Recent studies (Cretu et al., 2024; Koziarski et al., 2024) advanced this



Figure 1: **Overview of RXNFLOW.** (a) Synthetic action space which is represented in a continuous action space. Each colored box corresponds to a reaction template and the molecules in the box are reactant blocks. (b) Policy estimation using the action space subsampling in a manner of importance sampling. (c) Molecular generation process and model training.

approach by training the decision-making policy using the objective of generative flow networks (GFlowNets; Bengio et al., 2021). This objective encourages the policy to sample in proportion to the reward function, enabling the retrieval of samples from a diverse range of modes.

Unlike atom-based or fragment-based models, the synthetic action spaces are massive and composed 075 of millions of building blocks and tens of reaction templates. While the large action spaces offer 076 opportunities to discover novel hit candidates by expanding an explorable chemical space (Sady-077 bekov et al., 2022), it incurs significant computational overhead. Thus, prior works have restricted the action spaces to trade off the size of the search space for efficiency. However, reducing the 079 search space leads to a decrease in diversity and synthetic complexity. While one could compensate for the reduced number of building block candidates by adding more reaction steps, this leads to 081 the increase in synthetic complexity, negatively influencing synthesizability, yield, and cost (Coley 082 et al., 2018; Kim et al., 2023). 083

In response to this challenge, we propose RXNFLOW, a synthesis-oriented generative framework 084 that allows training generative flows over a large action space to generate synthetic pathways for 085 drug design. The distinctive features of this method are as follows. First, we introduce an action space subsampling (Figure 1) to handle massive action spaces without significant memory overhead, enabling us to explore a broader chemical space with fewer reaction steps than existing models. 088 Then, we train the generative policy with a GFlowNet objective to sample both diverse and potent molecules from the expanded search space. We demonstrate that RXNFLOW effectively generates 090 drug candidates, outperforming existing reaction-based, atom-based, and fragment-based baselines across various SBDD tasks, while ensuring the synthesizability of generated drug candidates. We 091 also achieve a new state-of-the-art Vina score, drug-likeness, and synthesizability on the Cross-092 Docked2020 pocket-conditional generation benchmark (Luo et al., 2021). 093

Furthermore, we formulate an adaptable MDP (Figure 2) for consistent flow estimation on modified building block libraries, which can be highly practical in real-world applications. By combining the proposed MDP with action embedding (Dulac-Arnold et al., 2015), which represents actions in a continuous space instead of a discrete space, RXNFLOW can achieve further objectives or incorporate newly discovered building blocks without retraining. We experimentally show that RXNFLOW can achieve an additional solubility objective and behave appropriately for unseen building blocks. This capability makes RXNFLOW highly adaptable to real-world drug discovery pipeline, where new objectives frequently arise (Fink et al., 2022) and building block libraries are continuously expanding (Grygorenko et al., 2020).

103

054

058

060

061

062

063

064 065 066

067

068

069

2 RELATED WORKS

104 105

> Structure-based drug discovery. The first type of SBDD involves pocket-specific optimization
> methods to enhance docking scores against a single pocket, including evolutionary algorithms (Reidenbach, 2024), reinforcement learning (RL) (Zhavoronkov et al., 2019), and GFlowNets (Bengio

et al., 2021). However, these require individual optimizations for each pocket, limiting scalability.
The second type is based on **pocket-conditional generation**, which generates molecules against arbitrary given pockets without additional training. This can be achieved by distribution-based generative models (Ragoza et al., 2022; Peng et al., 2022; Guan et al., 2023b; Schneuing et al., 2023; Qu et al., 2024) trained on protein-ligand complex datasets to model ligand distributions for given pockets. On the other hand, Shen et al. (2023) formulated a pocket-conditioned policy for GFlowNets to generate samples from reward-biased distributions in a *zero-shot* manner.

115 Syntheis-oriented generative models. To ensure the synthesizability of generated molecules, 116 synthesis-oriented de novo design approaches incorporate combinatorial chemistry principles into 117 generative models. Bradshaw et al. (2019) represented the synthetic pathways as directed acyclic 118 graphs (DAG) for generative modeling. Horwood & Noutahi (2020) formulated the synthetic pathway generation as an MDP and optimize the molecules with RL. Similarly, Gao et al. (2022b) 119 employed a genetic algorithm to optimize synthesis trees to generate molecules with the desired 120 properties. Seo et al. (2023) proposed a conditional generative model to directly sample molecules 121 with desired properties without optimization. Recently, Cretu et al. (2024); Koziarski et al. (2024) 122 proposed the reaction-based GFlowNet to generate diverse and potent molecules. 123

Action embedding for large action spaces. To handle large action spaces in the synthesis-oriented generation, Gottipati et al. (2020) employed action embedding (Dulac-Arnold et al., 2015) that represents building blocks in a continuous action space with their chemical information. Later, Seo et al. (2023); Koziarski et al. (2024) experimentally demonstrated that it can enhance the model training and generative performance. The continuous action space provides the benefit of reducing the computational complexity for sampling from large space of actions and the memory requirement for parameterizing the categorical distribution over the large action space.

131 Generative Flow Networks. GFlowNets are a learning framework for a stochastic generative policy that constructs an object through a series of decisions, where the probability of generating each 132 object is proportional to a given reward associated with that object (Bengio et al., 2021). Unlike other 133 optimization methods that maximize rewards and often converge to a single solution, GFlowNets 134 aim to sample a diverse set of high-rewarded modes, which is vital for novel drug design (Shen et al., 135 2023; Jain et al., 2022). To this end, the generative policy is trained using objectives such as flow 136 matching (Bengio et al., 2021), detailed balance (Bengio et al., 2023), and trajectory balance (Malkin 137 et al., 2022). Extending the GFlowNets to various applications is an active area of research, e.g., 138 GFlowNets have been applied to designing crystal structures (Nguyen et al., 2023b), phylogenetic 139 inference (Zhou et al., 2023), finetuning diffusion models (Venkatraman et al., 2024), and causal 140 inference (Nguyen et al., 2023a).

141 142 143

144

145

155 156 157

159 160 161

3 Method

3.1 GFLOWNET PRELIMINARIES

146 GFlowNets (Bengio et al., 2021) are the class of generative models that learn to sample objects 147 $x \in \mathcal{X}$ proportional to a given reward function, i.e., $p(x) \propto R(x)$. This is achieved by sequentially 148 constructing a compositional object x through a series of state transitions $s \to s'$, forming a trajectory $\tau = (s_0 \to \ldots \to s_n = x) \in \mathcal{T}$. The set of all complete trajectories from the initial state 149 s_0 can be represented as a directed acyclic graph $\mathcal{G} = (\mathcal{S}, \mathcal{A})$ with a reachable state space \mathcal{S} and 150 an action space A. Each action a induces a transition from the state s to the state s', expressed as 151 s' = T(s, a) and represented as $s \to s'$. Then, we define the *trajectory flow* $F(\tau)$, which flows 152 along the trajectory $\tau = (s_0 \to \ldots \to s_n = x)$, as the reward of the terminal state, R(x). The edge 153 flow $F(s \to s')$, or equivalently F(s, a), is defined as the total flow along the edge $a: s \to s'$: 154

$$F(s \to s') = F(s, a) = \sum_{\tau \in \mathcal{T} \text{ s.t. } (s \to s') \in \tau} F(\tau).$$
(1)

158 The state flow F(s) for the intermediate state is defined as the total flow through the state s:

$$F(s) = \sum_{\tau \in \mathcal{T} \text{ s.t. } s \in \tau} F(\tau) = \sum_{(s'' \to s) \in \mathcal{A}} F(s'' \to s) = \sum_{(s \to s') \in \mathcal{A}} F(s \to s').$$
(2)

Intermediate flow matching condition

169

175

185 186 187

188 189

190 191

205 206

207 208

In addition to the flow matching condition for intermediate states, there are two boundary conditions for the states. First, the flow of a terminal state x must equal the reward of the objective: F(x) = R(x). Second, the *partition function*, Z, is equivalent to the sum of all trajectory flows and the sum of all rewards: $Z = F(s_0) = \sum_{\tau \in \mathcal{T}} F(\tau) = \sum_{x \in \mathcal{X}} R(x)$. These three conditions—one for intermediate states and two for boundary states—are known as the *flow matching* conditions and ensure that GFlowNets generate objectives proportional to their rewards.

To convert the flow network into a usable policy, we define the *forward policy* as the forward transition probability $P_F(s'|s)$ and the *backward policy* as the backward transition probability $P_B(s|s')$:

$$P_F(s'|s) := P(s \to s'|s) = \frac{F(s \to s')}{F(s)}, \quad P_B(s|s') := P(s \to s'|s') = \frac{F(s \to s')}{F(s')}.$$
 (3)

3.2 ACTION SPACE FOR SYNTHETIC PATHWAY GENERATION

Following Cretu et al. (2024), we treated a chemical reaction as a forward transition and a synthetic 176 pathway as a trajectory for molecular generation. For the initial state s_0 , the model always chooses 177 AddFirstReactant to sample a building block b from the entire building block set \mathcal{B} as a starting 178 molecule. For the later states s, the model samples actions among ReactUni, ReactBi, or Stop. 179 When the action type is ReactUni, the model performs in silico uni-molecular reactions with an 180 assigned reaction template $r \in \mathcal{R}_1$. When the action type is ReactBi, the model performs bi-181 molecular reactions with a reaction template $r \in \mathcal{R}_2$ and a reactant block b in the possible reactant 182 set for the reaction template $r: \mathcal{B}_r \subseteq \mathcal{B}$. If Stop is sampled, the trajectory is terminated. To sum 183 up, the allowable action space $\mathcal{A}(s)$ for the state s is: 184

$$\mathcal{A}(s) = \begin{cases} \mathcal{B} & \text{if } s = s_0 \\ \{\text{Stop}\} \cup \mathcal{R}_1 \cup \{(r, b) | r \in \mathcal{R}_2, b \in \mathcal{B}_r\} & \text{otherwise} \end{cases}$$
(4)

where unavailable reaction templates to the molecule of the state *s* are masked.

3.3 FLOW NETWORK ON ACTION SPACE SUBSAMPLING

We propose a novel memory-efficient technique called the *action space subsampling*, that estimates the state flow $F_{\theta}(s)$ from a subset of the outgoing edge flows $F_{\theta}(s \to s')$ for forward policy estimation. First, we implement an auxiliary policy, termed *subsampling policy* $\mathcal{P}(\mathcal{A})$, which samples a subset of the action space $\mathcal{A}^* \subseteq \mathcal{A}$. This reduces both the memory footprint and the computational complexity from $\mathcal{O}(|\mathcal{B}||\mathcal{R}_2|)$ to $\mathcal{O}(|\mathcal{B}^*||\mathcal{R}_2|)$ with the controllable size $|\mathcal{B}^*|$. We then estimate the forward policy by importance sampling. In contrast to the parameterized forward policy, we formulate a fixed backward policy since it is hard to force invariance to molecule isomorphism (Malkin et al., 2022). Theoretical backgrounds are provided in Sec. A.

199 200 201 201 201 202 203 204 Subsampling policy. Subsampling policy $\mathcal{P}(\mathcal{A})$ performs uniform sampling for the initial state and importance sampling for the later states. For the initial state, the allowable action space $\mathcal{A}(s_0) = \mathcal{B}$ is homogeneous since all of them are AddFirstReactant actions. For the later state, the action space is comprised of one Stop, tens of ReactUni actions, and millions of ReactBi actions. To capture rare-type actions in the inhomogeneous space, we use all Stop and ReactUni actions. The partial action space $\mathcal{A}^*(s) \sim \mathcal{P}(\mathcal{A}(s))$ comprises the uniform subset $\mathcal{B}^* \subseteq \mathcal{B}$ or $\mathcal{B}_r^* \subseteq \mathcal{B}_r$:

$$\mathcal{A}^*(s) = \begin{cases} \mathcal{B}^* & \text{if } s = s_0\\ \{\text{Stop}\} \cup \mathcal{R}_1 \cup \{(r,b) | r \in \mathcal{R}_2, b \in \mathcal{B}_r^*\} & \text{otherwise} \end{cases}$$
(5)

Forward policy. To estimate the forward policy $P_F(s'|s) = F(s \to s')/F(s)$ from the partial action space \mathcal{A}^* , we estimate the state flow F(s) with a subset of outgoing edge flows $F(s \to s')$. Since we introduce importance sampling for action types, we weight the edge flow $F(s \to s')$ of edge $a: s \to s'$ according to the subsampling ratio:

213 214 $(|\mathcal{B}|/|\mathcal{B}^*| \quad \text{if } a \in \mathcal{B}$

214
215
$$w_a = w_{(s \to s')} = \begin{cases} |\mathcal{B}_r|/|\mathcal{B}_r^*| & \text{if } a \text{ is } (r, b) \text{ where } r \in \mathcal{R}_2 \\ 1 & \text{if } a \text{ is } \text{Stop or } a \in \mathcal{R}_1 \end{cases}$$
(6)



Figure 2: Comparison of using modified building block library for the generation: (a) a hierarchical MDP, and (b) a non-hierarchical MDP. More details are in Figure 8.

By weighting edge flows, we can estimate the state flow $\hat{F}_{\theta}(s; \mathcal{A}^*)$ as:

$$\hat{F}_{\theta}(s; \mathcal{A}^*) = \sum_{(s \to s') \in \mathcal{A}^*(s)} w_{(s \to s')} F_{\theta}(s \to s'), \tag{7}$$

which the estimated forward policy is $\hat{P}_F(s'|s; \mathcal{A}^*; \theta) = F_{\theta}(s \to s') / \hat{F}_{\theta}(s; \mathcal{A}^*)$.

Action embedding. In the standard implementation (Bengio et al., 2021) of the flow function F_{θ} and its neural network ϕ_{θ} , the edge flow of the edge $a: s \to s'$ is computed with the corresponding action-specific parameter $\theta_a: F_{\theta}(s \to s') = F_{\theta}(s, a) = \phi_{\theta_a}^{\text{flow}}(\phi_{\theta}^{\text{state}}(s)).$

However, large action spaces require numerous parameters which increase model complexity. To address this, we use an additional network $\phi_{\theta}^{\text{block}}$ for AddFirstReactant and ReactBi, which embeds the building block *b* into a continuous action space with its structural information, molecular fingerprints (see Sec. B.3):

$$F_{\theta}(s_0, b) = \phi_{\theta}^{\text{flow}}(\phi_{\theta}^{\text{state}}(s), \phi_{\theta}^{\text{block}}(b)), \quad F_{\theta}(s, (r, b)) = \phi_{\theta}^{\text{flow}}(\phi_{\theta}^{\text{state}}(s), \delta(r), \phi_{\theta}^{\text{block}}(b))$$
(8)

where $\delta(r)$ is the one-hot encoding for a bi-molecular reaction template r.

