

CONTROLLABLE MOLECULE GENERATION BY SAMPLING IN CONTINUOUS PARAMETER SPACE

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep generative models have made significant strides in continuous data generation, such as producing realistic images and 3D protein conformations. However, due to the sensitivity of topological graphs to noise and the constraints of long-range discrete relationships, the generation of purely discrete data—such as topological graphs—remains a long-standing challenge, with property control proving even more elusive. In this paper, we propose a novel molecular graph generative framework, called CtrlMol, to learn the topological graphs of molecules in a differentiable parameter space. Unlike diffusion models that iteratively refine samples, CtrlMol optimizes distribution parameters at different noise levels through a pre-defined Bayesian flow. At each sampling step, we leverage a property-guided output distribution to have fine-grained control of topological graphs toward the given property. Experimental results demonstrate that CtrlMol outperforms all the competing baselines in generating natural molecule graphs. In addition, CtrlMol advances the state of the art in producing the molecules with the desired properties.

1 INTRODUCTION

The molecular graph generation aims to generate valid and realistic molecules with desired properties, which is a fundamental task in drug discovery. However, this problem is very challenging. On the one hand, the sample space is discrete by nature, where the data distribution is rugged and not easy to learn based on continuous assumptions. On the other hand, the molecular graphs are primarily multi-modality, including the nodes (i.e., atoms) and the edges (i.e., bonds), which are jointly dependent but follow disparate distributions. Moreover, the edges in the molecular graph are not only discrete and sparse but also obey multiple constraints in chemical valence. These distinct characteristics of the different modalities further lift the complexity of molecular graph generation.

Traditional generative approaches transform the topological molecular structures into sequences (e.g., SMILES) and formulate molecular graph generation as the sequence generation. In this way, sequential models such as LSTM and Transformer are used to produce the molecule graph in an auto-regressive manner. For example, Gómez-Bombarelli et al. (2018) learn the molecular SMILES strings with a variational auto-encoder. But the lack of explicitly syntactic constraints for molecular graphs makes them inclined to generate invalid sequences. Dai et al. (2018) and Kusner et al. (2017) impose task-specific grammars into the sequence decoder to enhance the validity of the generated sequences. As a cost of introducing the syntactic constraints, the learning of inconsecutive sequences becomes challenging.

As another alternative, graph generative networks directly output molecule graphs by iteratively adding atoms and chemical bonds. In particular, GraphAF defines a feed-forward neural network from molecular graph structures to the base distribution and generates the nodes and edges based on existing sub-graphs (Shi et al., 2020). GraphDF generates molecular graphs by sequentially sampling discrete latent variables and mapping them to new nodes and edges (Luo et al., 2021). Although graph generative models achieve impressive advancement, the graph generative process is more complicated than the sequential one due to the high degree of freedom in generative ordering. With the success of the diffusion model (DM) in producing realistic continuous images, DM is also applied to molecule generation by denoising the token embeddings in the form of molecular sequences. Runcie & Mey (2023) present a selective iterative latent variable refinement method for conditioning an existing pre-trained equivariant diffusion model toward the molecule graph generation. However,

054 the simplified continuity assumptions of data distribution in DM (e.g., Gaussian distribution) do not
055 always hold in discrete data, which deepens the difficulties of multi-modality issues.

056
057 In this paper, we introduce a novel framework, CtrlMol to address the multi-modality issue in
058 molecular graph generation. The key idea is to establish a conditional Bayes flow network with
059 different base distributions for the sampling of the multi-modality. Concretely, CtrlMol adopts a
060 unified probabilistic modeling formulation for different modalities in the molecular graph, including
061 the atom identity (i.e., nodes) and the chemical bonds (i.e., edges). Second, we introduce a topo-
062 logical complete edge sampling strategy to address the huge complexity $O(N^2)$ in edge generation.
063 Further, CtrlMol employs the property-guided output distribution to have fine-grained control of
064 the topological structures. Finally, CtrlMol optimizes distribution parameters by minimizing the
065 discrepancy between the output distribution and the true data distribution.

066 We evaluate CtrlMol on both unconditional and conditional molecule generation tasks. Experimental
067 results demonstrate that CtrlMol outperforms all competing baselines in generating realistic molecular
068 graphs. For conditional generation, CtrlMol achieves an improvement of 26% and 64% over the
069 previous best generators in terms of QED and LogP properties, respectively. Moreover, CtrlMol
070 could sample with any number of steps by virtue of the continuous parameter sampling space, which
071 leads to a 100 \times speedup with SOTA performance on conditional molecule generation.

072 2 RELATED WORK

073
074 Early methods typically adopt SMILES/SMARTS strings as a text-based representation for atom-by-
075 atom molecule generation. By decoding the molecular strings token by token, the generative model
076 can learn prior knowledge that encapsulates the grammar and syntax of valid SMILES (Jensen, 2019;
077 Segler et al., 2018; Gómez-Bombarelli et al., 2018). With the development of graph neural networks
078 (GNNs), topological graphs emerge as a more prevailing representation of molecules. For example,
079 JTVAE (Jin et al., 2018b), GCPN (You et al., 2018), GraphAF (Shi et al., 2020), and GraphDF (Luo
080 et al., 2021) proposed several autoregressive models to generate new nodes and edges based on the
081 generated molecular graph.

