

Smooth Interpolation for Improved Discrete Graph Generative Models

Yuxuan Song^{*1,2} Juntong Shi^{*3} Jingjing Gong^{*1,2} Minkai Xu⁴
Stefano Ermon⁴ Hao Zhou¹ Wei-Ying Ma¹

Abstract

Though typically represented by the discrete node and edge attributes, the graph topological information can be sufficiently captured by the graph spectrum in a continuous space. It is believed that incorporating the continuity of graph topological information into the generative process design could establish a superior paradigm for graph generative modeling. Motivated by such prior and recent advancements in the generative paradigm, we propose Graph Bayesian Flow Networks (GraphBFN) in this paper, a principled generative framework that designs an alternative generative process emphasizing the dynamics of topological information. Unlike recent discrete-diffusion-based methods, GraphBFN employs the continuous counts derived from sampling infinite times from a categorical distribution as latent to facilitate a smooth decomposition of topological information, demonstrating enhanced effectiveness. To effectively realize the concept, we further develop an advanced sampling strategy and new time-scheduling techniques to overcome practical barriers and boost performance. Through extensive experimental validation on both generic graph and molecular graph generation tasks, GraphBFN could consistently achieve superior or competitive performance with significantly higher training and sampling efficiency. The code is available at <https://github.com/AlgoMole/GraphBFN>.

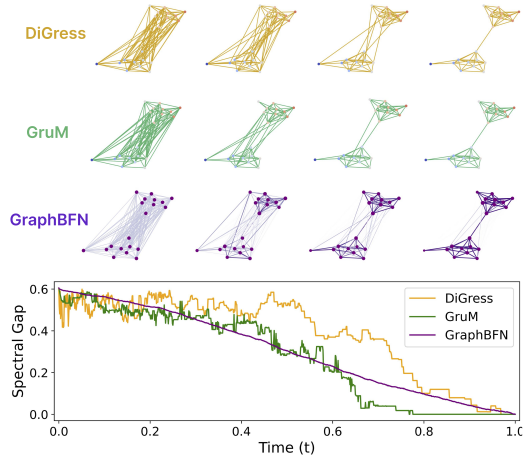


Figure 1: The top figures compares the sampling trajectories of GraphBFN to DiGress and GruM (Vignac et al., 2022; Jo et al., 2024). A smaller spectral gap reflects better graph clusterability (i.e. clearer sample). Both the visualizations and curves demonstrate that GraphBFN allows a much smoother transformation of graph topology information.

1. Introduction

Generative modeling of graph-structured data is an important task with applications in various important scenarios, *e.g.* molecule generation (Jin et al., 2018; Zang & Wang, 2020), traffic modeling (Yu & Gu, 2019) and protein design (Ingraham et al., 2019). With the development of deep generative models, there have been fruitful lines of research conducted to tackle the challenges, *e.g.*, Generative Adversarial Networks (Martinkus et al., 2022; De Cao & Kipf, 2018), Variational Auto-encoders (Jin et al., 2018; 2020), and Autoregressive models (You et al., 2018).

Recently, diffusion models (DMs) (Song & Ermon, 2019; Ho et al., 2020) have emerged as a powerful generative model across various fields (Zeng et al., 2022; Li et al., 2022). Learning DMs for graphs is challenging due to the discrete and complex relational nature of graphs. Niu et al. (2020); Jo et al. (2022) propose to dequantize the node vector and adjacent matrix with uniform/Gaussian noise and learn Gaussian DMs over the dequantized variables. However, these approaches face problems due to the in-

^{*}Equal contribution ¹Institute for AI Industry Research (AIR), Tsinghua University ²Department of Computer Science and Technology, Tsinghua University ³University of Southern California, USA ⁴Stanford University, USA. Correspondence to: Hao Zhou <zhouhao@air.tsinghua.edu.cn>.

compatibility of continuous Gaussian diffusion processes with discrete data. For instance, they involve noisy graph samples that lack well-defined topological information, *e.g.*, clusterability or connectivity. To alleviate these problems, Vignac et al. (2022) proposed discrete graph diffusion models based on discrete DMs (Austin et al., 2021). It utilizes a structured categorical corruption process (corresponding to successive graph edit operations such as additions, deletions, and replacements) to destroy the graph and learns to generate by reverting this process. Such discrete formulation ensures noisy samples are valid discrete graphs with well-defined graph topologies. This enables the discrete DMs to extract rich graph spectral features, such as graph spectrums, that aid the generation process, thus achieving better performance.

However, though always ensuring valid graph topology, the discrete graph diffusion process is very unbalanced. Some graph editing steps can significantly change the adjacency matrix and unexpectedly perturb the graph topology, making it difficult to learn a satisfactory model for the reverse diffusion process. Recently, utilizing certain kinds of topological information formed in continuous space to provide smooth and consistent information flow for the generative process has been demonstrated to be effective for graph generative model (Martinkus et al., 2022; Vignac et al., 2022; Jo et al., 2024). Based on such phenomena, we hold the intuition that stabilized dynamics of graph topological information during the generation process could be favorable for discrete graphs. Furthermore, we aim to design a new generative framework for discrete graphs that could conduct a smooth transition from data to noise, where the smoothness is quantified by generalized graph topological features.

Well motivated by this philosophy, we bring the concept of Bayesian Flow Networks, a recent advancement of generative paradigms, into the context of graph generation. Different from the discrete diffusion model whose latent variable is the discrete variable following a predefined forward categorical noisy distribution, the latent variable of the Bayesian Flow Network could be defined with the counts of discrete variables by sampling *infinite times* from the same forward categorical noisy distribution and is hence continuous. By projecting the continuous latent of BFN into probability simplex, we obtain a mixed-state representation of different discrete states which enables the smooth interpolation between an uninformative prior to the concrete observation. Based on such appealing property, we propose a new and general graph generative framework named Graph Bayesian Flow Networks (GraphBFN). By substituting the discrete variable to probability simplex in the intermediate state of the generative process, GraphBFN makes a probabilistic generalization of the adjacency matrix as *probabilistic adjacency matrix* whose entries represent the probability of corresponding edge type to emerge. Fortunately, graph topo-

logical features, like the graph spectrum, easily generalize to the *probabilistic adjacency matrix*, enabling straightforward analysis and the implementation of effective techniques such as spectrum conditioning. Quantitative measurements in Fig. 1 show that GraphBFN enables stabilized topological dynamics compared to diffusion-based methods with the target graph topology being smoothly recovered during the generation process. Specifically, our contribution can be summarized as follows:

- We firstly prove the concept of smooth transformation over topological information could facilitate the discrete graph generation with a new method termed GraphBFN based on Bayesian Flow Networks. We empirically demonstrated the effectiveness of smooth topological recovery. GraphBFN could fit more general graphs with rich node and edge features.
- To successfully implement the method and boost empirical performance, we introduce several innovations to sampling and training, including the adaptive sampling strategy which provides better tradeoffs between quality and diversity, a general time-scheduling that shows good compatibility with low entropy prior, and a novel graph spectrum conditioned strategy for large graphs.
- We conduct extensive experiments and ablation studies on several benchmarks including both abstract and 2D molecules. The empirical results show that the GraphBFN can consistently achieve superior or competitive performance in generating realistic and diverse graphs. Furthermore, GraphBFN also enjoys a high sampling efficiency with up to 10 \times speedup while matching the performance diffusion-based methods.

2. Related Works

From the fruitful lines of research conducted on generative modeling for graphs, we focus on the discussion over diffusion-based approaches as they are most related to our method. Recently, given the remarkable performance of diffusion models (Ho et al., 2020) in fields such as image generation, there has been a growing interest in graph diffusion models within the community. Research conducted by Jo et al. (2022); Niu et al. (2020); Luo et al. (2023) has successfully made the diffusion framework compatible with discrete graphs through the dequantization approach that adds continuous noise to discrete node and edge features. Jo et al. (2024) demonstrate the advantage of explicitly learning the final graph structures in the diffusion process; Liu et al. (2025) show the benefits of modeling the graph generation process within a bounded space; Vignac et al. (2022); Kong et al. (2023) highlight the benefits of modeling in the

original discrete data space; [Qin et al. \(2023\)](#); [Chen et al. \(2023\)](#) leverage the sparsity nature of graphs and show the possibility to scale the models up to larger graphs; [Karami \(2024\)](#) introduce hierarchical graph generative framework that captures the hierarchical nature of graphs. Furthermore, [Vignac et al. \(2022\)](#) highlight the potential issue of topology information getting distorted by continuous noise and thereby suggests utilizing more compatible discrete diffusions, resulting in noteworthy performance improvement. Similarly, [Liu et al. \(2025\)](#) emphasize the importance of maintaining critical graph topology and realizes this using a modulation technique.

In the literature of Bayesian Flow Network ([Graves et al., 2023](#)), there are several recent applications to the structured multi-modality data ([Song et al., 2023](#); [Qu et al., 2024](#)) while the topological information modeling is less explored. Drawing inspiration from [Graves et al. \(2023\)](#), we propose the concept of probabilistic graph space, where the representation of topological information can be generalized, allowing for a smoother decomposition of information of discrete graphs. We leave a more detailed discussion over discrete diffusion models in the Appendix K.2.

3. Preliminary

3.1. Notations and Background

We define the space of undirected graphs as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{E} represents a categorical edge space and \mathcal{V} a node space. For a graph \mathbf{G} with n nodes, the adjacency matrix $\mathbf{E} \in \mathbb{R}^{n \times n \times c_e}$ and the node vector $\mathbf{V} \in \mathbb{R}^{n \times c_n}$ are both represented using one-hot encoding. Here, c_e and c_n denote the respective cardinalities. The absence of an edge is treated as a specific edge type. This framework encompasses both abstract graphs and 2D molecules studied in this work.

Then, we introduce the basics of spectral graph theory for better illustration. Spectral graph theory explores the relationships between graph properties and the eigenvalues of associated matrices, such as the adjacency matrix or Laplacian matrix. Previous works ([Martinkus et al., 2022](#); [Vignac et al., 2022](#)) have demonstrated the effectiveness of conditioning on graph spectra in generating graphs. This motivates us to revisit the generative process from the perspective of building graph spectra, as shown in Fig. 1. Detailed discussions on spectral graph theory are presented in Appendix B.

3.2. Bayesian Flow Networks

In this section, we aim to provide a very intuitive introduction to the Bayesian Flow Networks (BFNs) ([Graves et al., 2023](#)) and leave an in-depth discussion and more details in Appendix A. Similar to the diffusion models, the BFNs also

take the form of latent variable models with a series of noisy samples $\mathbf{y}_{\{1:n\}} = \langle \mathbf{y}_1, \dots, \mathbf{y}_n \rangle$ of observation \mathbf{x} are introduced as latent. We note that \mathbf{y}_1 of the start of the sequence refers to the uninformative sample of pure noise, which is the opposite of the diffusion literatures. Besides, BFNs also optimize the *variational lower bound* of the likelihood:

$$\begin{aligned} \log p_\phi(\mathbf{x}) &\geq \mathbb{E}_{\mathbf{y}_{\{1:n\}} \sim q} \left[\log \frac{p_\phi(\mathbf{x} | \mathbf{y}_{\{1:n\}}) p_\phi(\mathbf{y}_{\{1:n\}})}{q(\mathbf{y}_{\{1:n\}} | \mathbf{x})} \right] \\ &= -D_{KL}(q \| p_\phi(\mathbf{y}_{\{1:n\}})) + \mathbb{E}_{\mathbf{y}_{\{1:n\}} \sim q} \log [p_\phi(\mathbf{x} | \mathbf{y}_{\{1:n\}})] \end{aligned} \quad (1)$$

where q is a predefined distribution for creating the noisy samples similar to the noise-application distribution in the forward process of diffusion models (sender distribution) and p_ϕ refers to the probabilistic distribution(receiver distribution) implied by BFNs. From the graphical model perspective, BFNs differs from diffusion models in the construction of latent variables $\mathbf{y}_{\{1:n\}} = \langle \mathbf{y}_1, \dots, \mathbf{y}_n \rangle$, which is defined as non-Markovian and could be extended as:

$$p_\phi(\mathbf{y}_1, \dots, \mathbf{y}_n) = p_\phi(\mathbf{y}_1) \prod_{t=2}^n p_\phi(\mathbf{y}_t | \mathbf{y}_{\{1:t-1\}}) \quad (2)$$