GFlowNet training. In this work, we use the trajectory balance (TB; Malkin et al., 2022) as the training objective of GFlowNets from Eq. (9) and train models following Sec. B.4. The action space subsampling is performed for each transition $s_t \rightarrow s_{t+1}$: $\mathcal{A}_t^* \sim \mathcal{P}(\mathcal{A})$.

$$\hat{\mathcal{L}}_{\text{TB}}(\tau) = \left(\log \frac{Z_{\theta} \prod_{t=1}^{n} \hat{P}_{F}(s_{t}|s_{t-1};\mathcal{A}_{t-1}^{*};\theta)}{R(x) \prod_{t=1}^{n} P_{B}(s_{t-1}|s_{t})}\right)^{2}$$
(9)

For online training, we use the sampling policy π_{θ} proportional to $\hat{P}_F(-|-; \mathcal{A}^*; \theta)$, given by:

$$\pi_{\theta}(s'|s;\mathcal{A}^*) = \frac{w_{(s\to s')}F_{\theta}(s\to s')}{\sum_{(s\to s'')\in\mathcal{A}^*(s)}w_{(s\to s'')}F_{\theta}(s\to s'')}$$
(10)

3.4 JOINT SELECTION OF TEMPLATES AND BLOCKS

For bi-molecular reactions, existing synthesis-oriented methods (Gao et al., 2022b; Cretu et al., 2024; Koziarski et al., 2024) formulated a hierarchical MDP which selects a reaction template r first and then the corresponding reactant block b sequentially from Eq. (11). However, as shown in Figure 2, the probability of selecting each reaction template r is fixed after training in a hierarchical MDP, and this rigidity can lead to incorrect policy estimates in modified block libraries. Therefore, we formulate a non-hierarchical MDP that jointly selects reaction templates and reactant blocks (r, b) at once, as given by Eq. (12), resulting in more consistent estimates of forward policy P_F .

$$P_F(T(s,(r,b))|s;\theta) = \frac{F_\theta(s,r)}{\sum_{r'\in\mathcal{R}_1\cup\mathcal{R}_2\cup\{\text{Stop}\}}F_\theta(s,r')} \times \frac{F_\theta(s,(r,b))}{\sum_{b'\in\mathcal{B}_r}F_\theta(s,(r,b'))}$$
(11)

 P_{I}

$$F(T(s,(r,b))|s;\theta) = \frac{F_{\theta}(s,(r,b))}{\sum_{r'\in\mathcal{R}_1\cup\{\text{Stop}\}}F_{\theta}(s,r') + \sum_{r'\in\mathcal{R}_2}\sum_{b'\in\mathcal{B}_{r'}}F_{\theta}(s,(r',b'))}$$
(12)

270 4 EXPERIMENTS

271 272

Overview. We validate the effectiveness of RXNFLOW in two common SBDD tasks: pocket-specific optimization (Sec. 4.1) and pocket-conditional generation (Sec. 4.2). To the best of our knowledge, this is the first synthesis-oriented approach for pocket-conditional generation. We also investigate the applicability of RXNFLOW in real-world drug discovery pipelines where new further objectives may be introduced (Sec. 4.3) and the building block libraries are constantly expanded (Sec. 4.4). Lastly, we conduct an ablation study in Sec. 4.5 and a theoretical analysis in Sec. D.8.

Setup. We use the reaction template set constructed by Cretu et al. (2024) including 13 uni- and 58 bi-molecular reaction templates. For the building blocks, we use 1.2M blocks from the Enamine comprehensive catalog. We use up to 3 reaction steps for generation following Enamine REAL Space (Grygorenko et al., 2020), while SynFlowNet and RGFN allow 4 steps. For the subsampling policy, we set a sampling ratio of 1%. The experimental details are provided in Sec. C.

Synthesizability estimation. To assess the synthesizability of the generated compounds, we used the computationally intensive retrosynthetic analysis tool AiZynthFinder (Genheden et al., 2020) with the Enamine building block library. We note that the molecule is identified as synthesizable only if it can be synthesized using the USPTO reactions (Lowe, 2017) and given building blocks.

287

4.1 POCKET-SPECIFIC OPTIMIZATION WITH GPU-ACCELERATED DOCKING

Setup. Since GFlowNets sample a large number of molecules for online training, we employed a GPU-accelerated UniDock (Yu et al., 2023) with Vina scoring (Trott & Olson, 2010). It is well known that docking can be hacked by increasing molecule size (Pan et al., 2003), so the appropriate constraints are required. We select QED (Bickerton et al., 2012) as a comprehensive molecular property constraint, QED>0.5, and set the reward function as $R(x) = w_1 \text{QED}(x) + w_2 \widehat{\text{Vina}}(x)$ where w_1, w_2 are used as the input of multi-objective GFlowNets (Jain et al., 2023) for all GFlowNets and are set to 0.5 for non-GFlowNet baselines. Vina is a normalized docking score (Eq. (32)).

297 Each method generates up to 64,000 molecules for each of the 15 proteins in the LIT-PCBA dataset 298 (Tran-Nguyen et al., 2020). We then filter the molecules with the property constraint and select 299 the top 100 diverse candidates based on the docking score, using a Tanimoto distance threshold of 300 0.5 to ensure structural diversity. The selected molecules are evaluated with the following metrics: 301 Hit ratio (%) measures the fraction of *hits*, defined as the molecules that are identified as synthesizable by AiZynthFinder and having better docking scores than known active ligands (Lee et al., 302 2023). Vina (kcal/mol) measures the average docking score. Synthesizability (%) is the fraction 303 of synthesizable molecules. Synthetic complexity, which is highly correlated to yield and cost, is 304 evaluated as the average number of synthesis steps (Coley et al., 2018). 305

306 Baselines. We perform comparisons to various synthetic-oriented approaches: genetic algorithm 307 (SynNet) (Gao et al., 2022b), conditional generative model (BBAR¹) (Seo et al., 2023), and GFlowNets (SynFlowNet, RGFN) (Cretu et al., 2024; Koziarski et al., 2024). For SynFlowNet 308 and RGFN, we used 6,000 and 350 blocks, respectively, and set the maximum reaction step to 4 309 following the original papers. Moreover, we consider fragment-based GFlowNets (FragGFN) to 310 analyze the effects of synthetic constraints on the performance. For FragGFN, we also consider 311 additional synthesizability objectives with commonly-used synthetic accessibility score (SA; Ertl & 312 Schuffenhauer, 2009) (FragGFN+SA). 313

Results. The results for the first five targets are shown in Tables 1 and 2, and additional results for 314 the 10 remaining targets are reported in Sec. D.1. The property distribution for each target are re-315 ported in Sec. D.2. RXNFLOW outperforms the baselines across all test proteins, demonstrating that 316 the expanded sample space with the large action space enabled the model to generate more potent 317 and diverse molecules. Additionally, as shown in Tables 3 and 4, RXNFLOW ensures the synthesiz-318 ability of the generated molecules more effectively than the other synthesis-oriented methods and 319 GFlowNets which employ the same reactions as ours. These results support our primary assertion 320 that existing synthesis-oriented approaches using more synthetic steps with smaller building block 321 libraries can increase overall synthesis complexity and reduce synthesizability. Furthermore, Frag-

¹Since BBAR requires labeled training data with QED and docking score, we perform docking with random 62,720 ZINC molecules for training and evaluate 1,280 samples according to the reported splitting ratio.

Hit Ratio (%, ↑)						
Category	Method	ADRB2	ALDH1	ESR_ago	ESR_antago	FEN1
Fragment	FragGFN	4.00 (± 3.54)	$3.75 (\pm 1.92)$	$0.25 (\pm 0.43)$	$0.25~(\pm 0.43)$	$0.25 (\pm 0.25)$
Tragment	FragGFN+SA	5.75 (± 1.48)	4.00 (± 1.58)	0.25 (± 0.43)	$0.00 (\pm 0.00)$	$0.00 (\pm 0.00)$
	SynNet	$45.83 \ (\pm \ 7.22)$	$25.00 \ (\pm 25.00)$	$0.00 (\pm 0.00)$	$0.00 \ (\pm \ 0.00)$	$50.00 (\pm$
	BBAR	$21.25 (\pm 5.36)$	$18.25 (\pm 1.92)$	$3.50(\pm 1.12)$	$2.25 (\pm 1.09)$	11.75 (±
Reaction	SynFlowNet	52.75 (± 1.09)	$57.00 (\pm 6.04)$	$30.75 (\pm 10.03)$	$11.25 (\pm 1.48)$	53.00 (±
	RGFN	$46.75 (\pm 6.87)$	$39.75 (\pm 8.17)$	$4.50 (\pm 1.66)$	$1.25 (\pm 0.43)$	19.75 (±
	RXNFLOW	60.25 (± 3.77)	63.25 (± 3.11)	71.25 (± 4.15)	46.00 (± 7.00)	65.50 (±
Tabl	e 2: Vina. Mea	an and standard	d deviation ove	r 4 runs. The b	est results are i	n bold.
Average Vina Docking Score (kcal/mol, \downarrow)						
Category	Method	ADRB2	ALDH1	ESR_ago	ESR_antago	FEN
D	FragGFN	-10.19 (± 0.33)	-10.43 (± 0.29)	-9.81 (± 0.09)	-9.85 (± 0.13)	-7.67 (±
Fragment	FragGFN+SA	$-9.70 (\pm 0.61)$	$-9.83 (\pm 0.65)$	$-9.27 (\pm 0.95)$	$-10.06 (\pm 0.30)$	-7.26 (±
	SynNet	-8.03 (± 0.26)	-8.81 (± 0.21)	-8.88 (± 0.13)	-8.52 (± 0.16)	-6.36 (±
Reaction	BBAR	$-9.95~(\pm 0.04)$	$-10.06 (\pm 0.14)$	$-9.97 (\pm 0.03)$	$-9.92 (\pm 0.05)$	-6.84 (±
	SynFlowNet	$-10.85~(\pm 0.10)$	$-10.69 \ (\pm 0.09)$	$-10.44 \ (\pm 0.05)$	$-10.27~(\pm 0.04)$	-7.47 (±
	RGFN	$-9.84 (\pm 0.21)$	$-9.93 (\pm 0.11)$	-9.99 (± 0.11)	$-9.72 (\pm 0.14)$	-6.92 (±
	RXNFLOW	-11.45 (± 0.05)	-11.26 (± 0.07)	-11.15 (± 0.02)	-10.77 (± 0.04)	-7.66 (±
Table 3: S	Synthesizability	y. Mean and st	andard deviation	on over 4 runs.	The best result	s are in b
			Percentage of S	Synthesizable M	olecules $(\%, \uparrow)$	
					EGD	
Category	Method	ADRB2	ALDH1	ESR_ago	ESR_antago	FEN
Category	Method FragGFN	ADRB2 4.00 (± 3.54)	ALDH1 3.75 (± 1.92)	ESR_ago 1.00 (± 1.00)	$\frac{\text{ESR}_\text{antago}}{3.75 (\pm 1.92)}$	0.25 (±
Category Fragment	Method FragGFN FragGFN+SA	ADRB2 4.00 (± 3.54) 5.75 (± 1.48)	ALDH1 3.75 (± 1.92) 6.00 (± 2.55)	ESR_ago 1.00 (± 1.00) 4.00 (± 2.24)	3.75 (± 1.92) 1.00 (± 0.00)	FEN 0.25 (± 0.00 (±
Category Fragment	Method FragGFN FragGFN+SA SynNet	ADRB2 4.00 (± 3.54) 5.75 (± 1.48) 54.17 (± 7.22)	ALDH1 3.75 (± 1.92) 6.00 (± 2.55) 50.00 (± 0.00)	$\frac{\text{ESR_ago}}{1.00 (\pm 1.00)}$ $\frac{4.00 (\pm 2.24)}{50.00 (\pm 0.00)}$	$\begin{array}{r} \text{ESR_antago} \\ 3.75 (\pm 1.92) \\ 1.00 (\pm 0.00) \\ 25.00 (\pm 25.00) \end{array}$	FEN 0.25 (± 0.00 (± 50.00 (±
Category Fragment	Method FragGFN FragGFN+SA SynNet BBAR	ADRB2 4.00 (± 3.54) 5.75 (± 1.48) 54.17 (± 7.22) 21.25 (± 5.36)	$\begin{array}{c} \text{ALDH1} \\ \hline 3.75 \ (\pm 1.92) \\ 6.00 \ (\pm 2.55) \\ \hline 50.00 \ (\pm 0.00) \\ 19.50 \ (\pm 3.20) \\ \end{array}$	$\frac{\text{ESR_ago}}{1.00 (\pm 1.00)}$ $\frac{4.00 (\pm 2.24)}{50.00 (\pm 0.00)}$ $17.50 (\pm 1.50)$	ESR_antago $3.75 (\pm 1.92)$ $1.00 (\pm 0.00)$ $25.00 (\pm 25.00)$ $19.50 (\pm 3.64)$	FEN 0.25 (± 0.00 (± 50.00 (± 20.00 (±
Category Fragment Reaction	Method FragGFN FragGFN+SA SynNet BBAR SynFlowNet	$\begin{array}{c} \text{ADRB2} \\ \hline 4.00 \ (\pm \ 3.54) \\ 5.75 \ (\pm \ 1.48) \\ \hline 54.17 \ (\pm \ 7.22) \\ 21.25 \ (\pm \ 5.36) \\ 52.75 \ (\pm \ 1.09) \\ \end{array}$	$\begin{array}{c} \text{ALDH1} \\ \hline 3.75 (\pm 1.92) \\ 6.00 (\pm 2.55) \\ \hline 50.00 (\pm 0.00) \\ 19.50 (\pm 3.20) \\ 57.00 (\pm 6.04) \\ \end{array}$	$\frac{\text{ESR_ago}}{1.00 (\pm 1.00)}$ $\frac{1.00 (\pm 2.24)}{50.00 (\pm 0.00)}$ $\frac{17.50 (\pm 1.50)}{53.75 (\pm 9.52)}$	$\begin{array}{r} \text{ESR_antago} \\ \hline 3.75 (\pm 1.92) \\ 1.00 (\pm 0.00) \\ \hline 25.00 (\pm 25.00) \\ 19.50 (\pm 3.64) \\ 56.50 (\pm 2.29) \end{array}$	FEN 0.25 (± 0.00 (± 50.00 (± 20.00 (± 53.00 (±
Category Fragment Reaction	Method FragGFN FragGFN+SA SynNet BBAR SynFlowNet RGFN	$\begin{array}{c} \text{ADRB2} \\ \hline 4.00 \ (\pm \ 3.54) \\ 5.75 \ (\pm \ 1.48) \\ \hline 54.17 \ (\pm \ 7.22) \\ 21.25 \ (\pm \ 5.36) \\ 52.75 \ (\pm \ 1.09) \\ 46.75 \ (\pm \ 6.86) \\ \end{array}$	$\begin{array}{c} \text{ALDH1} \\ \hline 3.75 (\pm 1.92) \\ 6.00 (\pm 2.55) \\ \hline 50.00 (\pm 0.00) \\ 19.50 (\pm 3.20) \\ 57.00 (\pm 6.04) \\ 47.50 (\pm 4.06) \\ \end{array}$	$\begin{array}{c} \text{ESR_ago} \\ \hline 1.00 \ (\pm \ 1.00) \\ 4.00 \ (\pm \ 2.24) \\ \hline 50.00 \ (\pm \ 0.00) \\ 17.50 \ (\pm \ 1.50) \\ 53.75 \ (\pm \ 9.52) \\ 50.25 \ (\pm \ 2.17) \\ \end{array}$	$\begin{array}{r} \text{ESR_antago} \\ \hline 3.75 (\pm 1.92) \\ 1.00 (\pm 0.00) \\ \hline 25.00 (\pm 25.00) \\ 19.50 (\pm 3.64) \\ 56.50 (\pm 2.29) \\ 49.25 (\pm 4.38) \end{array}$	FEN 0.25 (± 0.00 (± 50.00 (± 20.00 (± 53.00 (± 48.50 (±

0010.						
		Average Number of Synthesis Steps (\downarrow)				
Categor	y Method	ADRB2	ALDH1	ESR_ago	ESR_antago	FEN1
Fragme	nt FragGFN FragGFN+SA	$\begin{array}{c} 3.60 \ (\pm \ 0.10) \\ 3.73 \ (\pm \ 0.09) \end{array}$	$\begin{array}{c} 3.83 \ (\pm \ 0.08) \\ 3.66 \ (\pm \ 0.04) \end{array}$	$\begin{array}{c} 3.76 \ (\pm \ 0.20) \\ 3.66 \ (\pm \ 0.07) \end{array}$	$\begin{array}{c} 3.76 \ (\pm \ 0.16) \\ 3.67 \ (\pm \ 0.21) \end{array}$	$\begin{array}{c} 3.34 \ (\pm \ 0.18) \\ 3.79 \ (\pm \ 0.19) \end{array}$
Reactio	SynNet BBAR n SynFlowNet RGFN	$\begin{array}{c} 3.29 \ (\pm \ 0.36) \\ 3.60 \ (\pm \ 0.17) \\ 2.64 \ (\pm \ 0.07) \\ 2.88 \ (\pm \ 0.21) \end{array}$	$\begin{array}{c} 3.50 \ (\pm \ 0.00) \\ 3.62 \ (\pm \ 0.19) \\ 2.48 \ (\pm \ 0.07) \\ 2.65 \ (\pm \ 0.09) \end{array}$	$\begin{array}{c} 3.00 \ (\pm \ 0.00) \\ 3.76 \ (\pm \ 0.04) \\ 2.60 \ (\pm \ 0.25) \\ 2.78 \ (\pm \ 0.19) \end{array}$	$\begin{array}{c} 4.13 \ (\pm \ 0.89) \\ 3.72 \ (\pm \ 0.11) \\ 2.45 \ (\pm \ 0.09) \\ 2.91 \ (\pm \ 0.23) \end{array}$	$\begin{array}{c} 3.50 \ (\pm \ 0.00) \\ 3.59 \ (\pm \ 0.14) \\ 2.56 \ (\pm \ 0.29) \\ 2.76 \ (\pm \ 0.17) \end{array}$
	RXNFLOW	$2.42 (\pm 0.23)$	$2.19 (\pm 0.12)$	$1.95 (\pm 0.20)$	$2.15 (\pm 0.18)$	$2.23 (\pm 0.18)$

2.19 (± 0.12)

 $1.95 (\pm 0.20)$

 $2.15 (\pm 0.18)$

 $2.23 (\pm 0.18)$

Table 1:	Hit ratio.	Mean and	standard	deviation	over 4 runs.	The best re	esults are in b	oold.