082 Apart from the molecular strings (1D) and topological graphs (2D), the molecule can be also
083 represented as the atom types and the corresponding coordinates (3D). For the generation of molecular
084 3D structures, diffusion models have been widely adopted. DiffSBDD (Schneuing et al., 2022) and
085 TargetDiff (Guan et al., 2022) learn the distribution of atom types and positions from a standard
086 Gaussian prior based via the diffusion process. DecompDiff (Guan et al., 2023) decomposes the
087 ligand molecule and the prior into arms and scaffolds and then leverages the reverse process of
088 diffusion to generate molecules. In fact, these diffusion-based generative methods leverage the
089 continuous 3D coordinates to avoid the generation of the chemical bonds of atoms, which are sparse
090 and discrete. Considering that the 3D structures of molecules are usually not with high-quality (e.g.,
091 approximated by RDKit library) and available, the generation of atoms and coordinates can hardly
092 serve as a universal molecular graph generator, especially in the scenarios lack of 3D structures (e.g.,
093 generation toward in vivo and in vitro properties).

094 Recently, Bayes flow networks (BFNs) have been introduced to generate discrete structures in
095 continuous parameter space (Qu et al., 2024; Song et al., 2023). Similar to diffusion generators,
096 GeoBFN decomposes the molecular structure into atoms and the coordinates. Then GeoBFN adopts
097 BFN to generate both of them and transforms the atom coordinates to build the final molecular
098 graphs (Song et al., 2023). Different from GeoBFN, Our approach focuses on the 2D molecular
099 graph and proposes several novel components (such as the topological complete edge sampler and
100 property-guided output distribution) to address the multi-modality and edge sparsity issues.

101 3 METHODOLOGY

102 3.1 PROBLEM DEFINITION

103
104 We focus on the 2D molecular graphs due to the abundance of available data. The geometry of a
105 2D molecular graph is denoted as $\mathcal{M} = \{(\mathbf{V}, \mathbf{E})\}$, where $\mathbf{V} \in \mathbb{R}^{|\mathbf{V}| \times K_V}$ is the set of node (i.e.,
106 atoms) features and $\mathbf{E} \in \mathbb{R}^{|\mathbf{V}| \times |\mathbf{V}| \times K_E}$ is the set of edge (i.e., bonds) features. Each node and edge
107

has categorical attributes corresponding to their types. The i -th node is represented by a one-hot vector $\mathbf{v}^i \in \mathbf{V}$ with K_V dimension, where K_V represents the total number of atomic types. Similarly, the j -th edge $e^j \in \mathbf{E}$ is represented as a one-hot vector of dimension K_E , with K_E denoting the total number of bond types. Without loss of generality, there is an edge type to stand for no connection between two nodes. The primary goal of conditional molecular generation is to produce valid molecular structures \mathcal{M} that satisfy specific properties c (e.g., logP, QED (Bickerton et al., 2012), etc.). By conditioning the generation process on these properties, we aim to explore the vast chemical space effectively, facilitating the discovery of new compounds with desired characteristics.

3.2 THE OVERVIEW OF BAYESIAN FLOW NETWORKS

To effectively address the discrete nature of molecular topology, our CtrlMol model generates molecules by sampling in a continuous parameter space using the Bayesian Flow Networks (BFN) framework. The key distinction between BFN and diffusion models is that BFN refines the parameters of the data distribution rather than operating on noisy data as diffusion models do.

As illustrated in Figure 1, the generative process can be conceptualized as a message exchange between a sender and a receiver. The receiver begins with prior knowledge of the data distribution, referred to as the “input distribution” ($p_I(\mathcal{M} | \boldsymbol{\theta})$), where $\boldsymbol{\theta}$ is the parameters of the prior distribution for the data \mathcal{M} . Additionally, the receiver can estimate an “output distribution” ($p_O(\hat{\mathcal{M}} | \Phi(\boldsymbol{\theta}, t))$) based on the parameters $\boldsymbol{\theta}$ of the input distribution along with the processing time t , using a neural network to predict the data $\hat{\mathcal{M}}$.

During each message-passing step, the sender introduces noise (controlled by an accuracy parameter α) to the origin data, creating a “sender distribution” ($p_S(\mathbf{y} | \mathcal{M}; \alpha)$), and sending a sampled data point \mathbf{y} to the receiver. The receiver then constructs a “receiver distribution” ($p_R(\mathbf{y} | \boldsymbol{\theta}; t, \alpha)$) by adding noise to the estimated output distribution, generating the sender sample \mathbf{y} such that

$$p_R(\mathbf{y} | \boldsymbol{\theta}; t, \alpha) = \mathbb{E}_{\hat{\mathcal{M}} \sim p_O(\hat{\mathcal{M}} | \Phi(\boldsymbol{\theta}, t))} p_S(\mathbf{y} | \hat{\mathcal{M}}; \alpha). \quad (1)$$

When the sender transmits a sample from p_S to the receiver, the cost incurred is the KL divergence between receiver p_R and sender p_S . Subsequently, the receiver updates the parameters of its input distribution p_I based on the received sample \mathbf{y} by Bayesian inference, which can be expressed in closed form if the input distribution of each variable in the data is independent. Here we denote the Bayesian update function as h , which applies the rules of Bayesian inference to compute the updated parameters $\boldsymbol{\theta}'$, specifically, $\boldsymbol{\theta}' \leftarrow h(\boldsymbol{\theta}, \mathbf{y}, \alpha)$. Then the Bayesian update distribution is obtained by marginalizing out \mathbf{y} as follows:

$$p_U(\boldsymbol{\theta}' | \boldsymbol{\theta}, \mathcal{M}; \alpha) = \mathbb{E}_{\mathbf{y} \sim p_S(\mathbf{y} | \mathcal{M}; \alpha)} \delta(\boldsymbol{\theta}' - h(\boldsymbol{\theta}, \mathbf{y}, \alpha)), \quad (2)$$