The $p_\phi(\mathbf{y}_t | \mathbf{y}_{\{1:t-1\}})$ could be further represented with a predefined deterministic aggregation function of the process as $p_\phi(\mathbf{y}_t | f(\mathbf{y}_{\{1:t-1\}}))$. The f is related to the concrete Bayesian update rule for different distributions, and $f(\mathbf{y}_{\{1:t-1\}})$ is also notated as $\boldsymbol{\theta}_{t-1}$. Putting Eq. 2 into Eq. 1, the ELBO of BFNs could be further expressed as:

$$\begin{aligned} \mathcal{L}(\mathbf{x}) &= -\mathcal{L}_{\text{ELBO}} \\ &= \sum_{t=1}^n D_{KL}(q(\mathbf{y}_t | \mathbf{x}) \| p_\phi(\mathbf{y}_t | \boldsymbol{\theta}_t = f(\mathbf{y}_{\{1:t-1\}}))) \\ &\quad - \log p_O(\mathbf{x} | \boldsymbol{\theta}_n, \phi) \end{aligned} \quad (3)$$

The term $\log p_O(\mathbf{x} | \boldsymbol{\theta}_n, \phi)$ is also known as reconstruction loss, which is directly computed over the network prediction with $\boldsymbol{\theta}_n$. As the term is generally negligible as $\boldsymbol{\theta}_n$ is very close to \mathbf{x} , we do not involve the term in training loss following ([Graves et al., 2023](#)).

4. Methodology

In this section, we first identify the key difference between the GraphBFN and discrete diffusion-based approaches and discuss how the GraphBFN could generalize smooth dynamics of topological information for graph generation. Then, we provide a detailed formulation of our GraphBFN, along with a visualization of the training and testing procedure in Fig. 2. Lastly, we delve into several innovations for successfully realizing the concept and boosting the applications in graph generation.

4.1. Smoothing Topological Transformation via Infinite Noisy Discrete Graphs

We consider the following specific case of the perturbation distribution over discrete variables:

$$p_{\text{noise}}(\mathbf{y}_t | \mathbf{x}) = \text{Cat}(\mathbf{y}_t; \omega_t \mathbf{x} + (1 - \omega_t) \boldsymbol{\pi}) \quad (4)$$

Hence \mathbf{x} refers to the one-hot representation over the K categories and $\boldsymbol{\pi}$ is the uniform prior over the categories, i.e. $[\frac{1}{K}, \dots, \frac{1}{K}]$. For the discrete diffusion model (Vignac et al., 2022; Austin et al., 2021), the Eq. 4 could be seen as continuous-time variant of the variational distribution $q(\mathbf{y}_t | \mathbf{x})$ in Eq. 1, i.e. the extension of discrete-time uniform transition diffusion (Austin et al., 2021). In this scenario, the latent variable y_t exists in a discrete space and is generated through random substitution or deletion of elements from the clean sample x , with a probability of $1 - \omega_t$. In the context of graphs, these operations can significantly alter the topology; for example, deleting a single edge in the adjacency matrix may impact the entire graph’s connectivity.

Bayesian Flow Networks introduce a different concept of latent \mathbf{y}_t : Given the noisy distribution p_{noise} , we assume sample noisy samples from the noise distribution for m times instead of once, thus incurring a count variable over the different categories as $\mathbf{c} = [c_1, \dots, c_K]$. Based on the normalized counts $(\mathbf{c} - \frac{m}{K})$ and central limit theorem $\lim_{m \rightarrow \infty} \frac{\mathbf{c} - \frac{m}{K}}{\sqrt{mq_t}} \sim \mathcal{N}(0, \mathbf{I})$ (where $q_t = \omega_t \mathbf{x} + (1 - \omega_t) \boldsymbol{\pi}$), we could construct a continuous latent variable \mathbf{y} with the k -th dimension as as:

$$y_k \stackrel{\text{def}}{=} \lim_{m \rightarrow \infty, \omega \rightarrow 0} \left(c_k - \frac{m}{K} \right) \ln \left(1 + \frac{\omega K}{1 - \omega} \right), \quad (5)$$

$$\text{we have } q(\mathbf{y}_t | \mathbf{x}; \alpha_t) = \mathcal{N}(\alpha_t (K \mathbf{x} - \mathbf{1}), \alpha_t K \mathbf{I}) \quad (6)$$

where $m\omega^2 = \alpha$ is a condition and α is a predefined finite number. The above \mathbf{y} is obtained by assuming sampling **infinite samples** from the noise distribution. Detailed discussion is left as Appendix I.

In the context of graph generative modeling, the \mathbf{x} could refer to the unit component of the graph, e.g. a node in the node vector (\mathbf{V}_i) or an edge entry in the adjacency matrix ($\mathbf{E}_{i,j}$). For topological analysis, we focus on the discussion of the adjacency matrix. We use \mathbf{E}^{y_t} to denote the latent variable for the whole matrix \mathbf{E} , where the latent variable for each entry is therefore referred to as $\mathbf{E}_{i,j}^{y_t}$. Under the framework of GraphBFN, $\mathbf{E}_{i,j}^{y_t}$ is a count-based latent as shown in Eq. 5. Though the latent variable enjoys continuity and real-valued, it could be hard to view the \mathbf{E}^{y_t} as valid adjacency matrices and analyze the topological dynamics from the perspective of the latent variable. Fortunately, we could instead focus on the parameter variable $\boldsymbol{\theta}_t$ which is the aggregation of the latent variable as introduced in Sec. 3.2. The aggregation f in Eq. 3 for the discrete variable

is defined by applying the following update rule recursively:

$$\boldsymbol{\theta}_t = h(\boldsymbol{\theta}_{t-1}, y_t) = \frac{e^{y_t} \boldsymbol{\theta}_{t-1}}{\sum_{k=1}^K e^{y_t^{(k)}} \boldsymbol{\theta}_{t-1}^{(k)}} \quad (7)$$

Here $\boldsymbol{\theta}$ lies in the probability simplex, i.e. $\sum_{i=1}^K \boldsymbol{\theta}^i = 1$. In the above adjacency matrix setting, we use $\mathbf{E}_{i,j}^{\boldsymbol{\theta}_t}$ to denote the corresponding parameter. Given $\mathbf{E}_{i,j}^{\boldsymbol{\theta}_t}$ is on the probability simplex, it could be naturally interpreted as the probability over whether there is an edge between node i and node j . We design a *Probabilistic Adjacency Matrix* P based on $\mathbf{E}_{i,j}^{\boldsymbol{\theta}_t}$: we force $\mathbf{E}^{\boldsymbol{\theta}}$ to symmetric by $\mathbf{E}^{\boldsymbol{\theta}} = (\mathbf{E}^{\boldsymbol{\theta}} + (\mathbf{E}^{\boldsymbol{\theta}})^T)/2$; then, we derive the probability of whether there exist an edge at $E_{i,j}$ from $\mathbf{E}_{i,j}^{\boldsymbol{\theta}}$ and make it as the corresponding value for the $P_{i,j}$. Under this definition of P each entry $P_{i,j} \in [0, 1]$ could be seen as the parameter of the Bernoulli distribution.