GFN+SA does not show meaningful improvement in synthesizability, implying that optimization of a cheap synthesizability estimation function is suboptimal.

 $2.42 (\pm 0.23)$

Furthermore, it is noteworthy that RXNFLOW outperformed FragGFN which does not consider syn-thesizability. This improvement can be attributed to two key factors. First, the Enamine building block library is specifically curated for drug discovery, limiting the search space \mathcal{X} to a drug-like chemical space and simplifying the optimization complexity. Second, RXNFLOW needs shorter tra-jectories compared to FragGFNs since it assembles molecules with large building blocks instead of atoms or small fragments. This is beneficial for trajectory balance objectives where stochastic gradient variance tends to increase over longer trajectories (Malkin et al., 2022).

Table 5: CrossDocked2020 benchmark. We report the average and median values over the average properties for each test pocket. The best results are in **bold** and the second ones are in underlined. We denote the Generation Success as Succ., Synthesizability as Synthesiz, and Diversity as Div. Reference means the known binding active molecules in the test set. MolCRAFT-large (Qu et al., 2024) is the result when generating with more atoms than the reference ligands.

$ \begin{array}{c c} na (\downarrow) \\ \hline Med. \\ \hline -7.80 \\ \hline -7.16 \\ 7.56 \\ \hline -7.56 \\ \hline \end{array} $	QE Avg. 0.48 0.57 0.49	D (†) Med. 0.47 0.58 0.49	Synthe Avg. 36.1% <u>29.1%</u> 9.9%	$\frac{\text{siz. }(\uparrow)}{\text{Med.}}$ $-$ $\frac{22.0\%}{3.2\%}$	Div. (↑) Avg. - 0.83 0.87	$\frac{\text{Time } (\downarrow)}{\text{Avg.}}$
Med. -7.80 -7.16 -7.56	Avg. 0.48 0.57 0.49	Med. 0.47 0.58 0.49	Avg. 36.1% <u>29.1%</u> 9.9%	Med. - <u>22.0%</u> 3.2%	Avg.	Avg.
-7.80 -7.16 -7.56	0.48 0.57 0.49	0.47 0.58 0.49	36.1% <u>29.1%</u> 9.9%	- <u>22.0%</u> 3.2%	- 0.83	- 2504
) -7.16 -7.56	0.57 0.49	0.58 0.49	$\frac{29.1\%}{9.9\%}$	$\frac{22.0\%}{3.2\%}$	0.83	2504
-7.56	0.49	0.49	9.9%	3.2%	0.87	2420
				0.2 /0	0.07	3428
-8.25	0.37	0.35	0.9%	0.0%	0.84	6189
5 -7.10	0.47	0.48	2.9%	2.0%	0.88	135
-8.05	0.50	0.50	16.5%	9.1%	0.84	141
5 -9.24	0.45	0.44	3.9%	0.0%	0.82	>141
	0.67	0.67	1.3%	1.0%	0.67	4
-8.44			24.00	24 5 07	0.81	4
• /	.4 -0.44				25 0.03 0.67 0.67 34.8% 34.5%	35 -9.03 0.67 0.67 34.8% 34.5% 0.81



Figure 3: Visualization of generated molecules in a zero-shot manner. (a-b) Docking results of generated molecules and known reference ligands of TBK1 (PDB Id: 1FV, SU6). (c) Generative trajectory, which is the generated synthetic pathway of the left molecule in (a).

4.2 ZERO-SHOT SAMPLING VIA POCKET CONDITIONING

Setup. We extend our works to a pocket-conditional generation problem to design binders for arbi-trary pockets without additional training oracles (Ragoza et al., 2022; Liu et al., 2022; Peng et al., 2022; Guan et al., 2023a;b; Schneuing et al., 2023). To address this challenge, we follow TacoGFN (Shen et al., 2023), which is a fragment-based GFlowNet for pocket-conditional generation. Since it requires more training oracles to learn pocket-conditional policies than target-specific generation, TacoGFN used pre-trained proxies that predict docking scores against arbitrary pockets using pharmacophore representation (Seo & Kim, 2023). Since RXNFLOW explicitly considers synthesiz-ability, we exclude the SA score from the TacoGFN's reward function as described in Sec. C.2.

We generate 100 molecules for each of the 100 test pockets in the CrossDocked2020 benchmark (Francoeur et al., 2020) and evaluate them with the following metrics: Vina (kcal/mol) measures the average docking score from QuickVina (Alhossary et al., 2015). QED measures the average drug-likeness of molecules. Synthesizability (%) is the fraction of synthesizable molecules. Diversity measures an average pairwise Tanimoto distance of ECFP4 fingerprints. Moreover, we report the Generation Success (%) which is the percentage of unique RDKit-readable molecules without disconnected parts, and Time (sec.) which is the average sampling time to generate 100 molecules.

Baselines. We compare RXNFLOW with state-of-the-art distribution learning-based models trained on a synthesizable drug set: an autoregressive model (**Pocket2Mol** (Peng et al., 2022)), diffusion models (TargetDiff (Guan et al., 2023a), DiffSBDD (Schneuing et al., 2023), DecompDiff (Guan et al., 2023b)), and bayesian flow network (MolCRAFT (Qu et al., 2024)). We also perform a comparison with an optimization-based TacoGFN (Shen et al., 2023). For a fair comparison with the distribution learning-based approaches, we used the docking proxy trained on CrossDocked2020.

436

437

438

439

440

445 446

447

450

451



Figure 4: Property distribution of sampled 441 molecules with "all" building blocks and "low"-442 TPSA building blocks. Vina score was calculated 443 against the KRAS-G12C target. 444



QED reward distribution Figure 5: of generated molecules for each of the "seen", "unseen", and "all" blocks. Additional results are in Figure 14.

Result. As shown in Table 5, RXNFLOW achieves significant improvements in drug-related properties. In particular, RXNFLOW outperforms the docking score for TacoGFN and attains high druglikeness while showing a competitive docking score with the state-of-the-art model, MolCRAFT-448 large. Moreover, RXNFLOW ensures the synthesizability comparable to known active ligands, 449 outperforming the fragment-based TacoGFN trained on the SA score objective and the distributional learning-based models trained on synthesizable drug molecules. Figure 3 illustrates generated molecules against TANK-binding kinase 1 (TBK1) which is not included in the training set. 452

An important finding is that RXNFLOW maintains high structural diversity (0.81) despite the typ-453 ical trade-off between optimization power and diversity (Gao et al., 2022a). This is a significant 454 improvement over the fragment-based TacoGFN, which scored 0.67 and is comparable to the distri-455 butional learning-based models that range from 0.83 to 0.87. We attribute this enhancement to our 456 action space, which contains chemically diverse building blocks, in contrast to the small and lim-457 ited fragment sets used in fragment-based GFlowNets. This suggests that our model can effectively 458 balance the potency and diversity of generated molecules. 459

460 4.3 INTRODUCING ADDITIONAL OBJECTIVE WITHOUT RETRAINING 461

462 In drug discovery, new objectives often arise during the research process, such as enhancing solubil-463 ity, reducing toxicity, or improving selectivity (Fink et al., 2022; Joshi et al., 2021). These additional 464 objectives typically not only require retraining models but also increase the optimization complex-465 ity. In this context, RXNFLOW can achieve some additional objectives by simply introducing con-466 straints to MDP without retraining thanks to the non-hierarchical action space (Sec. 3.4). As shown 467 in Figure 4, we explore the scenario of adding a solubility objective to a pre-trained GFlowNet in Sec. 4.2. Specifically, we target the generation of hydrophobic molecules by restricting the build-468 ing blocks with topological polar surface area (TPSA) in the bottom 15% ("low") and sampled 500 469 molecules for the KRAS-G12C mutant (PDB Id: 60im). While there are slight differences in the 470 QED distributions due to the correlation between TPSA and QED, the generated molecules are more 471 hydrophobic and retain similar overall reward distributions. We also performed the ablation study 472 of non-hierarchical MDP in Sec. D.6. 473

- 474 475
- 4.4 SCALING ACTION SPACE WITHOUT RETRAINING

476 The building block libraries for drug discovery continue to grow, from 60,000 in 2020 to over 1.2 477 million blocks today, to enhance chemical diversity and novelty (Grygorenko et al., 2020). However, 478 existing generative models require retraining to accommodate newly discovered building blocks, 479 limiting their scalability and adaptability. On the other hand, RXNFLOW can integrate new building 480 blocks without retraining by understanding the chemical context of actions through action embed-481 ding. We first divide the 1M-sized block library ("all") into two sets: 500,000 blocks for training 482 ("seen") and the remaining blocks ("unseen"). After training with the QED objective on various 483 reward exponent settings (R^{β}) , we generate 5,000 molecules from each set ("seen", "unseen", and "all"). Figure 5 shows that the reward distributions of samples are nearly identical, demonstrating 484 that RXNFLOW performs robustly with unseen building blocks. This result highlights the general-485 ization ability and scalability of RXNFLOW, a significant advantage for real-world applications.



Figure 6: Optimization power, diversity, and generation time according to building block library
size. (a-c) Average and standard deviation of properties of the top-1000 high-affinity molecules over
4 runs on pocket-specific generation. (a) Average docking score. (b) The uniqueness of BemisMurcko scaffolds. (c) Average Tanimoto distance. (d) Average runtime to generate 100 molecules
over the CrossDocked2020 test pockets in a zero-shot manner.

514

528

529 530

4.5 ABLATION STUDY: THE EFFECT OF BUILDING BLOCK LIBRARY SIZE

501 Expanding the size of building block libraries provides an opportunity to discover more diverse 502 and potent drug candidates (Grygorenko et al., 2020). In Sec. 4.1, RXNFLOW outperforms Syn-503 FlowNet and RGFN which use smaller block libraries, but differences in model architectures may 504 have contributed to these results. To isolate the effect of the building block library size, we conduct 505 an ablation study using partial libraries with a pocket-specific optimization task against the kappa-506 opioid receptor (PDB Id: 6b73), as illustrated in Figure 6(a-c). The results indicate that increasing 507 the library size enhances both optimization power, in terms of docking scores, and diversity, in terms 508 of a higher number of unique Bemis-Murcko scaffolds (Bemis & Murcko, 1996) and an increased 509 Tanimoto diversity of the generated molecules. Additionally, as shown in Figure 6(d), the generation time only doubles on the 100-fold larger action space, highlighting the efficiency of RXNFLOW. 510 These results demonstrate the forte of RXNFLOW in navigating a broader chemical space to discover 511 novel drug candidates by overcoming the computational limitations for expanding the action space. 512 We also investigated the scaling laws with other reaction-based GFlowNets in Sec. D.3. 513

515 4.6 TOWARDS FURTHER DEVELOPMENT

516 Our framework has room for improvement regarding more efficient 517 learning and sampling. First, the 3D interaction modeling can en-518 hance the generative performance. Given the high correlation be-519 tween binding affinity and conformation, such spatial relationships 520 could provide meaningful information for our generative model to make a reasonable decision. In Figure 7, we observed that the 3D 521 interaction modeling using docking conformations boosts the dis-522 covery of candidate molecules against the beta-2-adrenergic recep-523 tor. Second, the current action space subsampling method forces 524 exploration due to uniform sampling to minimize bias. We can en-525 hance the exploitation by prioritizing the building blocks of action 526 space subsampling instead of uniform subsampling. 527



Figure 7: Affect of 3D interaction modeling. Mean and standard error over 3 seeds

5 CONCLUSION

In this work, we introduce RXNFLOW, a synthesis-oriented generative framework designed to ex-531 plore broader chemical spaces, thereby enhancing both diversity and potency for drug discovery. 532 Our framework efficiently handles massive action spaces to expand the search space without signifi-533 cant computational or memory overhead by employing a novel action space subsampling technique. 534 RXNFLOW can effectively identify diverse drug candidates with desired properties and synthetic 535 feasibility by learning the objective of generative flow networks on synthetic pathways. Addition-536 ally, by formulating a non-hierarchical MDP, RXNFLOW can model generative flows on modified 537 action spaces, allowing it to achieve additional objectives and incorporate newly discovered build-538 ing blocks without retraining. These results highlight the potential of RXNFLOW as a practical and versatile solution for real-world drug discovery.

REFERENCES

541 Amr Alhossary, Stephanus Daniel Handoko, Yuguang Mu, and Chee-Keong Kwoh. Fast, accurate, 542 and reliable molecular docking with quickvina 2. Bioinformatics, 31(13):2214–2216, 2015. 543 544 Guy W Bemis and Mark A Murcko. The properties of known drugs. 1. molecular frameworks. Journal of medicinal chemistry, 39(15):2887–2893, 1996. 546 Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow 547 network based generative models for non-iterative diverse candidate generation. Advances in 548 Neural Information Processing Systems, 34:27381–27394, 2021. 549 550 Yoshua Bengio, Salem Lahlou, Tristan Deleu, Edward J Hu, Mo Tiwari, and Emmanuel Bengio. 551 Gflownet foundations. The Journal of Machine Learning Research, 24(1):10006–10060, 2023. 552 G Richard Bickerton, Gaia V Paolini, Jérémy Besnard, Sorel Muresan, and Andrew L Hopkins. 553 Quantifying the chemical beauty of drugs. *Nature chemistry*, 4(2):90–98, 2012. 554 555 John Bradshaw, Brooks Paige, Matt J Kusner, Marwin Segler, and José Miguel Hernández-Lobato. A model to search for synthesizable molecules. Advances in Neural Information Processing 556 Systems, 32, 2019. 558 Alexander Button, Daniel Merk, Jan A Hiss, and Gisbert Schneider. Automated de novo molecu-559 lar design by hybrid machine intelligence and rule-driven chemical synthesis. Nature machine 560 intelligence, 1(7):307–315, 2019. 561 Connor W Coley, Luke Rogers, William H Green, and Klavs F Jensen. Scscore: synthetic com-562 plexity learned from a reaction corpus. Journal of chemical information and modeling, 58(2): 563 252-261, 2018. 564 565 Miruna Cretu, Charles Harris, Julien Roy, Emmanuel Bengio, and Pietro Lio. Synflownet: Towards 566 molecule design with guaranteed synthesis pathways. In ICLR 2024 Workshop on Generative and 567 Experimental Perspectives for Biomolecular Design, 2024. URL https://openreview. net/forum?id=kjcBA2I2My. 568 569 Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, 570 Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep rein-571 forcement learning in large discrete action spaces. arXiv preprint arXiv:1512.07679, 2015. 572 Joseph L Durant, Burton A Leland, Douglas R Henry, and James G Nourse. Reoptimization of mdl 573 keys for use in drug discovery. Journal of chemical information and computer sciences, 42(6): 574 1273-1280, 2002. 575 576 Peter Ertl and Ansgar Schuffenhauer. Estimation of synthetic accessibility score of drug-like 577 molecules based on molecular complexity and fragment contributions. Journal of cheminfor-578 matics, 1:1-11, 2009. 579 Elissa A Fink, Jun Xu, Harald Hübner, Joao M Braz, Philipp Seemann, Charlotte Avet, Veronica 580 Craik, Dorothee Weikert, Maximilian F Schmidt, Chase M Webb, et al. Structure-based dis-581 covery of nonopioid analysics acting through the α 2a-adrenergic receptor. Science, 377(6614): 582 eabn7065, 2022. 583 Paul G Francoeur, Tomohide Masuda, Jocelyn Sunseri, Andrew Jia, Richard B Iovanisci, Ian Snyder, 584 and David R Koes. Three-dimensional convolutional neural networks and a cross-docked data set 585 for structure-based drug design. Journal of chemical information and modeling, 60(9):4200-586 4215, 2020. 587 588 Bowen Gao, Minsi Ren, Yuyan Ni, Yanwen Huang, Bo Qiang, Zhi-Ming Ma, Wei-Ying Ma, and

Bowen Gao, Minsi Ren, Yuyan Ni, Yanwen Huang, Bo Qiang, Zhi-Ming Ma, Wei-Ying Ma, and
 Yanyan Lan. Rethinking specificity in SBDD: Leveraging delta score and energy-guided dif fusion. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver,
 Jonathan Scarlett, and Felix Berkenkamp (eds.), Proceedings of the 41st International Con ference on Machine Learning, volume 235 of Proceedings of Machine Learning Research,
 pp. 14811–14825. PMLR, 21–27 Jul 2024a. URL https://proceedings.mlr.press/
 v235/gao24k.html.