where δ is the Dirac delta distribution. The accuracy α has the additive property (as proven by Gra), that is, $p_U(\boldsymbol{\theta}'' | \boldsymbol{\theta}, \mathcal{M}; \alpha_a + \alpha_b) = \mathbb{E}_{\boldsymbol{\theta}' \sim p_U(\boldsymbol{\theta}' | \boldsymbol{\theta}, \mathcal{M}; \alpha_a)} p_U(\boldsymbol{\theta}'' | \boldsymbol{\theta}', \mathcal{M}; \alpha_b)$. Therefore the accuracy schedule $\beta(t)$ is introduced by accumulating α , that is

$$\beta(t) = \int_{t'=0}^t \alpha(t') dt'. \quad (3)$$

Then the Bayesian flow distribution can be derived by marginalizing distribution over input parameters at time t :

$$p_F(\boldsymbol{\theta} | \mathcal{M}; t) = p_U(\boldsymbol{\theta} | \boldsymbol{\theta}_0, \mathcal{M}; \beta(t)) \quad (4)$$

Considering a sequence of n steps where n sender samples $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$ is sampled at times t_1, t_2, \dots, t_n . At time step t_i , the sender distribution is $p_S(\cdot | \mathcal{M}; \alpha_i)$, where $\alpha_i = \beta(t_i) - \beta(t_{i-1})$. The receiver distribution is $p_R(\cdot | \boldsymbol{\theta}_{i-1}; t_{i-1}, \alpha_i)$, which is determined by the $\boldsymbol{\theta}_{i-1}$ and t_{i-1} since the receiver has not received the message to update its input parameters. The input parameter sequence $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_n$ is updated from Bayesian update function h : $\boldsymbol{\theta}_i = h(\boldsymbol{\theta}_{i-1}, \mathcal{M}, \alpha(i))$. Then the n -step discrete-time loss $L^n(\mathcal{M})$ can be calculated by the KL divergence between the sender and receiver distributions as follows:

$$L^n(\mathcal{M}) = \mathbb{E}_{p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{KL}(p_S(\cdot | \mathcal{M}; \alpha_i) || p_R(\cdot | \boldsymbol{\theta}_{i-1}; t_{i-1}, \alpha_i)) \quad (5)$$

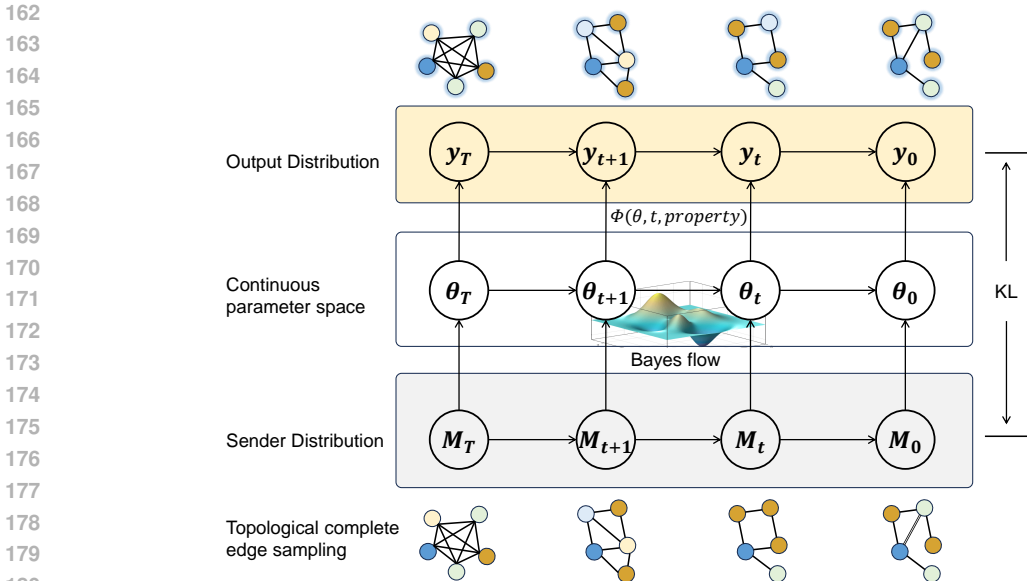


Figure 1: The generative process of the CtrlMol framework.

The summation term can be rewritten as:

$$\sum_{i=1}^n D_{KL} = n \times \frac{1}{n} \sum_{i=1}^n D_{KL} = n \times \sum_{i=1}^n \frac{1}{n} D_{KL} = n \mathbb{E}_{i \sim U\{1, n\}} D_{KL}. \quad (6)$$

where we abbreviate KL divergence as D_{KL} and $U\{1, n\}$ is the uniform distribution over the integers from 1 to n . Then the discrete-time loss (Eq 5) is reformulated as

$$\begin{aligned} L^n(\mathcal{M}) &= n \mathbb{E}_{i \sim U\{1, n\}} \mathbb{E}_{p_U(\theta_1 | \theta_0, \mathcal{M}; \alpha_1)} \mathbb{E}_{p_U(\theta_2 | \theta_1, \mathcal{M}; \alpha_2)} \dots \mathbb{E}_{p_U(\theta_{i-1} | \theta_{i-2}, \mathcal{M}; \alpha_{i-1})} D_{KL} \\ &= n \mathbb{E}_{i \sim U\{1, n\}, p_F(\theta | \mathcal{M}; t_{i-1})} D_{KL}(p_S(\cdot | \mathcal{M}; \alpha_i) || p_R(\cdot | \theta_{i-1}; t_{i-1}, \alpha_i)) \end{aligned} \quad (7)$$

which allows us approximate $L^n(\mathcal{M})$ via Monte-Carlo sampling without summation over n steps.