Remark 4.1. The graph spectrum, which is defined as the eigenvalues or eigenvectors of the adjacency matrix or Laplacian matrix could be directly extended by viewing P as the generalization of the adjacency matrix. We leave the detailed discussion over the graph spectrum in Appendix B.

The probabilistic adjacency matrix can be viewed as a special case of random graphs as in Athreya et al. (2018), providing bounds between the graph spectrum of sampled discrete graphs from P and the generalized spectrum directly over P (Chung & Radcliffe, 2011). The theoretical discussion is in Appendix C. The generation process in GraphBFN defines a transformation from $\mathbf{E}^{\boldsymbol{\theta}_0}$ to $\mathbf{E}^{\boldsymbol{\theta}_n}$, interpolating between $P^{\boldsymbol{\theta}_0}$, a prior probabilistic matrix with all entries of 0.5, and $P^{\boldsymbol{\theta}_n}$, representing authentic graphs with binary entries. This interpolation enables GraphBFN to smoothly decompose complex discrete structures, as illustrated in Fig. 1. The detailed formulation of interpolation and training objectives will be introduced in the following section.

4.2. Graph Bayesian Flow Networks

We elaborate on the components introduced in Sec. 3.2 to get a complete formulation of GraphBFN. Recall that the graph variable is denoted as $\mathbf{G} = (\mathbf{V}, \mathbf{E})$. Correspondingly, the latent variable at time step t is referred to as $\mathbf{G}^{y_t} = (\mathbf{V}^{y_t}, \mathbf{E}^{y_t})$ and the corresponding aggregation parameter $\mathbf{G}^{\boldsymbol{\theta}_t} = (\mathbf{V}^{\boldsymbol{\theta}_t}, \mathbf{E}^{\boldsymbol{\theta}_t})$. Then, we revisit the ELBO training objective introduced in Eq. 3. To improve conciseness, we use the adjacency matrix \mathbf{E} as the primary example for derivation. Below, we illustrate the concrete formulation for the two key elements in the Eq. 3, $q(\mathbf{y}_t | \mathbf{x})$ and $p_{\phi}(\mathbf{y}_t | \boldsymbol{\theta}_t = f(\mathbf{y}_{\{1:t-1\}}))$.

Sender distribution The $q(\mathbf{y}_t | \mathbf{x})$ in the context of BFN is termed sender distribution. To define the sender distribution, we need extra involve an accuracy parameter α_t as also mentioned above in Eq. 5. Note the sender distribution

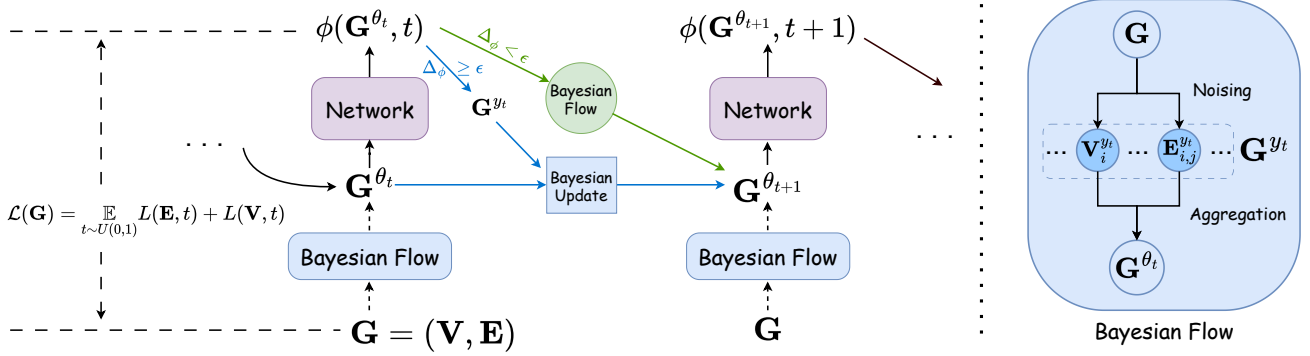


Figure 2: Diagram of GraphBFN. The solid lines and dashed lines denote the inference and training procedures, respectively. We note that $\Delta_\phi = \|\phi(\mathbf{G}^{\theta_t}, t) - \phi(\mathbf{G}^{\theta_{t-1}}, t-1)\|^2$.

for the adjacency matrix is set as independent for different entries, then the sender distribution here takes the form of

$$\begin{aligned} q(\mathbf{E}^{y_t} | \mathbf{E}, \alpha_t) &= \prod_{1 \leq i, j \leq n} q(\mathbf{E}_{i,j}^{y_t} | \mathbf{E}_{i,j}, \alpha_t) \\ &= \prod_{1 \leq i, j \leq n} \mathcal{N}(\mathbf{E}_{i,j}^{y_t} | \alpha_t (c_e \mathbf{E}_{i,j} - \mathbf{1}), \alpha_t c_e \mathbf{I}) \end{aligned} \quad (8)$$

Receiver distribution $p_\phi(\mathbf{y}_t | \theta_t = f(\mathbf{y}_{1:t-1}))$ here is also named as the receiver distribution. Intuitively, $p_\phi(\mathbf{y}_t | \theta_t)$ could be decomposed into first predict a distribution of x with neural network ϕ based on θ_t and combine predicted distribution with the sender distribution above to get a distribution over \mathbf{y} . In our case, the receiver distribution takes the form of:

$$\begin{aligned} p_\phi(\mathbf{E}^{y_t} | \mathbf{G}^{\theta_t}) &= \prod_{1 \leq i, j \leq n} p_\phi(\mathbf{E}_{i,j}^{y_t} | \mathbf{G}^{\theta_t}) \\ p_\phi(\mathbf{E}_{i,j}^{y_t} | \mathbf{G}^{\theta_t}) &= \sum_{k=1}^{c_e} \mathbf{E}_{i,j}^{\theta_t} [\phi(\mathbf{G}^{\theta_t}, t)]^{(k)} \mathcal{N}(\mathbf{E}_{i,j}^{y_t} | \alpha_t (c_e \mathbf{e}_k - \mathbf{1}), \alpha_t c_e \mathbf{I}) \end{aligned} \quad (9)$$

Here \mathbf{e}_k is the one-hot vector with k -th dimension as 1 and the other $c_e - 1$ dimensions are all 0. Given the \mathbf{G}^{θ_t} and t as network input, here $\mathbf{E}_{i,j}^{\theta_t} [\phi(\mathbf{G}^{\theta_t}, t)]^{(k)}$ is the probability on k -th category of the predicted categorical distributions for $\mathbf{E}_{i,j}$. Note the network ϕ usually has the final layer as softmax and hence output could be naturally seen as the probability of the predicted distribution, which is also referred to as output distribution.