594 595 596	Wenhao Gao and Connor W Coley. The synthesizability of molecules proposed by generative models. <i>Journal of chemical information and modeling</i> , 60(12):5714–5723, 2020.
597 598 599	Wenhao Gao, Tianfan Fu, Jimeng Sun, and Connor Coley. Sample efficiency matters: a benchmark for practical molecular optimization. <i>Advances in neural information processing systems</i> , 35: 21342–21357, 2022a.
600 601 602 603	Wenhao Gao, Rocío Mercado, and Connor W. Coley. Amortized tree generation for bottom-up synthesis planning and synthesizable molecular design. In <i>International Conference on Learning Representations</i> , 2022b. URL https://openreview.net/forum?id=FRxhHdnxt1.
604 605	Wenhao Gao, Shitong Luo, and Connor W Coley. Generative artificial intelligence for navigating synthesizable chemical space. <i>arXiv preprint arXiv:2410.03494</i> , 2024b.
606 607 608	Samuel Genheden, Amol Thakkar, Veronika Chadimová, Jean-Louis Reymond, Ola Engkvist, and Esben Bjerrum. Aizynthfinder: a fast, robust and flexible open-source software for retrosynthetic planning. <i>Journal of cheminformatics</i> , 12(1):70, 2020.
609 610 611 612 613	Sai Krishna Gottipati, Boris Sattarov, Sufeng Niu, Yashaswi Pathak, Haoran Wei, Shengchao Liu, Simon Blackburn, Karam Thomas, Connor Coley, Jian Tang, et al. Learning to navigate the synthetically accessible chemical space using reinforcement learning. In <i>International conference on machine learning</i> , pp. 3668–3679. PMLR, 2020.
614 615 616	David E Graff, Eugene I Shakhnovich, and Connor W Coley. Accelerating high-throughput virtual screening through molecular pool-based active learning. <i>Chemical science</i> , 12(22):7866–7881, 2021.
617 618 619 620	Oleksandr O Grygorenko, Dmytro S Radchenko, Igor Dziuba, Alexander Chuprina, Kateryna E Gu- bina, and Yurii S Moroz. Generating multibillion chemical space of readily accessible screening compounds. <i>Iscience</i> , 23(11), 2020.
621 622 623 624	Jiaqi Guan, Wesley Wei Qian, Xingang Peng, Yufeng Su, Jian Peng, and Jianzhu Ma. 3d equivariant diffusion for target-aware molecule generation and affinity prediction. In <i>The Eleventh Interna-</i> <i>tional Conference on Learning Representations</i> , 2023a. URL https://openreview.net/forum?id=kJqXEPXMsE0.
625 626 627	Jiaqi Guan, Xiangxin Zhou, Yuwei Yang, Yu Bao, Jian Peng, Jianzhu Ma, Qiang Liu, Liang Wang, and Quanquan Gu. Decompdiff: Diffusion models with decomposed priors for structure-based drug design. <i>ICML</i> , 2023b.
629 630 631	Markus Hartenfeller, Heiko Zettl, Miriam Walter, Matthias Rupp, Felix Reisen, Ewgenij Proschak, Sascha Weggen, Holger Stark, and Gisbert Schneider. Dogs: reaction-driven de novo design of bioactive compounds. <i>PLoS computational biology</i> , 8(2):e1002380, 2012.
632 633	Julien Horwood and Emmanuel Noutahi. Molecular design in synthetically accessible chemical space via deep reinforcement learning. <i>ACS omega</i> , 5(51):32984–32994, 2020.
635 636 637	Ruth Huey, Garrett M Morris, Stefano Forli, et al. Using autodock 4 and autodock vina with autodocktools: a tutorial. <i>The Scripps Research Institute Molecular Graphics Laboratory</i> , 10550 (92037):1000, 2012.
638 639 640 641	Moksh Jain, Emmanuel Bengio, Alex Hernandez-Garcia, Jarrid Rector-Brooks, Bonaventure FP Dossou, Chanakya Ajit Ekbote, Jie Fu, Tianyu Zhang, Michael Kilgour, Dinghuai Zhang, et al. Biological sequence design with gflownets. In <i>International Conference on Machine Learning</i> , pp. 9786–9801. PMLR, 2022.
642 643 644 645	Moksh Jain, Sharath Chandra Raparthy, Alex Hernández-Garcia, Jarrid Rector-Brooks, Yoshua Ben- gio, Santiago Miret, and Emmanuel Bengio. Multi-objective gflownets. In <i>International confer-</i> <i>ence on machine learning</i> , pp. 14631–14653. PMLR, 2023.
646 647	Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael John Lamarre Townshend, and Ron Dror. Learning from protein structure with geometric vector perceptrons. In <i>International Conference</i> <i>on Learning Representations</i> , 2020.

648 649 650	Tanuja Joshi, Priyanka Sharma, Tushar Joshi, Hemlata Pundir, Shalini Mathpal, and Subhash Chan- dra. Structure-based screening of novel lichen compounds against sars coronavirus main protease (mpro) as potentials inhibitors of covid-19. <i>Molecular Diversity</i> , 25:1665–1677, 2021.
652 653 654	John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. <i>Nature</i> , 596(7873):583–589, 2021.
655 656 657 658	Hyeongwoo Kim, Kyunghoon Lee, Chansu Kim, Jaechang Lim, and Woo Youn Kim. Dfrscore: deep learning-based scoring of synthetic complexity with drug-focused retrosynthetic analysis for high-throughput virtual screening. <i>Journal of Chemical Information and Modeling</i> , 64(7): 2432–2444, 2023.
659 660 661 662	Michał Koziarski, Andrei Rekesh, Dmytro Shevchuk, Almer van der Sloot, Piotr Gaiński, Yoshua Bengio, Cheng-Hao Liu, Mike Tyers, and Robert A Batey. Rgfn: Synthesizable molecular generation using gflownets. <i>arXiv preprint arXiv:2406.08506</i> , 2024.
663 664	Greg Landrum et al. Rdkit: A software suite for cheminformatics, computational chemistry, and predictive modeling. <i>Greg Landrum</i> , 8(31.10):5281, 2013.
665 666 667	Seul Lee, Jaehyeong Jo, and Sung Ju Hwang. Exploring chemical space with score-based out-of- distribution generation. In <i>International Conference on Machine Learning</i> , pp. 18872–18892. PMLR, 2023.
669 670	Meng Liu, Youzhi Luo, Kanji Uchino, Koji Maruhashi, and Shuiwang Ji. Generating 3d molecules for target protein binding. <i>arXiv preprint arXiv:2204.09410</i> , 2022.
671 672 673	Zhihai Liu, Yan Li, Li Han, Jie Li, Jie Liu, Zhixiong Zhao, Wei Nie, Yuchen Liu, and Renxiao Wang. Pdb-wide collection of binding data: current status of the pdbbind database. <i>Bioinformatics</i> , 31 (3):405–412, 2015.
674 675 676 677	Daniel Lowe. Chemical reactions from US patents (1976-Sep2016), 6 2017. URL https://figshare.com/articles/dataset/Chemical_reactions_from_US_patents_1976-Sep2016_/5104873.
678 679	Shitong Luo, Jiaqi Guan, Jianzhu Ma, and Jian Peng. A 3d generative model for structure-based drug design. <i>Advances in Neural Information Processing Systems</i> , 34:6229–6239, 2021.
680 681 682 683	Shitong Luo, Wenhao Gao, Zuofan Wu, Jian Peng, Connor W. Coley, and Jianzhu Ma. Projecting molecules into synthesizable chemical spaces. In <i>Forty-first International Conference on Machine Learning</i> , 2024. URL https://openreview.net/forum?id=scFlbJQdml.
684 685 686	Nikolay Malkin, Moksh Jain, Emmanuel Bengio, Chen Sun, and Yoshua Bengio. Trajectory balance: Improved credit assignment in gflownets. <i>Advances in Neural Information Processing Systems</i> , 35:5955–5967, 2022.
687 688	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representa- tions of words and phrases and their compositionality, 2013.
690 691 692	Harry L Morgan. The generation of a unique machine description for chemical structures-a technique developed at chemical abstracts service. <i>Journal of chemical documentation</i> , 5(2):107–113, 1965.
693 694 695	Trang Nguyen, Alexander Tong, Kanika Madan, Yoshua Bengio, and Dianbo Liu. Causal inference in gene regulatory networks with gflownet: Towards scalability in large systems. <i>arXiv preprint</i> <i>arXiv:2310.03579</i> , 2023a.
696 697 698 699	Tri Minh Nguyen, Sherif Abdulkader Tawfik, Truyen Tran, Sunil Gupta, Santu Rana, and Svetha Venkatesh. Hierarchical gflownet for crystal structure generation. In <i>AI for Accelerated Materials Design-NeurIPS 2023 Workshop</i> , 2023b.
700 701	Noel M O'Boyle, Michael Banck, Craig A James, Chris Morley, Tim Vandermeersch, and Geof- frey R Hutchison. Open babel: An open chemical toolbox. <i>Journal of cheminformatics</i> , 3:1–14, 2011.

702 703 704	Yongping Pan, Niu Huang, Sam Cho, and Alexander D Mackerell. Consideration of molecular weight during compound selection in virtual target-based database screening. <i>Journal of chemical</i>
704	information and computer sciences, 43(1):267–272, 2003.
706	Xingang Peng, Shitong Luo, Jiaqi Guan, Qi Xie, Jian Peng, and Jianzhu Ma. Pocket2mol: Effi-
707	cient molecular sampling based on 3d protein pockets. In International Conference on Machine
708	<i>Learning</i> , pp. 17644–17655. PMLR, 2022.
709	
710	Yanru Qu, Keyue Qiu, Yuxuan Song, Jingjing Gong, Jiawei Han, Mingyue Zheng, Hao Zhou, and Wai Ving Ma MalCD AFT: Structure based drug design in continuous perpendences. In Puslan
711	Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and
712	Felix Berkenkamp (eds.). Proceedings of the 41st International Conference on Machine Learn-
713	ing, volume 235 of Proceedings of Machine Learning Research, pp. 41749–41768. PMLR, 21–27
714	Jul 2024. URL https://proceedings.mlr.press/v235/qu24a.html.
715	
716 717	Matthew Ragoza, Tomohide Masuda, and David Ryan Koes. Generating 3d molecules conditional on receptor binding sites with deep generative models. <i>Chemical science</i> , 13(9):2701–2713, 2022.
718	Danny Reidenbach Evoshdd: Latent evolution for accurate and efficient structure-based drug de-
719 720	sign. In ICLR 2024 Workshop on Machine Learning for Genomics Explorations, 2024.
720	Arman A Sadybekov, Anastasiia V Sadybekov, Yongfeng Liu, Christos Iliopoulos-Tsoutsouvas, Xi-
721	Ping Huang, Julie Pickett, Blake Houser, Nilkanth Patel, Ngan K Tran, Fei Tong, et al. Synthon-
723	based ligand discovery in virtual libraries of over 11 billion compounds. Nature, 601(7893):
724	452–459, 2022.
725	Arna Schnauing, Yuangi Du, Charles Harris, Arian Jamash, Ilia Jaashay, Waitao Du, Tam Blundall
726	Pietro Lió Carla Gomes Max Welling Michael Bronstein and Bruno Correia Structure-based
727	drug design with equivariant diffusion models, 2023.
728	
729	Seonghwan Seo and Woo Youn Kim. Pharmaconet: Accelerating structure-based virtual screening
730	by pharmacophore modeling. arXiv preprint arXiv:2310.00681, 2023.
731	Seonghwan Seo, Jaechang Lim, and Woo Youn Kim. Molecular generative model via retrosynthet-
732	ically prepared chemical building block assembly. Advanced Science, 10(8):2206674, 2023.
733	
734	Tony Shen, Mohit Pandey, and Martin Ester. Tacogfn: Target conditioned gflownet for structure-
735	based drug design. arxiv preprint arxiv.2510.05225, 2025.
736	Viet-Khoa Tran-Nguyen, Célien Jacquemard, and Didier Rognan. Lit-pcba: an unbiased data set for
737	machine learning and virtual screening. Journal of chemical information and modeling, 60(9):
130	4263–4273, 2020.
740	Oleg Trott and Arthur I Olson Autodock vina: improving the speed and accuracy of docking with
7/11	a new scoring function, efficient optimization, and multithreading. <i>Journal of computational</i>
742	<i>chemistry</i> , 31(2):455–461, 2010.
743	
744	Siddarth Venkatraman, Moksh Jain, Luca Scimeca, Minsu Kim, Marcin Sendera, Mohsin Hasan,
745	Luke Rowe, Sarthak Mittal, Pablo Lemos, Emmanuel Bengio, et al. Amortizing intractable in-
746	2024
747	2027.
748	Shuzhe Wang, Jagna Witek, Gregory A Landrum, and Sereina Riniker. Improving conformer genera-
749	tion for small rings and macrocycles based on distance geometry and experimental torsional-angle
750	preferences. Journal of chemical information and modeling, 60(4):2044–2058, 2020.
751	Ronald I Williams Simple statistical gradient following algorithms for connectionist reinforcement
752	learning. Machine learning, 8:229–256, 1992.
753	
754	Yuejiang Yu, Chun Cai, Jiayue Wang, Zonghua Bo, Zhengdan Zhu, and Hang Zheng. Uni-dock:
755	Gpu-accelerated docking enables ultralarge virtual screening. <i>Journal of chemical theory and computation</i> , 19(11):3336–3345, 2023.

756 757 758	Seongjun Yun, Minbyul Jeong, Sungdong Yoo, Seunghun Lee, S Yi Sean, Raehyun Kim, Jaewoo Kang, and Hyunwoo J Kim. Graph transformer networks: Learning meta-path graphs to improve gnns. <i>Neural Networks</i> , 153:104–119, 2022.
759	Alex Zhousenhou Van A Juananhou Alex Aliner Mark & Vasalou Vladimin A Aladinahiu Anas
760	Alex Zilavololikov, Tali A Ivanelikov, Alex Aliper, Mark 5 veselov, vlaulilli A Alauliiskiy, Alas- tasiya V Aladinskaya Victor A Terentiev, Daniil A Polykovskiy, Maksim D Kuznetsov, Arin
761	Asadulaev, et al. Deep learning enables rapid identification of potent ddr1 kinase inhibitors. Na-
762	ture biotechnology, 37(9):1038–1040, 2019.
763	
764	Mingyang Zhou, Zichao Yan, Elliot Layne, Nikolay Malkin, Dinghuai Zhang, Moksh Jain, Mathieu
765	Blanchette, and Yoshua Bengio. Phylogin: Phylogenetic inference with generative flow networks.
700	<i>urxiv preprint urxiv.2510.06774, 2025.</i>
768	Wonho Zhung, Hyeongwoo Kim, and Woo Youn Kim. 3d molecular generative framework for
769	interaction-guided drug design. Nature Communications, 15(1):2688, 2024.
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
707	
700	
709	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
000	
808	
809	

A THEORETICAL ANALYSIS 811

In this section, we provide the theoretical background of action subsampling. We define \mathcal{U} as the uniform subsampling policy, i.e., $\mathcal{B}^* \sim \text{Uniform}(\{\mathcal{B}^* \subseteq \mathcal{B} \mid |\mathcal{B}^*| = k\})$.