3.3 THE CTRLMOL MODEL

In this section, we will elaborate on the details of different components of CtrlMol within the field of discrete molecular graph data. Furthermore, We introduce a property-guided output distribution (as described in section 3.3.2) that allows for finer-grained control over the generation of topological structures, aiding in the production of molecules with specific properties. Since both atoms and edges in the molecules are discrete data and are represented similarly, we will primarily use nodes as examples to illustrate our model. Any differences related to edges will be explicitly noted.

3.3.1 INPUT DISTRIBUTION

Input distribution of nodes. Given a molecular graph $\mathcal{M} = \{(\mathbf{V}, \mathbf{E})\}$, both nodes (atoms) and edges (bonds) are categorical (discrete) variables defined by their type indices. The input distribution characterized for them can be represented as a categorical distribution. It is important to note that the input distribution for each variable in a data point is independent. Take atoms \mathbf{V} as an example, for an atom (a variable) $\mathbf{v}^{(i)} \in \mathbf{V}$, the categorical distribution with K_V categories is denoted by a K_V -dimensional vector, $\theta^{(i)} = (\theta_1^{(i)}, \theta_2^{(i)}, \dots, \theta_{K_V}^{(i)}) \in [0, 1]^{K_V}$, where each entry $\theta_{k_v}^{(i)}$ corresponds to the probability assigned to class k_v for the variable $\mathbf{v}^{(i)}$. The input distribution of atoms \mathbf{V} is then expressed as the joint probability over the true categories of all N_v variables $\mathbf{v}^{(i)}$, i.e.,

$$p_I(\mathbf{V} | \theta) = \prod_{i=1}^{N_v} \theta_{v^{(i)}}^{(i)}. \quad (8)$$

The parameters of input distribution form a matrix $\theta = (\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(N_v)}) \in [0, 1]^{K_v \times N_v}$. For prior θ_0 , each entry is set to $\frac{1}{K_v}$ according to a uniform distribution.

Input distribution of edges. Similar to the node distribution, the input distribution of edges is also a categorical uniform distribution over all the edge types. Note that, the edge types and connection states can be changed during the sampling process. Therefore, we consider the possibility of no edge connection between two nodes and include a new edge type named no connections in addition to single bond, double bond, etc.

However, the parallel generation of edges $\mathbf{E} \in \mathbb{R}^{N_v \times N_v \times K_E}$ should be sampled from the full connections¹. It requires the complexity of $O(N_v^2)$ for each of the molecular graphs, requiring huge computational resources. While the total number of edges in true data distribution is far smaller than full connections. In addition, the input distribution is the true distribution for the learning of the output distribution. With the complexity of $O(N_v^2)$, the output distribution should describe the probability for each entry of the full connection, leading to computational redundancy for the edge prediction.

Topological complete edge sampling strategy. As demonstrated in Theorem 1, the connections of a K -regular graph can cover any arbitrary graph with the node degree at most K , regardless of the node identity. In other words, if we want to generate a given topological graph with N_v nodes and the degree at most K , the K -regular graph with N_v nodes is complete to be the initial topological graph for sampling. Considering that the node in molecular graphs with the maximum connections is the carbon atom, we set the maximum number of the edge connections in the initial graph as $K = 4$. Formally speaking, we will randomly add several edges with "no connection" type into the given molecular graph $\mathcal{M} = \{(\mathbf{V}, \mathbf{E})\}$ to satisfy that the degree equals K and then build the input distribution of edges. Note that the above topological edge sampling strategy is guaranteed by Theorem 1 and largely reduces the computational complexity to $O(KN_v)$, where K is a constant.

Theorem 1. For arbitrary positive integer N and D ($D < N$), undirected graph $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ that has $|\mathcal{V}| = N$ and $\deg(v) = D, v \in \mathcal{V}$, undirected graph $\mathcal{H} = \langle \mathcal{V}^H, \mathcal{E}^H \rangle$ that has $|\mathcal{V}^H| = N$ and $\deg(v) \leq D, v \in \mathcal{V}^H$, there exists at least one subgraph $\mathcal{G}' \subseteq \mathcal{G}$ such that \mathcal{H} is isomorphic to \mathcal{G}' . (Function $\deg(v)$ denotes the degree of the vertex v .)

Proof. Our goal is to find a subgraph $\mathcal{G}' \subseteq \mathcal{G}$ that is isomorphic to \mathcal{H} . We employ a combinatorial argument based on the properties of regular graphs and the degrees of vertices in \mathcal{H} . We will take the following steps to construct a subgraph \mathcal{G} of \mathcal{G} using a greedy algorithm: (1) Start with an empty subgraph \mathcal{G}' . (2) Iterate through the vertices of \mathcal{H} . For each vertex v_i in \mathcal{H} , we will select $\deg(v_i)$ neighbors from \mathcal{G} that are not yet included in \mathcal{G}' and add them to \mathcal{G}' . (3) For each vertex v_i in \mathcal{H} : Identify $\deg(v_i)$ vertices from \mathcal{G} that are not already included in \mathcal{G}' and that can serve as neighbors. This selection is possible because \mathcal{G} is K -regular. As long as the vertices chosen are distinct and do not exceed K in degree, we can ensure that the selections are valid. (4) To ensure that the resulting subgraph \mathcal{G}' is isomorphic to \mathcal{H} , we need to maintain the structure of \mathcal{H} . The selection of neighbors in \mathcal{G} should respect the edges present in \mathcal{H} . By the property of regular graphs, since \mathcal{G} has sufficient edges to accommodate the degree constraints of \mathcal{H} (as \mathcal{G} has K edges available per vertex), we can map the edges of \mathcal{H} to \mathcal{G}' .