Simulation-free Training However, as mentioned above and also shown in Eq. 3, the network input is θ_t which is determined by $\mathbf{y}_1, \dots, \mathbf{y}_n$ and the discrete update rule in Eq. 7. Fortunately, with the discrete update rule, we could do aggregation in the distribution level for $q(\mathbf{y}_1 | \mathbf{x}), \dots, q(\mathbf{y}_n | \mathbf{x})$ and directly obtain the analytic form of distribution over θ_t , i.e. $q(\theta_t | \mathbf{x})$. With the notation of \mathbf{E} , the distribution could

be illustrated as:

$$\begin{aligned} q(\mathbf{E}_{i,j}^{\theta_t} | \mathbf{E}_{i,j}, \mathbf{E}_{i,j}^{\theta_0}) &= \mathbb{E}_{\mathbf{E}_{i,j}^y \sim D(t)} \delta\left(\mathbf{E}_{i,j}^{\theta_t} - \frac{e^{\mathbf{E}_{i,j}^y \odot \mathbf{E}_{i,j}^{\theta_0}}}{\sum_{k=1}^{c_e} (e^{\mathbf{E}_{i,j}^y})_k (e^{\mathbf{E}_{i,j}^{\theta_0}})_k}\right) \end{aligned} \quad (10)$$

where $D(t) = N(\mathbf{E}_{i,j}^y | \beta(t)(c_e \mathbf{E}_{i,j} - \mathbf{1}), \beta(t)c_e \mathbf{I})$, $\beta(t) = \sum_{j=0}^t \alpha_j$. Note that \odot denotes the element-wise product between the vectors, and in $e^{\mathbf{E}_{i,j}^y}$, the exponential function is applied to each dimension individually. Such distribution is also noted as the Bayesian flow distribution. Putting Eq. 10, Eq. 9 and Eq. 8 in Eq. 3, we could obtain the final training objective.

Continuous-time Training Objective The above discussion is based on that the number of latent variables is n , which is finite. It could also be possible to put $n \rightarrow \infty$, which could correspond to the continuous-time training. For convenience, we consider the value of timestep t in $[0, 1]$. In this case, $\beta(t) = \int_0^t \alpha(t)$. For flexibility, we could use different scheduler of $\alpha(t)$ for edge and node, i.e. $\alpha_e(t)$ and $\alpha_v(t)$. With the following notations as:

$$\begin{aligned} q(\mathbf{E}^{\theta_t} | \mathbf{E}, \mathbf{E}^{\theta_0}) &= \prod_{1 \leq i, j \leq n} q(\mathbf{E}_{i,j}^{\theta_t} | \mathbf{E}_{i,j}, \mathbf{E}_{i,j}^{\theta_0}) \\ q(\mathbf{V}^{\theta_t} | \mathbf{V}, \mathbf{V}^{\theta_0}) &= \prod_{1 \leq i \leq n} q(\mathbf{V}_i^{\theta_t} | \mathbf{V}_i, \mathbf{V}_i^{\theta_0}) \\ q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0}) &= q(\mathbf{E}^{\theta_t} | \mathbf{E}, \mathbf{E}^{\theta_0}) q(\mathbf{V}^{\theta_t} | \mathbf{V}, \mathbf{V}^{\theta_0}) \end{aligned} \quad (11)$$

Under the continuous-time steps, there is a very concise analytic form for the KL divergence between the sender distribution (Eq. 8) and the receiver distribution (Eq. 9) over

time step t as:

$$L(\mathbf{E}, t) = c_e \frac{1}{2} \alpha_e(t) \mathbb{E}_{q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0})} \sum_{1 \leq i, j \leq n} \|\mathbf{E}_{i,j} - \mathbf{E}_{i,j}[\phi(\mathbf{G}^{\theta_t}, t)]\|^2$$

$$L(\mathbf{V}, t) = c_v \frac{1}{2} \alpha_v(t) \mathbb{E}_{q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0})} \sum_{1 \leq i \leq n} \|\mathbf{V}_i - \mathbf{V}_i[\phi(\mathbf{G}^{\theta_t}, t)]\|^2$$

Combining the above terms, our final objective for GraphBFN is as:

$$L(\mathbf{G}) = L((\mathbf{E}, \mathbf{V})) = \mathbb{E}_{t \sim U(0,1)} L(\mathbf{E}, t) + L(\mathbf{V}, t) \quad (12)$$

The detailed training and sampling algorithms are illustrated in Appendix E. And with perturbation equivariant architecture for ϕ , e.g., Graph Transformer, we have following properties:

Proposition 4.2. *If ϕ is a permutation equivariant function, then for any permutation π , we have:*

- The objective in Eq. 12 is permutation invariant: $L(\mathbf{G}) = L((\mathbf{E}, \mathbf{V})) = L(\pi^T \mathbf{E} \pi, \pi^T \mathbf{V})$
- The density function, p_ϕ , implied by the generative process of GraphBFN is also permutation invariant, i.e., $p_\phi(\mathbf{E}, \mathbf{V}) = p_\phi(\pi^T \mathbf{E} \pi, \pi^T \mathbf{V})$

We leave the proof of the Proposition 4.2 in the Appendix J.