Bias of log forward policy. In this section, we prove that the log forward policy estimation is unbiased (Eq. (7)) with the following weights:

$$w_a = \begin{cases} |\mathcal{B}|/|\mathcal{B}^*| & \text{if } a \in \mathcal{B} \\ |\mathcal{B}_r|/|\mathcal{B}_r^*| & \text{if } a \text{ is } (r,b) \text{ and } r \in \mathcal{R}_2 \\ 1 & \text{if } a \text{ is } \text{Stop or } a \in \mathcal{R}_1 \end{cases}$$

For readability, we express the edge flow F and forward policy P_F where the action a is $s \to s'$ as follows:

$$F_{\theta}(s \to s'; \theta) = F_{\theta}(s, a; \theta), \quad P_F(s'|s; \theta) = P_F(a|s; \theta)$$

Then, the forward policy P_F and the estimated policy \hat{P}_F (Eq. (7)) can be rewritten as follows:

$$P_F(s'|s;\theta) = P_F(a|s;\theta) = \frac{F_{\theta}(s,a)}{F_{\theta}(s)} = \frac{F_{\theta}(s,a)}{\sum_{a'\in\mathcal{A}(s)}F_{\theta}(s,a')}$$
$$\hat{P}_F(s'|s;\mathcal{A}^*;\theta) = \hat{P}_F(a|s;\mathcal{A}^*;\theta) = \frac{F_{\theta}(s,a)}{\hat{F}_{\theta}(s;\mathcal{A}^*)} = \frac{F_{\theta}(s,a)}{\sum_{a'\in\mathcal{A}^*(s)}w_{a'}F_{\theta}(s,a')}$$

The expectation of the estimated initial state flow $\hat{F}_{\theta}(s_0; \mathcal{A}^*) = \sum_{a \in \mathcal{A}^*} w_a F_{\theta}(s_0, a)$ is given by

$$\mathbb{E}_{\mathcal{A}^* \sim \mathcal{P}(\mathcal{A})}[\hat{F}_{\theta}(s_0; \mathcal{A}^*)] = \mathbb{E}_{\mathcal{B}^* \sim \mathcal{U}(\mathcal{B})} \left[\frac{|\mathcal{B}|}{|\mathcal{B}^*|} \sum_{b \in \mathcal{B}^*} F_{\theta}(s_0, b) \right] = \sum_{b \in \mathcal{B}} F_{\theta}(s_0, b) = F_{\theta}(s_0), \quad (13)$$

For a later state $s \neq s_0$, the expectation of the state flow $\hat{F}_{\theta}(s; \mathcal{A}^*)$ is given by

$$\mathbb{E}_{\mathcal{A}^* \sim \mathcal{P}(\mathcal{A})}[\hat{F}_{\theta}(s; \mathcal{A}^*)] = \sum_{a \in \mathcal{R}_1 \cup \{\text{Stop}\}} F_{\theta}(s, a) + \sum_{r \in \mathcal{R}_2} \mathbb{E}_{\mathcal{B}_r^* \sim \mathcal{U}(\mathcal{B}_r)} \left[\frac{|\mathcal{B}_r|}{|\mathcal{B}_r^*|} \sum_{b \in \mathcal{B}_r^*} F_{\theta}(s, (r, b)) \right]$$
$$= \sum_{a \in \mathcal{R}_1 \cup \{\text{Stop}\}} F_{\theta}(s, a) + \sum_{r \in \mathcal{R}_2} \sum_{b \in \mathcal{B}_r} F_{\theta}(s, (r, b))$$
$$= P_F(s'|s; \theta) \tag{14}$$

Variance of log forward policy. We define the standard deviation of the probability $P_F(-|s;\theta)$ as $\sigma_{\theta}(-|s)$. For the initial state s_0 , the variance of $\log \hat{P}_F(s'|s_0; \mathcal{A}^*; \theta)$ is given by

$$\operatorname{Var}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\log \hat{P}_{F}(s'|s_{0}; \mathcal{A}^{*}; \theta) \right]$$

$$= \operatorname{Var}_{\mathcal{B}^{*} \sim \mathcal{U}(\mathcal{B})} \left[\underbrace{\log \mathcal{E}_{\theta}(s_{0} \rightarrow s')}_{b \in \mathcal{B}^{*}} - \log \left(\frac{|\mathcal{B}|'}{|\mathcal{B}^{*}|} \sum_{b \in \mathcal{B}^{*}} F_{\theta}(s_{0}, b) \right) \right]$$

$$= \operatorname{Var}_{\mathcal{B}^{*} \sim \mathcal{U}(\mathcal{B})} \left[\log \sum_{b \in \mathcal{B}^{*}} P_{F}(b|s_{0}; \theta) \right] \qquad \text{(by normalizing with } F_{\theta}(s_{0}))$$

$$\approx \frac{|\mathcal{B}|^{2}(|\mathcal{B}| - |\mathcal{B}^{*}|)}{|\mathcal{B}^{*}|(|\mathcal{B}| - 1)} \sigma_{\theta}(-|s_{0})^{2} \quad \text{(by Eq. (24))} \qquad (15)$$

For the later state $s \neq s_0$, the variance of $\log \hat{P}_F(s'|s; \mathcal{A}^*; \theta)$ is:

 $\operatorname{Var}_{\mathcal{A}^* \sim \mathcal{P}(\mathcal{A})} \left[\log \hat{P}_F(s'|s; \mathcal{A}^*; \theta) \right]$

$$= \operatorname{Var}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\log \left(\sum_{a \in \mathcal{R}_{1} \cup \{ \text{Stop} \}} P_{F}(a|s;\theta) + \sum_{r \in \mathcal{R}_{2}} \frac{|\mathcal{B}_{r}|}{|\mathcal{B}_{r}^{*}|} \sum_{b \in \mathcal{B}_{r}^{*}} P_{F}((r,b)|s;\theta) \right) \right]$$

$$\approx \frac{\operatorname{Var}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\sum_{a \in \mathcal{R}_{4} \cup \{ \text{Stop} \}} P_{F}(a|s;\theta) + \sum_{r \in \mathcal{R}_{2}} \frac{|\mathcal{B}_{r}|}{|\mathcal{B}_{r}^{*}|} \sum_{b \in \mathcal{B}_{r}^{*}} P_{F}((r,b)|s;\theta) \right]}{\mathbb{E}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\sum_{a \in \mathcal{R}_{1} \cup \{ \text{Stop} \}} P_{F}(a|s;\theta) + \sum_{r \in \mathcal{R}_{2}} \frac{|\mathcal{B}_{r}|}{|\mathcal{B}_{r}^{*}|} \sum_{b \in \mathcal{B}_{r}^{*}} P_{F}((r,b)|s;\theta) \right]^{2}}$$

$$= \frac{1}{\left(\sum_{a \in \mathcal{A}(s)} P_{F}(a|s;\theta) \right)^{2}} \operatorname{Var}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A}(s))} \left[\sum_{r \in \mathcal{R}_{2}} \frac{|\mathcal{B}_{r}|}{|\mathcal{B}_{r}^{*}|} \sum_{b \in \mathcal{B}_{r}^{*}} P_{F}((r,b)|s;\theta) \right]$$

$$= \sum_{r \in \mathcal{R}_{2}} \frac{|\mathcal{B}_{r}|^{2} (|\mathcal{B}_{r}| - |\mathcal{B}_{r}^{*}|)}{|\mathcal{B}_{r}^{*}| (|\mathcal{B}_{r}| - 1)} \sigma_{\theta}((r, -)|s)^{2} \quad \text{(by Eq. (22))} \quad (16)$$

Finally, we get the variance of $\log \hat{P}_F(\tau; \theta)$ where $\tau = (s_0 \to ... \to s_n = x)$:

$$\operatorname{Var}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\log \hat{P}_{F}(\tau; \mathcal{A}^{*}; \theta) \right] \approx \frac{|\mathcal{B}|^{2}(|\mathcal{B}| - |\mathcal{B}^{*}|)}{k|\mathcal{B}^{*}|(|\mathcal{B}| - 1)} \sigma_{\theta}(-|s_{0})^{2} + \sum_{t=1}^{n-1} \sum_{r \in \mathcal{R}_{2}} \frac{|\mathcal{B}_{r}|^{2}(|\mathcal{B}_{r}| - |\mathcal{B}^{*}_{r}|)}{|\mathcal{B}^{*}_{r}|(|\mathcal{B}_{r}| - 1)} \sigma_{\theta}((r, -)|s_{t})^{2}$$
(17)

Expectation of trajectory balance loss. Due to the variance of forward policy, there is a bias of the trajectory balance loss equal to the variance of $\log \hat{P}_F(\tau)$:

$$\mathbb{E}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} [\mathcal{L}_{\mathrm{TB}}(\tau)] = \mathbb{E}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\log \frac{Z_{\theta} \hat{P}_{F}(\tau; \theta)}{R(x) P_{B}(\tau | x)} \right]^{2} + \operatorname{Var}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\log \frac{Z_{\theta} \hat{P}_{F}(\tau; \theta)}{R(x) P_{B}(\tau | x)} \right] = \mathcal{L}_{\mathrm{TB}}(\tau) + \operatorname{Var}_{\mathcal{A}^{*} \sim \mathcal{P}(\mathcal{A})} \left[\log \hat{P}_{F}(\tau; \theta) \right]$$
(18)

$$\approx \mathcal{L}_{\mathrm{TB}}(\tau) + \frac{|\mathcal{B}|^2(|\mathcal{B}| - |\mathcal{B}^*|)}{|\mathcal{B}^*|(|\mathcal{B}| - 1)} \sigma_{\theta}(-|s_0)^2 + \sum_{t=1}^{n-1} \sum_{r \in \mathcal{R}_2} \frac{|\mathcal{B}_r|^2(|\mathcal{B}_r| - |\mathcal{B}_r^*|)}{|\mathcal{B}_r^*|(|\mathcal{B}_r| - 1)} \sigma_{\theta}((r, -)|s_t)^2 \quad (19)$$

Therefore, the gradient of trajectory balance loss is:

$$\mathbb{E}_{\mathcal{A}^* \sim \mathcal{P}(\mathcal{A})} \left[\nabla_{\theta} \hat{\mathcal{L}}_{\text{TB}}(\tau) \right]$$

$$\approx \nabla_{\theta} \mathcal{L}_{\text{TB}}(\tau) + \frac{|\mathcal{B}|^2 (|\mathcal{B}| - |\mathcal{B}^*|)}{|\mathcal{B}^*| (|\mathcal{B}| - 1)} \nabla_{\theta} \sigma_{\theta} (-|s_0)^2 + \sum_{t=1}^{n-1} \sum_{r \in \mathcal{R}_2} \frac{|\mathcal{B}_r|^2 (|\mathcal{B}_r| - |\mathcal{B}^*_r|)}{|\mathcal{B}^*_r| (|\mathcal{B}_r| - 1)} \nabla_{\theta} \sigma_{\theta} ((r, -)|s_t)^2$$
(20)

Given that $\sigma_{\theta}(\cdot)$, which is the standard deviation of probability to select actions, is inversely pro-portional to $|\mathcal{B}|$ or $|\mathcal{B}_r|$, the bias is highly dependent on the subsampling size (e.g. $|\mathcal{B}^*|$) rather than the subsampling ratio (e.g. $|\mathcal{B}^*|/|\mathcal{B}|$). Moreover, we can decrease the bias by MC sampling for state flow estimation, and the bias is reversely proportional to the number of samples k. However, in the same computational cost (proportional to $\mathcal{O}(k \times \sum_r |\mathcal{B}_r^*|)$), using the larger partial action spaces without MC sampling is more precise than using smaller partial action spaces with multiple MC samples. In Sec. D.8, we experimentally show that the bias is relatively small to trajectory balance loss with the toy experiment.

918 A.1 THEORETICAL BACKGROUNDS 919

920 Variance of uniformly sampled subset For the set \mathcal{X} with a size of n, we define its uniformly 921 sampled subset with size m as $\mathcal{X}' \sim \mathcal{P}(\mathcal{X})$. For the function $F(\mathcal{X}) = \sum_{x \in \mathcal{X}} f(x)$, we define the 922 unbiased estimation:

$$\hat{F}(\mathcal{X}) = \frac{n}{m} F(\mathcal{X}') = \frac{n}{m} \sum_{x \in \mathcal{X}'} f(x)$$

When the mean and standard deviation of f(x) is μ_f and σ_f , variance of $\hat{F}(\mathcal{X})$ is:

$$\begin{aligned} \operatorname{Var}_{\mathcal{X}' \sim \mathcal{P}(\mathcal{X})} \left[\hat{F}(\mathcal{X}) \right] \\ &= \mathbb{E}_{\mathcal{X}' \sim \mathcal{P}(\mathcal{X})} \left[\hat{F}(\mathcal{X})^2 \right] - \mathbb{E}_{\mathcal{X}' \sim \mathcal{P}(\mathcal{X})} \left[\hat{F}(\mathcal{X}) \right]^2 \\ &= \frac{1}{nC_m} \sum_{\mathcal{X}' \in \mathcal{P}(\mathcal{X})} \left(\frac{n}{m} \sum_{x \in \mathcal{X}'} f(x) \right)^2 - F(\mathcal{X})^2 \\ &= \frac{1}{nC_m} \frac{n^2}{m^2} \left({}^{n-1}C_{m-1} \sum_{i=1}^n f(b_i)^2 + {}^{n-2}C_{m-2} \sum_{i=1}^n \sum_{j>i}^n 2f(b_i)f(b_j) \right) - \left(\sum_{x \in \mathcal{X}} f(x) \right)^2 \\ &= \frac{n}{m} \sum_{i=1}^n f(x_i)^2 + \frac{2n(m-1)}{m(n-1)} \sum_{i=1}^n \sum_{j>i}^n f(x_i)f(x_j) - \left(\sum_{i=1}^n f(x_i)^2 + 2 \sum_{i=1}^n \sum_{j>i}^n f(x_i)f(x_j) \right) \\ &= \frac{n-m}{m(n-1)} \left((n-1) \sum_{i=1}^n f(x_i)^2 - 2 \sum_{i=1}^n \sum_{j>i}^n f(x_i)f(x_j) \right) \\ &= \frac{n-m}{m(n-1)} \sum_{i=1}^n \sum_{j>i}^n (f(x_i) - f(x_j))^2 \\ &= \frac{n^2(n-m)}{m(n-1)} \sigma_f^2 \end{aligned} \tag{21}$$

For MC sampling with k samples, the variance is:

$$\operatorname{Var}_{\mathcal{X}'\sim\mathcal{P}(\mathcal{X})}\left[\hat{F}(\mathcal{X})\right] = \frac{n^2(n-m)}{km(n-1)}\sigma_f^2$$
(22)

The variance of $\log \hat{F}(\mathcal{X})$ is:

$$\operatorname{Var}_{\mathcal{X}'\sim\mathcal{P}(\mathcal{X})}\left[\log\hat{F}(\mathcal{X})\right] \approx \frac{\operatorname{Var}_{\mathcal{X}'\sim\mathcal{P}(\mathcal{X})}\left[\hat{F}(\mathcal{X})\right]}{\mathbb{E}_{\mathcal{X}'\sim\mathcal{P}(\mathcal{X})}\left[\hat{F}(\mathcal{X})\right]^2} = \frac{(n-m)}{km(n-1)}(\sigma_f/\mu_f)^2$$
(23)

When the $F(\mathcal{X}) = 1$, the variance is:

$$\operatorname{Var}_{\mathcal{X}'\sim\mathcal{P}(\mathcal{X})}\left[\log\hat{F}(\mathcal{X})\right] \approx \frac{m^2(n-m)}{km(n-1)}\sigma_f^2 \tag{24}$$

972 B RXNFLOW ARCHITECTURE 973

974 B.1 FORWARD POLICY 975

We use the model architecture inspired from Cretu et al. (2024); Koziarski et al. (2024). We used a graph transformer (Yun et al., 2022) as the backbone f_{θ} following Bengio et al. (2021) and a multilayer perceptron (MLP) for action embedding g_{θ} . The graph embedding dimension is d_1 , and the building block embedding dimension is d_2 . The molecular graph for a state is *s*, and the GFlowNet condition vector is *c* which includes a reward exponent and multi-objective optimization weights (Jain et al., 2023). $a \parallel b$ means the concatenation of two feature vectors *a* and *b*.