After selecting neighbors for all vertices of \mathcal{H} , we have ensured that \mathcal{G}' has n vertices. The edges selected match the edges of \mathcal{H} in terms of structure. Thus, we have constructed a subgraph \mathcal{G}' of \mathcal{G} such that \mathcal{G}' is isomorphic to \mathcal{H} . \square

3.3.2 PROPERTY-GUIDED OUTPUT DISTRIBUTION

The property-guided output distribution receives the parameters θ at step t and the condition information c (e.g., the LogP property) to approximate the sender distribution at each step. The output distribution differs from the input distribution in that the parameters of each variable's input distribution are independent and do not incorporate contextual information (such as neighboring information). In contrast, the output distribution is formed by jointly processing the parameters of the input distributions for all variables, resulting in an interdependent distribution that integrates the contextual relationship among these variables. To achieve this, we can employ graph neural networks

¹In principle, there are only $\frac{N_v \times N_v - 1}{2}$ edges needed to be determined since the edge matrix is symmetry.

(denoted by Φ) to aggregate all the parameters from the input distribution and produce the output distribution.

Specifically, the input distribution parameters θ and the process time t are fed into a neural network (denoted by Φ). The network then generates the output distribution parameters $\Phi(\theta, t) = (\Phi^{(1)}(\theta, t), \Phi^{(2)}(\theta, t), \dots, \Phi^{(D)}(\theta, t))$. For a molecular graph $\mathcal{M} = \{(\mathbf{V}, \mathbf{E})\}$, the output parameters $\Phi^{(i)}(\theta, t)$ are passed through a softmax function to produce the output categorical distribution, with the probability corresponding to the true label $v^{(i)}$ for atom $v^{(i)} \in \mathbf{V}$ given by:

$$p_O^{(i)}(v^{(i)} | \theta; t; \mathbf{c}) = (\text{softmax}(\Phi^{(i)}(\theta, t, \mathbf{c})))_{v^{(i)}}. \quad (9)$$

Then the output distribution of atoms \mathbf{V} is calculated as follows:

$$P_O(\mathbf{V} | \theta; t; \mathbf{c}) = \prod_{i=1}^{N_v} p_O^{(i)}(v^{(i)} | \theta; t; \mathbf{c}). \quad (10)$$

Note that the output distribution parameters $\Phi^{(i)}(\theta, t)$ depend on all of θ through the neural network Φ , enabling access to contextual information.

Similarly, the output distribution of edges \mathbf{E} can be derived in the same manner as above. Here, we utilize a graph attention neural network Φ to jointly process all the input distribution parameters. The output node embedding and edge embedding are projected to the K_V and K_E dimensions for the output categorical distribution of atoms \mathbf{V} and edges \mathbf{E} , respectively.

3.3.3 SENDER DISTRIBUTION

The sender distribution is used to construct observation samples \mathbf{y} , which can convey information about the real data and can be considered as the messenger of the data.

Following Graves et al. (2023), the sender distribution of atom $v^i \in \mathbf{V}$ could be derived with the central limit theorem, which lies in the form of

$$p_S(y_V^{(i)} | v_{(i)}; \alpha) = \mathcal{N}(y_V^{(i)} | \alpha(K_V \mathbf{o}_{v^{(i)}} - \mathbf{1}), \alpha K_V \mathbf{I}), \quad (11)$$

where $\mathbf{1}$ is a vector of ones, \mathbf{I} is the identity matrix, and $\mathbf{o}_j \in \mathbb{R}^{K_V}$ is a vector defined as the projection from the class index j to a length K_V one-hot vector. Each entry in \mathbf{o}_j is defined as $(\mathbf{o}_j)_k = \delta_{jk}$, where δ_{jk} is the Kronecker delta function. Then the sender distribution of atoms \mathbf{V} is expressed as follows,

$$p_S(\mathbf{y}_V | \mathbf{V}; \alpha) = \mathcal{N}(\mathbf{y}_V | \alpha(K_V \mathbf{o}_V - \mathbf{1}), \alpha K_V \mathbf{I}), \quad (12)$$

where $\mathbf{o}_V = (\mathbf{o}_{v^{(1)}}, \mathbf{o}_{v^{(2)}}, \dots, \mathbf{o}_{v^{(N_v)}}) \in \mathbb{R}^{K_V \times N_v}$.

3.3.4 RECEIVER DISTRIBUTION

To obtain the receiver distribution for atom $v^i \in \mathbf{V}$, we substitute Eq 10 and Eq 12 into Eq 1 resulting in the following receiver distribution

$$p_R^{(i)}(y_V^{(i)} | \theta; t, \alpha) = \sum_{v^{(i)}=1}^{K_V} p_O^{(i)}(v^{(i)} | \theta; t) \cdot p_S(y_V^{(i)} | v_{(i)}; \alpha), \quad (13)$$

Then distribution distribution of atoms \mathbf{V} is produced by

$$p_R(\mathbf{y}_V | \theta; t, \alpha) = \prod_{i=1}^{N_v} p_R^{(i)}(y_V^{(i)} | \theta; t, \alpha). \quad (14)$$

3.3.5 BAYESIAN UPDATES

Taking atoms \mathbf{V} as an example, to update the input distribution parameters θ_i given θ_{i-1} and the atoms sample \mathbf{y}_V drawn from $p_S(\mathbf{y}_V | \mathbf{V}; \alpha_i)$ (Eq 12), we utilize the Bayesian update function h defined in Graves et al. (2023) as follows

$$\theta_i \leftarrow h(\theta_{i-1}, \mathbf{y}_V, \alpha_i), \quad (15)$$

$$h(\boldsymbol{\theta}_{i-1}, \mathbf{y}_V, \alpha_i) \stackrel{\text{def}}{=} \frac{e^{\mathbf{y}_V \boldsymbol{\theta}_{i-1}}}{\sum_{k=1}^{K_V} e^{\mathbf{y}_V^k (\boldsymbol{\theta}_{i-1})_k}}. \quad (16)$$