4.3. Improved Techniques for GraphBFN

Sampling with Adaptive Flowback The Bayesian Flow Networks conducted an autoregressive sampling procedure. For example, the t -th step vanilla sampling of BFN could be introduced as: first sample \hat{G}_{t-1} according to the network output distribution $\phi(\mathbf{G}^{\theta_{t-1}}, t-1)$; and assuming \hat{G} as the final generated targets to sample $\hat{G}_{t-1}^{y_t}$ based on the sender distribution, and then apply the Bayesian update function to get \mathbf{G}^{θ_t} . However, there exists a notable discrepancy between training and sampling: During training the noisy samples y are sampled based on real data; while at inference, the y are sampled based on model prediction which brings extra errors, especially for the initial steps where there lacks information in the model inputs. The error will always be kept in the inputs due to the autoregressive Bayesian update, which could potentially hurt the performance.

To the end, we introduce another sampling procedure to make the sampling focus on the current prediction by using the Bayesian Flow distribution:

$$p_\phi(\mathbf{G}^{\theta_t} | \mathbf{G}^{\theta_{t-1}}) = q(\mathbf{G}^{\theta_t} | \phi(\mathbf{G}^{\theta_{t-1}}, t-1), \mathbf{G}^{\theta_0}) \quad (13)$$

$q(\mathbf{G}^{\theta_t} | \phi(\mathbf{G}^{\theta_{t-1}}, t-1), \mathbf{G}^{\theta_0})$ here is the Bayesian flow distribution analogy to Eq. 10. We termed this approach as Flowback. Though the error accumulation is fixed, Flowback sampling could collapse due to that the network output distribution $\phi(\mathbf{G}^{\theta_t}, t)$ could hold low entropy practically. To balance exploration and exploitation, we combine the two sampling procedure and introduce the adaptive Flowback: with a predefined constant ϵ , for any time t , when $\|\phi(\mathbf{G}^{\theta_t}, t) - \phi(\mathbf{G}^{\theta_{t-1}}, t-1)\|^2 \geq \epsilon$, we conduct vanilla sampling; when $\|\phi(\mathbf{G}^{\theta_t}, t) - \phi(\mathbf{G}^{\theta_{t-1}}, t-1)\|^2 < \epsilon$, the Flowback (Eq. 13) is conducted. Empirically, we find the adaptive Flowback could be approximately implemented by setting a threshold over t to determine which kind of sampling to conduct to improve efficiency. The detailed sampling process is presented in Algorithm 2.

Accuracy Scheduling for Any Prior The accuracy scheduling is defined by the function $\beta(t) = \int_0^t \alpha(t)$, which represents the level the sender distribution corrupted the original information of the data samples. The original BFN paper (Graves et al., 2023) defines the accuracy scheduling as $\beta(t) = \beta(1)t^2$ and, correspondingly, the accuracy rate as $\alpha(t) = \beta(1)2t$. Such design relies on the intuition that, with a uniform prior, the expected entropy of θ could linearly decrease along the Bayesian flow distributions. The scheduling also implies the fact that the accuracy is 0 for the prior θ_0 as it is fully uninformative. However, in some scenarios of graph generation, an informative prior has been demonstrated to be effective (Vignac et al., 2022). To this end, we introduce a more general accuracy scheduling with a simple yet effective formulation:

$$\alpha(t) = a + 2(\beta(1) - a)t \quad (14)$$

$$\beta(t) = at + (\beta(1) - a)t^2 \quad (15)$$

Here, a is a hyperparameter. One degenerate case of Eq. 14 is when $a = \beta(1)$, which assigns a constant accuracy at each time. As shown in Fig. 3, our general accuracy scheduling exhibits more steady paces of information reconstruction. This scheduling could be compatible with low-entropy prior, i.e. assigning a nonzero accuracy for informative prior, and also have a faster input entropy decrease.

Conditioned Graph Features on Output Distributions

Previous literature (Beaini et al., 2021; Vignac et al., 2022) has demonstrated that incorporating extra features of input graphs can enhance the representation power of graph neural networks. For example, Vignac et al. (2022) computes various graph descriptors of the noisy graph and includes them in the input of the denoising network during both training and sampling. In GraphBFN, we could also calculate the graph descriptors of the probabilistic adjacency matrix for the inclusion of the interdependency modeling network, i.e. ϕ . We note that the extra graph feature is not necessarily used, and the choice is determined by the data property and

complexity. For the usage of extra features, we apply a different strategy compared to previous literature (Vignac et al., 2022) by conditioning the network output distribution. We denote the module for computing the extra features as $M(\cdot)$, for example, obtaining the eigenvalues or eigenvectors of the adjacency matrix. The computation of the output distribution with the spectral feature conditioning is then formulated as:

$$\phi(\mathbf{G}^{\theta_t}, t) = \phi(\mathbf{G}^{\theta_t}, t, M(\text{sg } \phi(\mathbf{G}^{\theta_t}, t))) \quad (16)$$

where sg is the stop gradient operator. The key motivation of the proposed method is related to the self condition as proposed in Chen et al. (2022), which helps to improve the inner consistency of the generated samples. Such techniques are only used in large graphs such as the **SBM** dataset.

5. Experiments

We empirically study the effectiveness of the proposed framework. We conduct extensive experiments on both the abstract graph datasets and 2D molecule datasets to compare the performance of GraphBFN against several competitive graph generation baselines, including diffusion/score matching based methods, such as Digress (Vignac et al., 2022), GruM (Jo et al., 2024), GBD (Liu et al., 2025), GDSS (Jo et al., 2022) and EDP-GNN (Niu et al., 2020); latent variable models, such as JT-VAE (Jin et al., 2018). GAN-based method SPECTRE (Martinkus et al., 2022); flow-based methods such as GraphAF (Shi et al., 2020), MoFlow (Zang & Wang, 2020) and GraphDF (Luo et al., 2021). Please refer to Appendix L for details about our experiment settings, the link to our implementation, and additional experiment results.

5.1. Abstract Graph Generation

The abstract graph generation experiments test the model’s ability to model complex graph topology. We conduct experiments on the benchmarks introduced in Vignac et al. (2022); Martinkus et al. (2022). It should be noted that we only need to model edge connections for the generation tasks of abstract graphs. The benchmark consists of two datasets: **Planar** with the constant number of 64 nodes per graph and **Stochastic Block Model (SBM)** with up to 200 nodes per graph. We evaluate the maximum mean discrepancy (MMD) of four graph statistics, *i.e.* degree (Deg.), clustering coefficient (Clus.), count of orbits with 4 nodes (Orbit), and the eigenvalues of the graph Laplacian (Spec.) between the generated graphs and test sets. Besides, we also report the percentage of valid, unique, and novel (V.U.N.) graphs, where a graph is considered valid if it satisfies the statistical properties of the **SBM** model or being planar and connected for the **Planar** dataset (Martinkus et al., 2022). The experiment result can be found in Tab. 1. GraphBFN

has shown superior or competitive performance compared to all the baselines on both datasets. Specifically, on the **Planar** dataset, GraphBFN obtains a V.U.N. value approaching the limits, demonstrating the advantages of modeling the discrete graph in the probabilistic matrix space. We include more details of visualizations and analysis in Appendix R and Appendix G.