982 983 984 Initial block selection. For the first action, the model always selects AddFirstReactant action which selects $b \in \mathcal{B}$ for the starting molecule with MLP_{AddFirstReactant} : $\mathbb{R}^{d_1+d_2} \to \mathbb{R}$.

$$F_{\theta}(s_0, b, c) = \mathsf{MLP}_{\mathsf{AddFirstReactant}}(f_{\theta}(s_0, c) \| g_{\theta}(b)) \tag{25}$$

986 **Reaction selection.** For the later states $s \neq s_0$, the model calculates the logits for the Stop action, 987 ReactUni actions $r_1 \in \mathcal{R}_1$, and ReactBi actions $(r_2, b) \in \mathcal{R}_2 \times \mathcal{B}$.

988 The logit for Stop is calculated by $MLP_{Stop} : \mathbb{R}^{d_1} \to \mathbb{R}$: 989

$$F_{\theta}(s, \text{Stop}, c) = \text{MLP}_{\text{Stop}}(f_{\theta}(s, c)).$$
(26)

991 The logit for a ReactUni action $r_1 \in \mathcal{R}_1$ is calculated by $MLP_{\text{ReactUni}}^{r_1} : \mathbb{R}^{d_1} \to \mathbb{R}$: 992 $F_{\theta}(s, r_1, c) = MLP_{\text{ReactUni}}^{r_1}(f_{\theta}(s, c)).$ (27)

Finally, the logit for a ReactBi action (r_2, b) is calculated by the one-hot embedding of reaction template $\delta(r_2): \{0, 1\}^{|\mathcal{R}_2|}$ and MLP_{ReactBi}: $\mathbb{R}^{d_1+|\mathcal{R}_2|+d_2} \to \mathbb{R}$:

$$F_{\theta}(s, (r_2, b), c) = \mathsf{MLP}_{\mathsf{ReactBi}}(f_{\theta}(s, c) \| \delta(r_2) \| g_{\theta}(b)).$$
(28)

Pocket conditioning. For a pocket-conditional generation, the model uses a *K*-NN pocket residual graph \mathcal{G}^P and encodes GVP-GNN (Jing et al., 2020) according to Shen et al. (2023). The pocket conditions are included in the GFlowNet condition vector *c*.

3D interaction modeling. Instead of using a 2D molecular graph of the ligand, the model can utilize a 3D binding complex graph to represent the state as the input of the graph transformer backbone f_{θ} . In this graph, each ligand atom connects to its K_1 nearest protein atoms as well as to its neighboring ligand atoms, while each protein atom connects to its K_2 nearest protein atoms. The edges encode spatial relationship information between the connected nodes.

Furthermore, we lightweight the MDP formulation to implement computationally intensive 3D interaction modeling with modified layers $MLP^*_{AddFirstReactant}$: $\mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$ and $MLP^*_{ReactBi}$: $\mathbb{R}^{d_1+|\mathcal{R}_2|} \to \mathbb{R}^{d_2}$ as follows:

 $F_{\theta}(s_0, b, c) = \mathsf{MLP}^*_{\mathsf{AddFirstReactant}}(f_{\theta}(s_0, c)) \odot g_{\theta}(b), \tag{29}$

$$F_{\theta}(s, (r_2, b), c) = \mathsf{MLP}^*_{\mathsf{ReactBi}}(f_{\theta}(s, c) \| \delta(r_2)) \odot g_{\theta}(b), \tag{30}$$

1011 1012 1013

1023 1024 1025

1010

985

990

996 997

1014 B.2 BACKWARD POLICY.

1015 Malkin et al. (2022) introduced a uniform backward policy for small-scale drug discovery. However, 1016 in a directed acyclic graph (DAG) on synthetic pathways, where each state has numerous outgoing 1017 edges but few incoming ones, there is a significant imbalance in the number of trajectories along 1018 each incoming edge, depending on the distance from the initial state. In a uniform backward policy, 1019 where the flow of all incoming edges is equal, this imbalance diminishes the flow of trajectories that 1020 reach a state via shorter routes, i.e., those with fewer reaction steps. To facilitate shorter synthetic pathways as trajectories, we set the backward transition probability proportional to the expected 1021 number of trajectories along each incoming edge of the state, $a: s'' \rightarrow s$: 1022

$$P_B(s''|s) = P((s'' \to s) \in \tau | s \in \tau) = \frac{\sum_{\tau \in \mathcal{T}: \tau = (s_0 \to \dots \to s'' \to s)} |\mathcal{A}(s)|^{N-|\tau|}}{\sum_{\tau \in \mathcal{T}: \tau = (s_0 \to \dots \to s)} |\mathcal{A}(s)|^{N-|\tau|}}$$
(31)

where N is the maximum length of trajectories.

1026 **B.3** BUILDING BLOCK REPRESENTATION FOR ACTION EMBEDDING. 1027

1028 To represent building blocks, we used both physicochemical properties, chemical structural properties, and topological properties. For physicochemical properties, we used 8 molecular features: 1029 molecular weight, the number of atoms, the number of H-bond acceptors/donors (HBA, HBD), the 1030 number of aromatic/non-aromatic rings, LogP, and TPSA. For chemical properties, we used the 1031 MACCS fingerprint (Durant et al., 2002), which represents the composition of the chemical func-1032 tional groups in the molecule. For topological information, we used the Morgan ECFP4 fingerprint 1033 (Morgan, 1965) with a dimension of 1024, which is widely used in fingerprint-based deep learning 1034 researches. All properties are calculated with RDKit (Landrum et al., 2013). 1035

1036 1037

B.4 GFLOWNET TRAINING.

The training algorithm with the trajectory balance objective is described at Algorithm 1. The nota-1039 tions are in Sec. 3.3. 1040

Algorithm 1 Training GFlowNets with action space subsampling					
1:	input Entire action space \mathcal{A} , Maximum trajectory length N				
2:	repeat				
3:	Sample partial action spaces $\mathcal{A}_0^*, \mathcal{A}_1^*,, \mathcal{A}_{N-1}^*$ from subsampling policy $\mathcal{P}(\mathcal{A})$				
4:	Sample trajectory τ from sampling policy $\pi_{\theta}(- -; \mathcal{A}^*_{(-)})$				
5:	Update model $\theta \leftarrow \theta - \eta \nabla \hat{\mathcal{L}}_{TB}(\tau)$				
6:	until model converges				

1049 **B.5** COMPARISON BETWEEN THREE GFLOWNET METHODS

1050 Table 6: Comparison with SynFlowNet and RGFN. If the method includes each technique, it is 1051 denoted by an \bigcirc , otherwise by an \times . RGFN investigates scaling up to 64,000 building blocks, but 1052 their experimental validation and proof-of-principle implementations use only 350.

1		1 1 1	1		
Methods	Hierarchical	Action Embedding	ReactUni	Num Blocks	Max Reaction Steps
Enamine REAL	-	-	0	>1M	3
RGFN	0	0	×	350-64,000	4
SynFlowNet	Ō	×	\bigcirc	6,000	4
SynFlowNet(upd)	Õ	\bigcirc	Ŏ	10,000-220,000	3/4
RXNFLOW	×	Ŏ	Ŏ	>1M	3

1058 1059

1048

While RXNFLOW shares similar rules with SynFlowNet (Cretu et al., 2024) and RGFN (Koziarski 1061 et al., 2024), there are several major differences as described in Table 6. Since, there are two versions 1062 of SynFlowNet, the version presented t the ICLR 2024 GEM workshop and the version updated on 1063 October 16, 2024, we refer to the updated version as SynFlowNet(upd). For the benchmark studies, 1064 we use the workshop version of SynFlowNet.

Backward Trajectory. In contrast to the simple fragment-based or atom-based GFlowNet, some 1066 backward transitions do not reach the initial state s_0 . While RGFN and SynFlowNet do not con-1067 sider the invalid transitions, SynFlowNet(upd) and RXNFLOW address this issue. To prevent the 1068 invalid backward transition, SynFlowNet(upd) introduces additional maximum likelihood and RE-1069 INFORCE (Williams, 1992) training objectives for parameterized backward policy to prefer the 1070 backward transitions which can reach the initial state. In contrast, we implement explicit retrosynthetic analysis within maximum reaction steps for each state to collect valid backward trajectories. 1071

1072 Flow Network Framework. In contrast to SynFlowNet which uses discrete action space, RGFN, 1073 SynFlowNet(upd), and RXNFLOW include the action embedding (Dulac-Arnold et al., 2015) for 1074 chemical building blocks. While every building block must be sampled and trained at least once 1075 to prevent unsafe actions in a discrete action space, action embedding expresses actions based on 1076 structural information-domain knowledge. In particular, we formulate the non-hierarchical MDP 1077 instead of the hierarchical MDP for bi-molecular reactions. It is adaptable to modified building block libraries as described in Sec. B.6 but increases the computational cost and memory require-1078 ment dramatically. Thanks to action space subsampling, RXNFLOW can use non-hierarchical MDP 1079 without a computational burden problem.

1080 Sampling algorithm. For each forward transition step, RGFN and SynFlowNet (both versions) 1081 model the flow for the entire action space to estimate the forward transition probability. Therefore, 1082 the computational cost and memory requirement is proportional to the number of using building 1083 blocks \mathcal{B} , i.e., $\mathcal{O}(|\mathcal{B}|)$. Instead of modeling all actions, we propose an action space subsampling 1084 technique, which is similar to the negative sampling technique (Mikolov et al., 2013) in natural language processing, to estimate the forward transition probability from the subset of action space. It makes a cost-variance trade-off to decrease the computational cost complexity to $\mathcal{O}(r|\mathcal{B}|)$, where 1086 r is a subsampling ratio. As a result, RXNFLOW can handle millions of building blocks on training 1087 and inference with the computational cost of handling 10,000 building blocks. 1088

1089 1090

1091

1092

1093

1094

1095

1133

B.6 NON-HIERARCHICAL MARKOV DECISION PROCESS

To describe the difference, we supposed the situation that the toxicity is observed in molecules containing a functional group $-NR_3$, thereby excluding blocks with the function group from the building block library. In trained hierarchical MDP, the modification of the block library cannot change the probability of reaction templates, leading to the overestimation or underestimation of edge flows (the red color of Figure 8(a)). However, in the non-hierarchical MDP, the exclusion of 1096 some reactants does not affect actions that share the same reaction template (Figure 8(b)).



Figure 8: **Illustration of a situation** where an additional objective is introduced: excluding building 1130 blocks (reactants) containing the functional group $-NR_3$. The Gray dashed lines mean masked 1131 actions. (a) GFlowNet on hierarchical MDP. (b) GFlowNet on non-hierarchical MDP. 1132

1134 C **EXPERIMENTAL DETAILS** 1135

1136 C.1 ACTION SPACE DETAILS 1137

1138 **Reaction templates.** In this work, we used the reaction template set constructed by Cretu et al. (2024) from two public collections Hartenfeller et al. (2012); Button et al. (2019). The entire reaction 1139 template set includes 71 reaction templates, 13 for uni-molecular reactions and 58 for bi-molecular 1140 reactions. In *in silico* reactions with bi-molecular reaction templates, the products depend on the 1141 order of the input reactants. To ensure the consistency of the action, we consider two templates for 1142 each bi-molecular template according to the order of reactants, i.e. $|\mathcal{R}_1| = 13$, $|\mathcal{R}_2| = 116$. We note 1143 that our template set does not contain templates that have the same first and second reactant patterns. 1144

1145 **Building blocks.** We used the Enamine comprehensive catalog with 1,309,385 building blocks 1146 released on 2024.06.10 (Grygorenko et al., 2020). We filtered out building blocks that are not 1147 RDKit-readable, have no possible reactions, or contain unallowable atoms², resulting in 1,193,871 1148 remaining blocks.

1149

1155

1150 C.2 GFLOWNET TRAINING DETAILS 1151

1152 To minimize optimization performance influencing factors, we mostly followed the standard GFlowNet's model architecture and hyperparameters³ except for some parameters in Table 7. All 1153 experiments were performed on a single NVIDIA RTX A4000 GPU. 1154

Table 7: **Default hyperparameters** used in RXNFLOW training.

1156	Table 7: Default hyperparamet	ers used in RXNFLOW training.
1157	Hyperparameters	Values
1158	Minimum trajectory length	2 (minimum reaction steps: 1)
1159	Maximum trajectory length	4 (maximum reaction steps: 3)
1160	GFN temperature β	Uniform(0, 64)
1161	Train random action probability	0.05 (5%)
1162	Action space subsampling ratio	1%
1163	Building block embedding size	64

1164

For action space subsampling, we randomly subsample 1% actions for AddFirstReactant and 1165 each bi-molecular reaction template $r \in \mathcal{R}_2$. However, for bi-molecular reactions with small pos-1166 sible reactant block sets $\mathcal{B}_r \in \mathcal{B}$, the memory benefit from the action space subsampling is small 1167 while a variance penalty is large. Therefore, we set the minimum subsampling size to 100 for each 1168 bi-molecular reaction, and the action space subsampling is not performed when the number of ac-1169 tions is smaller than 100. 1170

The number of actions for each action type is imbalanced, and the number of reactant blocks (\mathcal{B}_r) for 1171 each bi-molecular reaction template r is also imbalanced. This can make some rare action categories 1172 not being sampled during training. We empirically found that ReactBi action were only sampled 1173 during 20,000 iterations (1.28M samples) in a toy experiment that uses one bi-molecular reaction 1174 template and 10,000 building blocks in some random seeds. Therefore, we set the random action 1175 probability as the default of 5%, and the model uniformly samples each action category in the 1176 random action sampling. This prevents incorrect predictions by ensuring that the model experiences 1177 trajectories including rare actions. We note that this random selection is only performed during 1178 model training.

1179 1180

1181

1184 1185

C.3 EXPERIMENTAL DETAILS

Pocket-specific optimization. For pocket-specific optimization with GPU-accelerated UniDock, 1182 we normalize the Vina and SA scores for multi-objective optimization: 1183

$$\widehat{\text{Vina}}(x) = -0.1 \max(\text{Vina}(x), 0), \quad \widehat{\text{SA}}(x) = (10 - \text{SA}(x))/9.$$
 (32)

²The allowable atom types are B, C, N, O, F, P, S, Cl, Br, I. 1186

³Default hyperparameters in https://github.com/recursionpharma/gflownet/blob/ 1187 trunk/src/gflownet/tasks/seh_frag.py

For FragGFN, we set the maximum trajectory length to 9, and for SynFlowNet and RGFN, we used the same hyperparameters as our framework except for the maximum trajectory length and building block library size. According to their default setting, we set a the maximum trajectory length to 5 rather than 4, and we randomly sampled 350 building blocks for RGFN and 6000 building blocks for SynFlowNet. Moreover, we do not perform action subsampling for SynFlowNet and RGFN.

Pocket-conditional generation in a zero-shot manner. We used the modified version of the TacoGFN's reward function and training set. Since we don't need to optimize SA (Ertl & Schuffenhauer, 2009), we excluded the SA term from the reward functions:

$$\begin{array}{c} 1197 \\ 1198 \\ 1199 \\ 1199 \\ 1199 \\ 1200 \\ 1200 \\ 1201 \\ 1201 \\ 1201 \\ 1202 \\ 1202 \\ 1202 \\ 1202 \\ 1202 \\ 1202 \\ 1202 \\ 1203 \\ 1204 \\ 1205 \\ 1204 \\ 1205 \\ 1204 \\ 1205 \\ 1206 \\ 1206 \\ 1206 \\ 1206 \\ 1206 \\ 1206 \\ 1206 \\ 1207 \\ 1208 \\ 1209 \\ 1208 \\ 1209 \\ 1209 \\ 1208 \\ 1209 \\ 1209 \\ 1209 \\ 1200 \\ 1200 \\ 1200 \\ 1210 \\ 1210 \\ 1210 \\ 1210 \\ 1211 \\ 1$$

$$\operatorname{RxnFlow-Reward}(x) = \frac{r_{\operatorname{affinity}}(x) \times r_{\operatorname{QED}}(x)}{\sqrt[3]{\operatorname{HeavyAtomCounts}(x)}}$$
(34)

For hyperparameters, we set the pocket embedding dimension to 128 and the training GFN temperature to Uniform(0, 64) which are used in TacoGFN. We trained the model with 40,000 oracles whereas TacoGFN is trained for 50,000 oracles.