Once we have the Bayesian update function h , the Bayesian update distribution can be derived by substituting Eq 12 and Eq 16 into Eq 2, that is

$$p_U(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{i-1}, V; \alpha) = \mathbb{E}_{\mathcal{N}(\mathbf{y}_V \mid \alpha(K_V \mathbf{e}_V - \mathbf{1}), \alpha K_V \mathbf{I})} \delta\left(\boldsymbol{\theta} - \frac{e^{\mathbf{y}_V \boldsymbol{\theta}_{i-1}}}{\sum_{k=1}^{K_V} e^{\mathbf{y}_V^k (\boldsymbol{\theta}_{i-1})_k}}\right) \quad (17)$$

Then we substitute Eq 17 into Eq 4, generating the Bayesian flow distribution

$$p_F(\boldsymbol{\theta} \mid \mathbf{V}; t) = \mathbb{E}_{\mathcal{N}(\mathbf{y}_V \mid \beta(t)(K_V \mathbf{e}_V - \mathbf{1}), \beta(t) K_V \mathbf{I})} \delta\left(\boldsymbol{\theta} - \frac{e^{\mathbf{y}_V \boldsymbol{\theta}_0}}{\sum_{k=1}^{K_V} e^{\mathbf{y}_V^k (\boldsymbol{\theta}_0)_k}}\right) \quad (18)$$

Recall that the prior is set to uniform distribution with parameters $\boldsymbol{\theta}_0 = \frac{\mathbf{1}}{K_V}$. The above equation can be transformed into

$$p_F(\boldsymbol{\theta} \mid \mathbf{V}; t) = \mathbb{E}_{\mathcal{N}(\mathbf{y}_V \mid \beta(t)(K_V \mathbf{e}_V - \mathbf{1}), \beta(t) K_V \mathbf{I})} \delta(\boldsymbol{\theta} - \text{softmax}(\mathbf{y}_V)) \quad (19)$$

4 EXPERIMENT

4.1 DATASET

Our experiments are conducted on ZINC-250K (Irwin et al., 2012), a dataset comprising 250,000 drug-like molecules selected from the ZINC database. We performed experiments on both unconditional and conditional molecule generation tasks. For conditional molecule generation, we select two properties (drug-likeness QED and hydrophobicity LogP) to guide the generation process.

4.2 COMPETING MODELS

We compare our model against several state-of-art models for molecule generation. For unconditional molecule generation, our model is evaluated alongside JT-VAE (Jin et al., 2018a), GCPN (You et al., 2018), MRNN (Popova et al., 2019), GraphNVP (Madhawa et al., 2019), GRF (Honda et al., 2019), GraphAF (Shi et al., 2020), MoFlow (Zang & Wang, 2020), GraphCNF (Lippe & Gavves, 2020), and GraphDF (Luo et al., 2021). For conditional molecule generation, we compare our model with GDSS (Jo et al., 2022), DiGress (Vignac et al., 2022), GLDM, SSSVAE (Mailoa et al., 2023) and ConGen Mailoa et al. (2023).

4.3 METRICS

We use several widely used metrics to evaluate our model on both conditional and unconditional molecular graph generation. For unconditional molecule generation, we focus on three metrics: validity, uniqueness, and novelty. Validity measures the percentage of molecules that comply with chemical valency rules. In addition, since some methods check the valency during the generation process, we also report Validity w/o check, which reflects the validity percentage when this valency correction is disabled. Uniqueness is defined as the percentage of unique molecules among the generated valid molecules. Novelty represents the percentage of the generated molecules that do not appear in the training data. For conditional molecule generation, we evaluate the performance of each model by comparing the mean absolute error (MAE) between the target properties and those of the generated molecules. All the metrics are calculated based on 10000 generated molecules.

4.4 IMPLEMENTATION DETAILS

Our CtrlMol uses a graph attention network to model the interdependence of the different variables and produce the output distribution. The network consists of 9 graph attention layers with 8 attention heads. Both the hidden dimension and the output dimension are 256. The node and edge features output by the graph network are then mapped to the dimensions corresponding to the number of atom types $K_V = 9$ and edge types $K_E = 5$ to generate the output distribution. The model is trained through the Adam optimizer for 1000 epochs, where the batch size is 256 and the learning rate is 0.0005.

Table 1: Random generation performance on ZINC250K dataset.

Method	Validity	Validity w/o check	Uniqueness	Novelty
JT-VAE	100%	-	100%	100%
GCPN	100%	20%	99.97%	100%
MRNN	100%	65%	99.89%	100%
GraphNVP	42.6%	-	94.8%	100%
GRF	73.4%	-	53.7%	100%
GraphAF	100%	68%	99.1%	100%
MoFlow	100%	50.3%	99.99%	100%
GraphCNF	96.35%	-	99.98%	99.98%
GraphDF	100%	89.03%	99.16%	100%
Ours	100%	91.6%	100%	100%

Table 2: Mean Absolute Error for molecular property guided generation. A lower number indicates a better controllable generation result.