5.2. 2D Molecule Generations

For the molecule datasets, we conduct experiments on 2D molecule datasets with various scales: **QM9** which includes 100k molecules with up to 9 heavy atoms in the training set; **ZINC250k** which contains 250,000 drug-like molecules with a maximum atom number of 38; **MOSEs** with 1,936,962 drug-like molecules whose atom number is up to 38. GraphBFN jointly models the node features \mathbf{V} and the graph topologies \mathbf{E} of the molecular graphs. Besides, for the experiments on **QM9**, we consider both the without-hydrogen version and the more challenging version where the hydrogen will be explicitly modeled. For **QM9** with explicit hydrogen, we evaluate GraphBFN’s performance with two additional metrics, the Atom Stability, *i.e.* the percentage of the atoms with valid valency, and Molecule Stability for the percentage of generated molecules whose atoms all have valid valency. For other datasets, we follow the evaluation of previous works (Jo et al., 2024; Vignac et al., 2022) and report the corresponding metrics, including FCD, SNN, Scaffold similarity(Scaf), NSPDK, Novelty, and Validity, etc.

The experiment results of **QM9** can be found in Tab. 2 and Tab. 4. It could be found that GraphBFN shows superior or competitive performance compared with the previous method. In the challenging setting of **QM9** with explicit hydrogen, the GraphBFN could beat the baselines with a large margin and approach the empirical limits in the sense of molecule stability, which justifies the effectiveness of jointly modeling the distribution of node features and edge feature under the probabilistic parameter space.

The empirical results of **ZINC250k** and **MOSEs** in Tab. 2 and Tab. 3 show GraphBFN’s consistently superior performance compared to the baselines, demonstrating its scalability to the large complex datasets. It should be noted that in Tab. 3, compared with other one-shot models, GraphBFN shows a slight decrease in novelty while maintaining a high uniqueness. One possible explanation lies in that the probability distribution learned with GraphBFN is so close to the training distribution that causes some overlap of the support. Despite this overfitting, the novelty still lies at a high level (89%). It also demonstrates the generalization ability of GraphBFN. We leave the visualization of generated samples in Appendix R.

Table 1: **Main results on general graph generation.** Best results are **bolded**, and the second-best results are underlined. GraphBFN consistently achieves either the best or competitive performance across all datasets.

	Planar					SBM				
	Synthetic, $ V = 64$					Synthetic, $44 \leq V \leq 187$				
	Deg. ↓	Clus. ↓	Orbit ↓	Spec. ↓	V.U.N. ↑	Deg. ↓	Clus. ↓	Orbit ↓	Spec. ↓	V.U.N. ↑
Training set	0.0002	0.0310	0.0005	0.0052	100.0	0.0008	0.0332	0.0255	0.0063	100.0
GraphRNN	0.0049	0.2779	1.2543	0.0459	0.0	0.0055	0.0584	0.0785	0.0065	5.0
GRAN	0.0007	0.0426	<u>0.0009</u>	0.0075	0.0	0.0113	0.0553	0.0540	0.0054	25.0
SPECTRE	<u>0.0005</u>	0.0785	0.0012	0.0112	25.0	0.0015	0.0521	0.0412	0.0056	52.5
EDP-GNN	0.0044	0.3187	1.4986	0.0813	0.0	0.0011	0.0552	0.0520	0.0070	35.0
GDSS	0.0041	0.2676	0.1720	0.0370	0.0	0.0212	0.0646	0.0894	0.0128	5.0
DiGress	0.0003	0.0372	<u>0.0009</u>	0.0106	75.0	0.0013	0.0498	<u>0.0434</u>	0.0400	74.0
GBD	0.0003	<u>0.0353</u>	<u>0.0135</u>	<u>0.0059</u>	<u>92.5</u>	0.0013	<u>0.0493</u>	0.0445	0.0047	75.0
GruM	<u>0.0005</u>	<u>0.0353</u>	<u>0.0009</u>	0.0062	90.0	<u>0.0007</u>	0.0492	0.0448	<u>0.0050</u>	<u>85.0</u>
GraphBFN	<u>0.0005</u>	0.0294	0.0002	0.0046	96.7	0.0005	0.0560	0.0370	0.0053	87.5

Table 2: **Main results on the QM9 with implicit hydrogens and ZINC250k.** Best results are **bolded**, and the second-best results are underlined. GraphBFN consistently achieves the best or competitive results across all datasets.

Method	QM9 ($ V \leq 9$)				ZINC250k ($ V \leq 38$)			
	Valid (%)↑	FCD↓	NSPDK↓	Scaf.↑	Valid (%)↑	FCD↓	NSPDK↓	Scaf.↑
Training set	100.0	0.0398	0.0001	0.9719	100.0	0.0615	0.0001	0.8395
MoFlow (Zang & Wang, 2020)	91.36	4.467	0.0169	0.1447	63.11	20.931	0.0455	0.0133
GraphAF (Shi et al., 2020)	74.43	5.625	0.0207	0.3046	68.47	16.023	0.0442	0.0672
GraphDF (Luo et al., 2021)	93.88	10.928	0.0636	0.0978	90.61	33.546	0.1770	0.0000
EDP-GNN (Niu et al., 2020)	47.52	2.680	0.0046	0.3270	82.97	16.737	0.0485	0.0000
GDSS (Jo et al., 2022)	95.72	2.900	0.0033	0.6983	97.01	14.656	0.0195	0.0467
DiGress (Vignac et al., 2022)	98.19	<u>0.095</u>	<u>0.0003</u>	0.9353	94.99	3.482	0.0021	0.4163
GBD (Liu et al., 2025)	99.88	0.093	0.0002	0.9510	97.87	<u>2.248</u>	0.0018	0.5042
GruM (Jo et al., 2024)	99.69	0.108	0.0002	<u>0.9449</u>	<u>98.65</u>	2.257	<u>0.0015</u>	<u>0.5299</u>
GraphBFN (Ours)	<u>99.73</u>	0.101	0.0002	0.9386	99.22	2.116	0.0013	0.5304

5.3. Ablation Studies

Sampling Steps We conduct extensive ablation studies of GraphBFN. We first consider GraphBFN’s property of sampling any steps and then examine the improved techniques we claimed in Sec. 4.3. One appealing property of GraphBFN is that with continuous-time objectives, the GraphBFN could sample with any steps. We conduct an ablation study of sampling steps on the dataset of **QM9** with *explicit* hydrogen and **Planar** in Tab. 4. It could be found that the GraphBFN with as few as 25 and 50 sampling steps could be on par with diffusion models with 500 and 1000 steps. This result implies a more than $10\times$ increase in sampling efficiency.