1217 Introducing further objectives without retraining. For the restricted block library (TPSA<30),
1218 we set the action space subsampling ratio as 10% for both AddFirstReactant and ReactBi,
1219 and we set that as 1% for an entire library.

Scaling action space without retraining. For action space subsampling, we set the subsampling ratio as 2% for both AddFirstReactant and ReactBi for the "seen" and "unseen" libraries.
For the "all" library ("seen" + "unseen"), we set the subsampling ratio as 1%.

Ablation study. We set different subsampling ratios according to the building block library size.
For 100-sized, 1k-sized, and 10k-sized libraries, we do not perform the action space subsampling.
We set an action space subsampling ratio of 10% for a 100k-sized one and 1% for a 1M-sized one.

Further improvement with 3D interaction modeling. To better analyze the effects of interaction modeling, we limit the maximum reaction step to 1. Since the initial state is an empty graph, we can analyze the impact of interaction modeling on a single decision process for selecting a reaction to perform. To obtain the 3D binding conformation, we used GPU-accelerated UniDock (Yu et al., 2023) under the *fast* search mode. We used the same neural network structure for both 2D-based generation and 3D-based generation. For the reward setting, we used the Vina docking score as a reward function and filtered out molecules according to the Lipinski Rule.

1234 1235 C.4 SOFTWARES

1212 1213

Molecular docking software. For a fair comparison with the baseline model, we used UniDock (Yu et al., 2023) for target-specific generation and QuickVina 2.1 (Alhossary et al., 2015) for SBDD. The initial ligand conformer is generated with srETKDG3 (Wang et al., 2020) in RDKit (Landrum et al., 2013). For QuickVina, we converted the molecule format to pdbqt with OpenBabel (O'Boyle et al., 2011) and AutoDock Tools (Huey et al., 2012). To set up an exhaustive search, we set the search mode to *balance* for UniDock and the exhaustiveness to 8 for QuickVina. We kept the seed fixed at 1 throughout the ETKDG and the whole docking process.

1242 Docking proxy. We used the Quick Vina proxy proposed by Shen et al. (2023) which is implemented in PharmacoNet (Seo & Kim, 2023). We used a proxy model trained on the CrossDocked2020 training set rather than the model trained on the ZINCDock15M training set.

Synthetic accessibility estimation. To evaluate the synthetic accessibility of molecules, we used the retrosynthesis planning tool AiZynthFinder (Genheden et al., 2020). AiZynthFinder uses MCTS to find synthesis paths and estimate the number of steps, search time, success rate, and synthetically accessible score as metrics to indicate synthesis complexity or synthesizability.

1250 C.5 BASELINES

SynNet. For SynNet (Gao et al., 2022b), we perform multi-objective optimization with the following
 reward function:

$$R(x) = 0.5 \text{QED}(x) + 0.5 \widehat{\text{Vina}}(x). \tag{35}$$

According to the standard setting for optimization, we set the number of offspring to 512 and the number of oracles to 125. To use pre-trained models, we used SynNet's template set (91 templates) instead of a template set (71 templates) used in our work.

RGFN. Since the code of RGFN is not released, we reimplement the RGFN.

BBAR. Since BBAR (Seo et al., 2023) allows multi-conditional generation, we directly used QED and docking scores without any processing. We split 64,000 ZINC20 molecules according to the reported splitting of BBAR: 90% for the training set, 8% for the validation set, and 2% for the test set (in our case, the number of sampling molecules). We performed UniDock for training and validation set to prepare the label of the molecules. Since BBAR requires the desired property value, we used the average docking score of the top 100 diverse modes from our model.

Pocket2Mol, TargetDiff, DecompDiff, TacoGFN. We followed the reported generative setting to generate 100 molecules for each CrossDocked test pocket. We set the center of the pockets with the reference ligands in the CrossDocked2020 database. We reuse reported runtime in Shen et al. (2023), which is measured on NVIDIA A100 for Pocket2Mol, TargetDiff, and DecompDiff, and NVIDIA RTX3090 for TacoGFN.

1271 DiffSBDD, MolCRAFT. We used the generated samples from their official GitHub repository.

1273 C.6 LIT-PCBA POCKETS

1275Table 8 describes the protein information used in pocket-specific optimization with UniDock, which
is performed on Sec. 4.1.

1277 1278

1272

1274

Table 8: **The basic target information** of the LIT-PCBA dataset and PDB entry used in this work.

1279		DDD II	-
1280	Target	PDB Id	Target name
1281	ADRB2	4ldo	Beta2 adrenoceptor
1282	ALDH1	512m	Aldehyde dehydrogenase 1
1283	ESR_ago	2p15	Estrogen receptor α with agonist
1200	ESR_antago	2iok	Estrogen receptor α with antagonist
1204	FEN1	5fv7	FLAP Endonuclease 1
1285	GBA	2v3d	Acid Beta-Glucocerebrosidase
1286	IDH1	4umx	Isocitrate dehydrogenase 1
1287	KAT2A	5h86	Histone acetyltransferase KAT2A
1288	MAPK1	4zzn	Mitogen-activated protein kinase 1
1289	MTORC1	4dri	PPIase domain of FKBP51, Rapamycin
1290	OPRK1	6b73	Kappa opioid receptor
1291	PKM2	4jpg	Pyruvate kinase muscle isoform M1/M2
1292	PPARG	5y2t	Peroxisome proliferator-activated receptor γ
1293	TP53	3zme	Cellular tumor antigen p53
1294	VDR	3a2i	Vitamin D receptor
1295			

1296 D ADDITIONAL RESULTS

 D.1 ADDITIONAL RESULTS FOR POCKET-SPECIFIC GENERATION TASK

We reported the additional results of Sec. 4.1 for the remaining 10 pockets on the LIT-PCBA benchmark.

Table 9: **Hit ratio** (%). Average and standard deviation for 4 runs. The best results are in bold.

		Hit ratio (%, ↑)				
Category	Method	GBA	IDH1	KAT2A	MAPK1	MTOR
Fragment	FragGFN FragGFN+SA	$\begin{array}{c} 5.00 \ (\pm \ 4.24) \\ 3.00 \ (\pm \ 1.00) \end{array}$	$\begin{array}{c} 4.50 \ (\pm \ 1.66) \\ 4.50 \ (\pm \ 4.97) \end{array}$	$\begin{array}{c} 1.25 \ (\pm \ 0.83) \\ 1.50 \ (\pm \ 0.50) \end{array}$	$\begin{array}{c} 0.75 \ (\pm \ 0.83) \\ 2.00 \ (\pm \ 1.73) \end{array}$	$\begin{array}{c} 0.00 \ (\pm \\ 0.00 \ (\pm \end{array}) \end{array}$
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$\begin{array}{c} 50.00 \ (\pm \ 0.00) \\ 17.75 \ (\pm \ 2.28) \\ 58.00 \ (\pm \ 4.64) \\ 48.00 \ (\pm \ 1.22) \\ \textbf{66.00} \ (\pm \ 1.58) \end{array}$	$\begin{array}{c} 50.00 \ (\pm \ 0.00) \\ 19.50 \ (\pm \ 1.50) \\ 59.00 \ (\pm \ 4.06) \\ 43.00 \ (\pm \ 2.74) \\ \textbf{64.00} \ (\pm \ 5.05) \end{array}$	$\begin{array}{c} 29.17 (\pm 18.16) \\ 18.75 (\pm 1.92) \\ 55.50 (\pm 10.23) \\ 49.00 (\pm 1.22) \\ \textbf{66.50} (\pm 2.06) \end{array}$	$\begin{array}{c} 37.50 \ (\pm \ 21.65) \\ 16.25 \ (\pm \ 3.49) \\ 47.25 \ (\pm \ 6.61) \\ 38.00 \ (\pm \ 4.12) \\ \textbf{63.00} \ (\pm \ 4.64) \end{array}$	$\begin{array}{c} 0.00\ (\pm \\ \end{array})$
		OPRK1	PKM2	PPARG	TP53	VDF
Fragment	FragGFN FragGFN+SA	$\begin{array}{c} 0.50 \ (\pm \ 0.50) \\ 0.50 \ (\pm \ 0.87) \end{array}$	$\begin{array}{c} 7.25 \ (\pm \ 1.92) \\ 4.50 \ (\pm \ 1.50) \end{array}$	$\begin{array}{c} 0.75 \ (\pm \ 0.43) \\ 1.00 \ (\pm \ 0.71) \end{array}$	$\begin{array}{c} 4.25 \ (\pm \ 1.64) \\ 2.25 \ (\pm \ 1.92) \end{array}$	$\begin{array}{c} 0.00 \ (\pm \\ 0.00 \ (\pm \end{array}$
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$\begin{array}{c} 0.00 \ (\pm \ 0.00) \\ 2.50 \ (\pm \ 1.12) \\ 23.50 \ (\pm \ 5.94) \\ 2.50 \ (\pm \ 2.06) \\ \textbf{72.25} \ (\pm \ 2.05) \end{array}$	$\begin{array}{c} 0.00 \ (\pm \ 0.00) \\ 20.00 \ (\pm \ 0.71) \\ 50.75 \ (\pm \ 1.09) \\ 34.75 \ (\pm \ 6.57) \\ \textbf{62.00} \ (\pm \ 3.24) \end{array}$	$\begin{array}{c} 33.33 (\pm 20.41) \\ 10.50 (\pm 2.69) \\ 53.50 (\pm 5.68) \\ 29.00 (\pm 6.52) \\ \textbf{65.50} (\pm 4.03) \end{array}$	$\begin{array}{c} 8.33 \ (\pm 14.43) \\ 14.00 \ (\pm 3.94) \\ 55.50 \ (\pm 9.94) \\ 37.00 \ (\pm 6.60) \\ \textbf{67.50} \ (\pm 2.96) \end{array}$	$\begin{array}{c} 0.00\ (\pm \\ 0.00\ (\pm \\ 0.00\ (\pm \\ 0.00\ (\pm \\ 1.75\ (\pm \end{array}$

Table 10: Vina. Average and standard deviation for 4 runs. The best results are in bold.

		Average Vina Docking Score (kcal/mol, \downarrow)				
Category	Method	GBA	IDH1	KAT2A	MAPK1	MTORC1
Fragment	FragGFN FragGFN+SA	$\begin{array}{c} -8.76 \ (\pm \ 0.46) \\ -8.92 \ (\pm \ 0.27) \end{array}$	$\begin{array}{c} -9.91 \ (\pm \ 0.32) \\ -9.76 \ (\pm \ 0.64) \end{array}$	$\begin{array}{c} -9.27 \ (\pm \ 0.20) \\ -9.14 \ (\pm \ 0.43) \end{array}$	$\begin{array}{c} \textbf{-8.93} \ (\pm \ 0.18) \\ \textbf{-8.28} \ (\pm \ 0.40) \end{array}$	$-10.51 (\pm 0.3)$ -10.14 (± 0.3)
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$\begin{array}{c} -7.60 \ (\pm \ 0.09) \\ -8.70 \ (\pm \ 0.05) \\ -9.27 \ (\pm \ 0.06) \\ -8.48 \ (\pm \ 0.06) \\ \textbf{-9.62} \ (\pm \ 0.04) \end{array}$	$\begin{array}{c} -8.74 \ (\pm \ 0.08) \\ -9.84 \ (\pm \ 0.09) \\ -10.40 \ (\pm \ 0.08) \\ -9.49 \ (\pm \ 0.13) \\ \textbf{-10.95} \ (\pm \ 0.05) \end{array}$	$\begin{array}{c} -7.64 \ (\pm \ 0.38) \\ -8.54 \ (\pm \ 0.06) \\ -9.41 \ (\pm \ 0.04) \\ -8.53 \ (\pm \ 0.11) \\ \textbf{-9.73} \ (\pm \ 0.03) \end{array}$	$\begin{array}{c} -7.33 \ (\pm \ 0.14) \\ -8.49 \ (\pm \ 0.07) \\ -8.92 \ (\pm \ 0.05) \\ -8.22 \ (\pm \ 0.15) \\ \textbf{-9.30} \ (\pm \ 0.01) \end{array}$	$\begin{array}{c} -9.30 \ (\pm \ 0.4 \\ -10.07 \ (\pm \ 0.4 \\ -10.84 \ (\pm \ 0.4 \\ -9.89 \ (\pm \ 0.0 \\ \textbf{-11.39} \ (\pm \ 0.4 \ (\pm \ 0$
		OPRK1	PKM2	PPARG	TP53	VDR
Fragment	FragGFN FragGFN+SA	$\begin{array}{c} -10.28 \ (\pm \ 0.15) \\ -9.58 \ (\pm \ 0.44) \end{array}$	$\begin{array}{c} -11.24 \ (\pm \ 0.27) \\ -10.83 \ (\pm \ 0.34) \end{array}$	$\begin{array}{c} \textbf{-9.54} \ (\pm \ 0.12) \\ \textbf{-9.19} \ (\pm \ 0.29) \end{array}$	$\begin{array}{l} \textbf{-7.90} \ (\pm \ 0.02) \\ \textbf{-7.61} \ (\pm \ 0.27) \end{array}$	-10.96 (± 0. -10.66 (± 0.
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$\begin{array}{c} -8.70 \ (\pm \ 0.36) \\ -9.84 \ (\pm \ 0.10) \\ -10.34 \ (\pm \ 0.07) \\ -9.61 \ (\pm \ 0.11) \\ \textbf{-10.84} \ (\pm \ 0.03) \end{array}$	$\begin{array}{c} -9.55 \ (\pm \ 0.14) \\ -11.39 \ (\pm \ 0.08) \\ -11.98 \ (\pm \ 0.12) \\ -10.96 \ (\pm \ 0.18) \\ \textbf{-12.53} \ (\pm \ 0.02) \end{array}$	$\begin{array}{c} -7.47 \ (\pm \ 0.34) \\ -8.69 \ (\pm \ 0.10) \\ -9.40 \ (\pm \ 0.05) \\ -8.53 \ (\pm \ 0.07) \\ \textbf{-9.73} \ (\pm \ 0.02) \end{array}$	$\begin{array}{c} -5.34 \ (\pm \ 0.23) \\ -7.05 \ (\pm \ 0.09) \\ -7.90 \ (\pm \ 0.10) \\ -7.07 \ (\pm \ 0.06) \\ \textbf{-8.09} \ (\pm \ 0.06) \end{array}$	-10.98 (± 0. -11.07 (± 0. -11.62 (± 0. -10.86 (± 0. -12.30 (± 0.