Property	QED	LogP
GDSS	0.44	-
DiGress	0.53	-
GLDM	0.41	-
HGLDM	0.35	-
SSVAE	0.81	2.38
ConGen	0.82	2.13
Ours	0.26	0.77
Improvement over SOTA	25.71%	63.85%

4.5 RESULTS

Unconditional molecule generation We evaluate the ability of the CtrlMol to generate realistic molecules by calculating metrics validity, uniqueness, and novelty. Table 1 shows that our model outperforms the previous methods across these metrics. While several models utilize valency checks during each sampling step to ensure the validity of generated molecules, this checking process relies on external chemical rules that the model itself does not learn, which does not accurately reflect the model’s true capabilities. For instance, GCPN may achieve 100% validity through valency checking, yet it fails to capture the underlying distribution of molecules (20% validity without check).

Therefore, we focus on the validity without check, which can truly reflect the model’s understanding of molecular characteristics. CtrlMol achieves the best results in validity without checking, highlighting the superiority of our model in generating realistic molecules through a one-shot sampling process without re-sampling by valency check. Additionally, CtrlMol demonstrates 100% uniqueness and novelty, highlighting its capability to produce new molecular structures that do not appear in the training data. This indicates that the model has genuinely captured the underlying distribution characteristics of the molecules, rather than merely overfitting to the training data.

Conditional molecule generation For conditional molecule generation, we evaluate the CtrlMol model based on 2 properties: hydrophobicity LogP and drug-likeness QED. The model is provided with a specific property value to generate molecules that align with the desired property c . The property of generated molecule \hat{c} is measured by RDKit. The mean absolute error (MAE) between c and \hat{c} is calculated to determine the relevance of the generated molecules to the conditioned property.

The results are shown in Table 2. The MSE of our model is considerably lower than that of other baseline models by an obvious margin in both the QED and LogP conditional generation tasks, suggesting that it generates molecules that more effectively align with the property conditions.

Table 3: The performance of CtrlMol with respect to different sampling steps

Sample steps	10	100	300	500	700	900	1000 (reference)
Validity w/o check	85.9%	88.9%	89.8%	88.0%	87.6%	91.5%	91.7%
LogP MSE	0.93	0.85	0.74	0.76	0.81	0.75	0.77

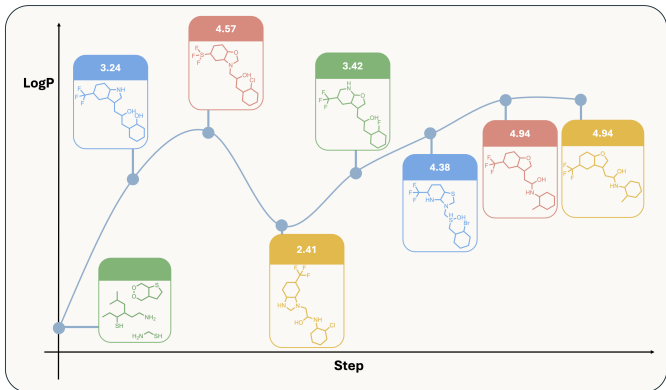


Figure 2: An example of the conditional molecule generation process.

Few-step sampling Furthermore, the CtrlMol model exhibits excellent performance with a limited number of sampling steps. We show the performance in unconditional generation (validity w/o check metric) with different sampling steps in Table 3. Our model requires only 300 sampling steps to surpass the "validity w/o check" metric of GraphDF using 1000 sampling steps. Notably, with just 10 sampling steps, the model achieved a validity w/o check of 85.9%, already exceeding the results of other models with 1000 sampling steps, resulting in a $100\times$ speed-up during the sampling generation process. Similarly, in conditional generation, CtrlMol also demonstrates superiority, the property LogP of generated molecules with only 10 sampling steps is more accurate compared to other models.

Case study We illustrate a conditional molecule generation process in Figure2. Given a target LogP value of 5.0, the CtrlMol model successfully generates a molecule that aligns with the desired property within 10 sampling steps. In the first step, the output samples consist of some molecular fragments. Subsequently, the model begins to output valid molecules. As the time steps progress, the logP values of the molecules gradually converge toward the target. After seven steps, the output molecules stabilize and no longer change. The logP value of the final molecule is 4.9, showcasing the model's proficiency in closely aligning with the specified property while maintaining structural validity.

5 CONCLUSION

Deep generative models have achieved remarkable progress in generating continuous data, such as lifelike images and 3D protein structures. Nonetheless, purely discrete data like topological molecular graphs suffer the serious issues of the multi-modality and the connection sparsity, making molecular graph generation a challenging task for the existing methods. In this paper, we introduce CtrlMol, an innovative framework for generating molecular graphs within a differentiable parameter space. In contrast to diffusion models that enhance samples through iterative refinement, CtrlMol adjusts distribution parameters at different noise levels via a non-parametric Bayesian flow. CtrlMol also integrates a topological complete edge sampling strategy and property-guided output distribution to address the above generative issue. The experimental results indicate that CtrlMol outperforms all current benchmarks in creating realistic molecular graphs, also establishing a new benchmark for generating molecules with pre-defined characteristics.

REFERENCES

- 486
487
488 G Richard Bickerton, Gaia V Paolini, Jérémy Besnard, Sorel Muresan, and Andrew L Hopkins.
489 Quantifying the chemical beauty of drugs. *Nature chemistry*, 4(2):90–98, 2012.
- 490
491 Hanjun Dai, Yingtao Tian, Bo Dai, Steven Skiena, and Le Song. Syntax-directed variational
492 autoencoder for structured data. In *International Conference on Learning Representations*, 2018.
- 493
494 Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato,
495 Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel,
496 Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous
497 representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- 498
499 Alex Graves, Rupesh Kumar Srivastava, Timothy Atkinson, and Faustino Gomez. Bayesian flow
500 networks. *arXiv preprint arXiv:2308.07037*, 2023.
- 501
502 Jiaqi Guan, Wesley Wei Qian, Xingang Peng, Yufeng Su, Jian Peng, and Jianzhu Ma. 3d equivariant
503 diffusion for target-aware molecule generation and affinity prediction. In *The Eleventh International
504 Conference on Learning Representations*, 2022.
- 505
506 Jiaqi Guan, Xiangxin Zhou, Yuwei Yang, Yu Bao, Jian Peng, Jianzhu Ma, Qiang Liu, Liang Wang,
507 and Quanquan Gu. Decompdiff: Diffusion models with decomposed priors for structure-based
508 drug design. In *International Conference on Machine Learning*, pp. 11827–11846. PMLR, 2023.
- 509
510 Shion Honda, Hirotaka Akita, Katsuhiko Ishiguro, Toshiki Nakanishi, and Kenta Oono. Graph
511 residual flow for molecular graph generation. *arXiv preprint arXiv:1909.13521*, 2019.
- 512
513 John J Irwin, Teague Sterling, Michael M Mysinger, Erin S Bolstad, and Ryan G Coleman. Zinc: a
514 free tool to discover chemistry for biology. *Journal of chemical information and modeling*, 52(7):
515 1757–1768, 2012.
- 516
517 Jan H Jensen. A graph-based genetic algorithm and generative model/monte carlo tree search for the
518 exploration of chemical space. *Chemical science*, 10(12):3567–3572, 2019.
- 519
520 Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for
521 molecular graph generation. In *International conference on machine learning*, pp. 2323–2332.
522 PMLR, 2018a.
- 523
524 Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for
525 molecular graph generation. In *International conference on machine learning*, pp. 2323–2332.
526 PMLR, 2018b.
- 527
528 Jaehyeong Jo, Seul Lee, and Sung Ju Hwang. Score-based generative modeling of graphs via the
529 system of stochastic differential equations. In *International conference on machine learning*, pp.
530 10362–10383. PMLR, 2022.
- 531
532 Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder.
533 *Proceedings of the 34th International Conference on Machine Learning*, 70:1945–1954, 2017.
- 534
535 Phillip Lippe and Efstratios Gavves. Categorical normalizing flows via continuous transformations.
536 *arXiv preprint arXiv:2006.09790*, 2020.
- 537
538 Youzhi Luo, Keqiang Yan, and Shuiwang Ji. Graphdf: A discrete flow model for molecular graph
539 generation. In *International conference on machine learning*, pp. 7192–7203. PMLR, 2021.
- 539
540 Kaushalya Madhawa, Katushiko Ishiguro, Kosuke Nakago, and Motoki Abe. Graphnvp: An invertible
541 flow model for generating molecular graphs. *arXiv preprint arXiv:1905.11600*, 2019.
- 542
543 Jonathan P Mailoa, Xin Li, Jiezhong Qiu, and Shengyu Zhang. Multi-constraint molecular generation
544 using sparsely labelled training data for localized high-concentration electrolyte diluent screening.
545 *Digital Discovery*, 2(5):1390–1403, 2023.
- 546
547 Mariya Popova, Mykhailo Shvets, Junier Oliva, and Olexandr Isayev. Molecularrrn: Generating
548 realistic molecular graphs with optimized properties. *arXiv preprint arXiv:1905.13372*, 2019.

- 540 Yanru Qu, Keyue Qiu, Yuxuan Song, Jingjing Gong, Jiawei Han, Mingyue Zheng, Hao Zhou, and
541 Wei-Ying Ma. Molcraft: Structure-based drug design in continuous parameter space. *arXiv*
542 *preprint arXiv:2404.12141*, 2024.
- 543
544 Nicholas T. Runcie and Antonia S.J.S. Mey. Silvr: Guided diffusion for molecule generation. *Journal*
545 *of Chemical Information and Modeling*, 63(19):5996–6005, 2023. doi: 10.1021/acs.jcim.3c00667.
546 PMID: 37724771.
- 547 Arne Schneuing, Yuanqi Du, Charles Harris, Arian Jamasb, Ilia Igashov, Weitao Du, Tom Blundell,
548 Pietro Lió, Carla Gomes, Max Welling, et al. Structure-based drug design with equivariant diffusion
549 models. *arXiv preprint arXiv:2210.13695*, 2022.
- 550 Marwin HS Segler, Thierry Kogej, Christian Tyrchan, and Mark P Waller. Generating focused
551 molecule libraries for drug discovery with recurrent neural networks. *ACS central science*, 4(1):
552 120–131, 2018.
- 553
554 Chence Shi, Minkai Xu, Zhaocheng Zhu, Weinan Zhang, Ming Zhang, and Jian Tang. Graphaf: a
555 flow-based autoregressive model for molecular graph generation. *arXiv preprint arXiv:2001.09382*,
556 2020.
- 557 Yuxuan Song, Jingjing Gong, Hao Zhou, Mingyue Zheng, Jingjing Liu, and Wei-Ying Ma. Unified
558 generative modeling of 3d molecules with bayesian flow networks. In *The Twelfth International*
559 *Conference on Learning Representations*, 2023.
- 560
561 Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard.
562 Digress: Discrete denoising diffusion for graph generation. *arXiv preprint arXiv:2209.14734*,
563 2022.
- 564 Jiaxuan You, Bowen Liu, Zhitao Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy
565 network for goal-directed molecular graph generation. *Advances in neural information processing*
566 *systems*, 31, 2018.
- 567
568 Chengxi Zang and Fei Wang. Moflow: an invertible flow model for generating molecular graphs. In
569 *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data*
570 *mining*, pp. 617–626, 2020.
- 571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593