	Synthesizabilit	y. Average an	u stanuaru uev	1au011 101 4 1u11	s. The best lest	ints are in boi	
			AiZynthFinder Success Rate (%, ↑)				
Category	Method	GBA	IDH1	KAT2A	MAPK1	MTORC1	
Fragment	FragGFN FragGFN+SA	$\begin{array}{c} 5.00 \ (\pm \ 4.24) \\ 3.00 \ (\pm \ 1.00) \end{array}$	$\begin{array}{c} 4.50 \ (\pm \ 1.66) \\ 4.50 \ (\pm \ 4.97) \end{array}$	$\begin{array}{c} 1.25 \ (\pm \ 0.83) \\ 1.50 \ (\pm \ 0.50) \end{array}$	$\begin{array}{c} 0.75 \ (\pm \ 0.83) \\ 3.25 \ (\pm \ 1.48) \end{array}$	$\begin{array}{c} 2.75 \ (\pm \ 1.30) \\ 3.50 \ (\pm \ 2.50) \end{array}$	
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$\begin{array}{c} 50.00 \ (\pm \ 0.00) \\ 17.75 \ (\pm \ 2.28) \\ 58.00 \ (\pm \ 4.64) \\ 48.00 \ (\pm \ 1.22) \\ \textbf{66.00} \ (\pm \ 1.58) \end{array}$	$\begin{array}{c} 50.00 \ (\pm \ 0.00) \\ 19.50 \ (\pm \ 1.50) \\ 59.00 \ (\pm \ 4.06) \\ 43.00 \ (\pm \ 2.74) \\ \textbf{64.00} \ (\pm \ 5.05) \end{array}$	$\begin{array}{c} 45.83 \ (\pm\ 27.32) \\ 18.75 \ (\pm\ 1.92) \\ 55.50 \ (\pm\ 10.23) \\ 49.00 \ (\pm\ 1.22) \\ \textbf{66.50} \ (\pm\ 2.06) \end{array}$	$\begin{array}{c} 50.00 \ (\pm \ 0.00) \\ 16.25 \ (\pm \ 3.49) \\ 47.25 \ (\pm \ 6.61) \\ 42.00 \ (\pm \ 3.00) \\ \textbf{63.00} \ (\pm \ 4.64) \end{array}$	$\begin{array}{c} 54.17 (\pm 7.22 \\ 18.75 (\pm 3.90 \\ 57.00 (\pm 7.58 \\ 44.50 (\pm 4.03 \\ \textbf{70.50} (\pm 2.87 \\ \end{array} \right.$	
		OPRK1	PKM2	PPARG	TP53	VDR	
Fragment	FragGFN FragGFN+SA	2.50 (± 2.29) 3.25 (± 1.79)	8.75 (± 3.11) 9.75 (± 2.28)	$\begin{array}{c} 0.75 \ (\pm \ 0.43) \\ 1.25 \ (\pm \ 1.09) \end{array}$	4.25 (± 1.64) 2.25 (± 1.92)	3.50 (± 2.18) 3.75 (± 2.77)	
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$54.17 (\pm 7.22) \\ 13.75 (\pm 3.11) \\ 56.50 (\pm 7.63) \\ 48.00 (\pm 2.55) \\ \textbf{72.25} (\pm 2.05) \\ \end{array}$	$\begin{array}{c} 50.00 \ (\pm \ 0.00) \\ 20.00 \ (\pm \ 0.71) \\ 50.75 \ (\pm \ 1.09) \\ 48.50 \ (\pm \ 3.20) \\ 62.00 \ (\pm \ 3.24) \end{array}$	$54.17 (\pm 7.22) \\ 15.50 (\pm 2.29) \\ 53.50 (\pm 5.68) \\ 47.00 (\pm 5.83) \\ \textbf{65.50} (\pm 4.03) \\ \end{cases}$	$\begin{array}{c} 29.17 (\pm 18.16) \\ 18.50 (\pm 3.28) \\ 55.50 (\pm 9.94) \\ 53.25 (\pm 3.63) \\ \textbf{67.50} (\pm 2.96) \end{array}$	$\begin{array}{c} 45.83 (\pm 7.22 \\ 12.25 (\pm 3.34 \\ 53.50 (\pm 1.80 \\ 46.50 (\pm 2.69 \\ \textbf{66.75} (\pm 2.28 \\ \end{array}$	

Table 11: Synthesizability. Average and standard deviation for 4 runs. The best results are in bold.

Table 12: Synthetic complexity. Average and standard deviation for 4 runs. The best results are in bold.

		Average Number of Synthesis Steps (\downarrow)				
Category	Method	GBA	IDH1	KAT2A	MAPK1	MTORC1
Fragment	FragGFN FragGFN+SA	$\begin{array}{c} 3.94 \ (\pm \ 0.11) \\ 3.94 \ (\pm \ 0.15) \end{array}$	$\begin{array}{c} 3.74 \ (\pm \ 0.10) \\ 3.84 \ (\pm \ 0.23) \end{array}$	$\begin{array}{c} 3.78 \ (\pm \ 0.09) \\ 3.66 \ (\pm \ 0.18) \end{array}$	$\begin{array}{c} 3.72 \ (\pm \ 0.18) \\ 3.69 \ (\pm \ 0.21) \end{array}$	$\begin{array}{c} 3.84 \ (\pm \ 0.18) \\ 3.94 \ (\pm \ 0.08) \end{array}$
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$\begin{array}{c} 3.38 \ (\pm \ 0.22) \\ 3.71 \ (\pm \ 0.12) \\ 2.48 \ (\pm \ 0.18) \\ 2.77 \ (\pm \ 0.20) \\ \textbf{2.10} \ (\pm \ 0.08) \end{array}$	$\begin{array}{c} 3.38 \ (\pm \ 0.22) \\ 3.68 \ (\pm \ 0.02) \\ 2.61 \ (\pm \ 0.13) \\ 2.97 \ (\pm \ 0.15) \\ \textbf{2.16} \ (\pm \ 0.11) \end{array}$	$\begin{array}{c} 3.46 \ (\pm \ 0.95) \\ 3.63 \ (\pm \ 0.05) \\ 2.45 \ (\pm \ 0.37) \\ 2.78 \ (\pm \ 0.10) \\ \textbf{2.29} \ (\pm \ 0.05) \end{array}$	$\begin{array}{c} 3.50 \ (\pm \ 0.00) \\ 3.73 \ (\pm \ 0.05) \\ 2.81 \ (\pm \ 0.24) \\ 2.86 \ (\pm \ 0.19) \\ \textbf{2.29} \ (\pm \ 0.11) \end{array}$	$\begin{array}{c} 3.29 \ (\pm \ 0.36) \\ 3.77 \ (\pm \ 0.09) \\ 2.44 \ (\pm \ 0.27) \\ 2.92 \ (\pm \ 0.06) \\ \textbf{2.05} \ (\pm \ 0.09) \end{array}$
		OPRK1	PKM2	PPARG	TP53	VDR
Fragment	FragGFN FragGFN+SA	$\begin{array}{c} 3.82 \ (\pm \ 0.13) \\ 3.62 \ (\pm \ 0.12) \end{array}$	$\begin{array}{c} 3.71 \ (\pm \ 0.12) \\ 3.84 \ (\pm \ 0.21) \end{array}$	$\begin{array}{c} 3.73 \ (\pm \ 0.24) \\ 3.71 \ (\pm \ 0.04) \end{array}$	$\begin{array}{c} 3.73 \ (\pm \ 0.23) \\ 3.66 \ (\pm \ 0.05) \end{array}$	$\begin{array}{c} 3.75 \ (\pm \ 0.06) \\ 3.67 \ (\pm \ 0.25) \end{array}$
Reaction	SynNet BBAR SynFlowNet RGFN RXNFLOW	$\begin{array}{c} 3.29 \ (\pm \ 0.36) \\ 3.70 \ (\pm \ 0.17) \\ 2.49 \ (\pm \ 0.33) \\ 2.81 \ (\pm \ 0.12) \\ \textbf{2.00} \ (\pm \ 0.09) \end{array}$	$\begin{array}{c} 3.50 \ (\pm \ 0.00) \\ 3.61 \ (\pm \ 0.05) \\ 2.62 \ (\pm \ 0.10) \\ 2.82 \ (\pm \ 0.10) \\ \textbf{2.34} \ (\pm \ 0.19) \end{array}$	$\begin{array}{c} 3.29 \ (\pm \ 0.36) \\ 3.72 \ (\pm \ 0.13) \\ 2.56 \ (\pm \ 0.12) \\ 2.82 \ (\pm \ 0.18) \\ \textbf{2.21} \ (\pm \ 0.06) \end{array}$	$\begin{array}{c} 3.67 \ (\pm \ 0.91) \\ 3.65 \ (\pm \ 0.05) \\ 2.51 \ (\pm \ 0.27) \\ 2.64 \ (\pm \ 0.10) \\ \textbf{2.12} \ (\pm \ 0.12) \end{array}$	$\begin{array}{c} 3.63 \ (\pm \ 0.22) \\ 3.77 \ (\pm \ 0.16) \\ 2.55 \ (\pm \ 0.09) \\ 2.84 \ (\pm \ 0.18) \\ \textbf{2.12} \ (\pm \ 0.12) \end{array}$

1404 D.2 PROPERTY DISTRIBUTION FOR POCKET-SPECIFIC GENERATION TASK







1456 1457



Figure 10: The property distribution of the generated samples for the last 5 LIT-PCBA targets.

SCALING LAWS WITH BASELINE GFLOWNETS D.3

In this section, we investigate the scaling laws of our model and the baseline GFlowNets, Syn-FlowNet and RGFN, focusing on performance (Figure 11) and computational cost (Figure 12). Our action space subsampling method reduces the computational cost and memory consumption via cost-variance trade-off and memory-variance trade-off.



Figure 11: **Optimization power and diversity**. Average of standard deviation over the 4 runs. (a) Average docking score. (b) The uniqueness of Bemis-Murcko scaffolds. (c) Average Tanimoto distance.



Figure 12: **Runtime** according to the building block library size. Average of standard deviation over the 100 batches. (a) The training runtime. (b) The sampling runtime without model training.

Performance. We conducted an optimization for the kappa-opioid receptor using RXNFLOW and baseline GFlowNets. To prevent the hacking of docking by increasing the molecular size, we performed Vina-QED multi-objective using the multi-objective GFlowNet framework (Jain et al., 2023). As shown in Figure 11, all three models exhibit similar trends in docking score optimization. However, conventional GFlowNets face memory resource constraints that scale with the action space size. As a result, SynFlowNet and RGFN were restricted to library sizes up to 100,000 and 500,000 sizes, respectively, due to memory limitations. In contrast, RXNFLOW can accommodate larger action spaces leveraging the memory-variance trade-off.

Speed. While SynFlowNet and RGFN consider all actions in the massive action space to estimate the forward transition probability, our architecture estimate Our architecture estimates forward tran-sition probabilities over a subset of the action space, unlike SynFlowNet and RGFN, which compute all flows for the entire action space. This design allows RXNFLOW to handle larger action spaces without encountering computational bottlenecks. To evaluate the cost-efficiency of our technique, we measured runtime during both training and generation. To unify the generation environments across different methods, we used a constant reward function R(x) = 1 and a maximum reaction step of 1. Training times were measured on random trajectories, and sampling times were measured with initialized models, according to the building block library size. Due to memory limitations, SynFlowNet was restricted to library sizes up to 400,000.

For RXNFLOW, we tested two configurations: RXNFLOW (1k) and RXNFLOW (10k), which sample up to 1,000 and 10,000 building blocks, respectively. As demonstrated in Figure 12, RXNFLOW achieves significantly better cost efficiency in both training and generation compared to the base-line models. Additionally, computational costs can be easily reduced by adjusting the action-space subsampling ratio.

D.4 STATISTICAL INFORMATION FOR POCKET-CONDITIONAL GENERATION TASK

1568	Table 13. Statistical Information for nocket-conditional generative models. Mean and standar
1569	deviation for 5 sample sets
1570	

		Vina	ı (↓)	QEI) (†)
Category	Model	Avg.	Std.	Avg.	Std.
	Pocket2Mol	-7.603	0.087	0.567	0.007
	TargetDiff	-7.367	0.028	0.487	0.006
Atom	DiffSBDD	-6.949	0.079	0.467	0.002
	DecompDiff	-8.350	0.033	0.368	0.004
	MolCRAFT	-8.053	0.033	0.500	0.004
	MolCRAFT-large	-9.302	0.033	0.448	0.002
Fragment	TacoGFN	-8.237	0.268	0.671	0.002
Reaction	RXNFLOW	-8.851	0.031	0.666	0.001

TARGET SPECIFICITY OF GENERATED SAMPLES D.5

To investigate the target specificity of pocket-conditional generation, we measured delta score (Gao et al., 2024a) for the top-10 molecules for each pocket. The delta score evaluates the pocket speci-ficity of a proposed molecule by comparing the docking scores difference in how well each molecule binds to other proteins compared to the target protein.

Table 14: Delta Score for each methods.					
Category	Model	Delta Score			
Atom	DecompDiff	-1.29			
Fragment	TacoGFN	-1.13			
Reaction	RXNFLOW	-1.13			



ABLATION STUDY FOR NON-HIERARCHICAL MARKOV DECISION PROCESS STRUCTURE





k

To investigate the effectiveness of the non-hierarchical MDP structure, we performed an ablation study. We randomly selected 100,000 general building blocks from the Enamine building block library for model training (general). To simulate a scenario where specific functional groups are unallowed, we filtered out all aromatic building blocks to create a nonaromatic building block set.
We trained the GFlowNets under Vina-QED multi-objective settings (Jain et al., 2023) against the beta-2 adrenergic receptor for 1,000 training oracles. After training, we generated 100 molecules using both the general block set and the nonaromatic block set without additional training.

As shown in Figure 13, the proposed non-hierarchical MDPs closely align with the identified reward distributions for both objectives, Vina and QED, on the general and nonaromatic building block sets. In contrast, hierarchical MDPs, as utilized in existing methodologies, demonstrate a shift in the reward distribution when the building block set is restricted. This indicates that non-hierarchical MDPs are more robust in changes in the building block set compared to hierarchical MDPs.

1633 D.7 Additional results for scaling action space without retraining



We reported the additional results for additional reward exponent settings (R^{β}) .

Figure 14: **QED reward distribution** of generated molecules.

Moreover, we investigated the generalization ability of our method to structurally different building blocks under the Vina-QED objectives against the beta-2 adrenergic receptor. For model training, we randomly selected 100,000 building blocks from the entire library (seen). Additionally, we selected 100,000 building blocks with a Tanimoto similarity of less than 0.5 to all training building blocks (unseen). The model was trained using the seen blocks first, and then the trained model subsequently generate molecules with seen blocks, unseen blocks, and a combination of both sets (all), respectively. As illustrated in Figure 15, the model can also generate samples with similar reward distributions from building blocks with different distributions than the ones used for training.



1674 D.8 THEORETICAL ANALYSIS

1676 To assess the impact of action space subsampling in GFlowNet training, we conduct a toy experiment using a simplified setup with 10,000 blocks, one uni-molecular reaction template, one bi-molecular 1677 reaction template, and the QED objective. We used the Hell-Volhard-Zelinsky reaction as a uni-1678 molecular reaction and the Amide reaction as a bi-molecular reaction, which are illustrated in Fig. 1679 16. We used a minimum trajectory length of 1, max trajectory length of 2, constant GFN temperature 1680 of 1.0, and learning rate decay of 3,000 for P_F and $\log Z$. For the GFlowNet sampler, we used the 1681 same weights of the proxy model, i.e. EMA factor of 0. We performed optimization for 30,000 1682 oracles with a batch size of 64. 1683

As shown in Figures 17(a) and 17(b), we compare a baseline GFlowNet trained without subsam-1684 pling ("base") to models using various subsampling ratios and Monte Carlo (MC) sampling. The 1685 differences in log Z_{θ} are relatively small (<0.005) across all settings, and increasing MC samples 1686 for state flow estimation F_{θ} further reduced the bias. In Figures 17(c) and 17(d), we also evaluate 1687 the bias in the trajectory balance loss $(\hat{\mathcal{L}}_{TB})$ and its gradient norm $(\|\nabla_{\theta}\hat{\mathcal{L}}_{TB}\|)$ during training, find-1688 ing negligible differences compared to the true values. These results indicate that our importance 1689 sampling reweighting approach effectively mitigates bias from action space subsampling, enabling 1690 efficient and accurate policy estimation. 1691



Figure 16: Reaction templates employed in toy experiments. (a) Hell-Volhard-Zelinsky reaction.
(b) Amide reaction.



Figure 17: **Bias estimation.** (a) $\log Z_{\theta}$ according to the action space subsampling ratio (left) and the number of MC samples where the subsampling ratio is 1/9 (right). (b) The trajectory balance loss ($\mathcal{L}_{TB}, \hat{\mathcal{L}}_{TB}$) where the subsampling ratio is 1/9 under 4 MC samples. (c) The loss gradient norms ($\|\nabla_{\theta} \mathcal{L}_{TB}\|, \|\nabla_{\theta} \hat{\mathcal{L}}_{TB}\|$) where the subsampling ratio is 1/9 under 4 MC samples.

1720

1693 1694 1695

1697 1698

1699

1700

1701

1702

- 1721 1722
- 1723
- 172/
- 1725
- 1726
- 1727