Adaptive Flowback Sampling We verified the effectiveness of the Adaptive Flowback sampling on the abstract graph dataset **Planar**. The experiment results can be found in Tab. 5. We note that with vanilla sampling (i.e., $\epsilon = 0.0$, where the vanilla BFN sampling is used throughout the sampling process), the performance of V.U.N. is only 14.8%. With full Flowback, we could obtain a competitive performance with V.U.N. consistently outperforming the baseline, while the distribution-level metrics are still sub-optimal due to the collapse phenomenon. With adaptive Flowback, we could obtain the best quality-diversity tradeoff which makes adaptive Flowback a more appealing strategy.

Due to space limits, we defer the additional ablation studies on **Time Scheduler** and **Output Conditioned Feature** to Appendix H. Furthermore, we conducted an empirical

Table 3: **Main results on MOSES.** We **bold** the best and underline the second-best results across all one-shot methods. GraphBFN outperforms other one-shot baselines on most metrics.

Model	Class	Val \uparrow	Unique \uparrow	Novel \uparrow	Filters \uparrow	FCD \downarrow	SNN \uparrow	Scaf \uparrow
VAE	SMILES	97.7	99.8	69.5	99.7	0.57	0.58	5.9
JT-VAE	Fragment	100	100	99.9	97.8	1.00	0.53	10.0
GraphINVENT	Autoreg.	96.4	99.8	–	95.0	1.22	0.54	12.7
ConGress		83.4	<u>99.9</u>	96.4	94.8	1.48	0.50	16.4
DiGress	One-shot	<u>85.7</u>	100	<u>95.0</u>	<u>97.1</u>	<u>1.19</u>	<u>0.52</u>	<u>14.8</u>
GraphBFN		88.5	99.8	89.0	98.3	1.07	0.59	10.0

Table 4: **Ablation study** on GraphBFN’s ability to tradeoff sampling *quality* and *efficiency*. Our experiment on QM9 with *explicit* hydrogen and Planar demonstrates that GraphBFN can produce equally good samples as the previously best models with only 10% of their sampling steps (500 for QM9 and 1000 for Planar). A few entries are left as blank because the corresponding experiment results are not implemented in the original works (Vignac et al., 2022; Jo et al., 2024).

% Sample Steps	QM9 with H $ V \leq 29$				Planar $ V \leq 64$
	Valid (%) \uparrow	Unique (%) \uparrow	Atom Stab. (%) \uparrow	Mol. Stab. (%) \uparrow	V.U.N. (%) \uparrow
100	99.2	94.9	99.4	94.7	96.7
20	98.8	95.6	99.1	91.4	93.8
<u>10</u>	97.6	95.9	98.7	87.4	71.1
5	94.2	96.4	98.1	79.7	14.8
100 (DiGress)	95.4	97.6	98.1	79.8	75.0
100 (GruM)	–	–	–	–	90.0

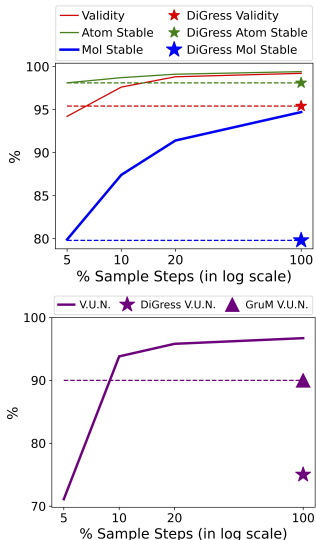


Table 5: **Ablation study** on the impact of the adaptive Flowback threshold ϵ on generation *quality* and *diversity*. According to our experiment on Planar Graph, increasing the use of Flowback sampling strategy always yields better sampling quality (as reflected in the strictly increasing V.U.N.). However, overusing it might compromise sample diversity (as reflected in the sub-optimal graph statistics of $\epsilon = 1.0$ in comparison to $\epsilon = 0.5$).

ϵ	Planar				
	Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Spec. \downarrow	V.U.N. \downarrow
Training Set	0.0002	0.0310	0.0005	0.0052	100.0
Flowback (1.0)	0.0015	0.0703	0.0017	0.0047	98.2
Adaptive Flowback (0.5)	0.0005	0.0294	0.0002	0.0046	96.7
Vanilla (0.0)	0.0051	0.1793	0.0278	0.0099	14.1

study that quantifies the benefits of smooth graph topology transformation, with details presented in Appendix G.2.

6. Conclusions

We introduce GraphBFN for the generative modeling of discrete graphs. GraphBFN is motivated by the intuitive idea of smoothly decomposing and modeling the information of the graph topology by introducing a probabilistic

matrix space. We first identify the train-test discrepancy issues in the vanilla formulation of Bayesian Flow Networks and propose several innovations under the framework. These appealing properties and improved techniques help GraphBFN achieve consistently superior performance on extensive graph generation benchmarks. Furthermore, GraphBFN can maintain decent sample quality while requiring much fewer sampling steps, thus largely improving the sample efficiency.

Acknowledgments

This work is supported the Natural Science Foundation of China (Grant No. 62376133) and sponsored by Beijing Nova Program (20240484682) and the Wuxi Research Institute of Applied Technologies, Tsinghua University (20242001120).

Impact Statement

We confirm that our work complies with the ICML Publication Ethics, and we have carefully considered potential ethics concerns relating to the development and use of our proposed method GraphBFN for generative modeling of graph-structured data.

Our model is designed for general graph generation tasks and does not involve the use of sensitive personal data. However, we acknowledge that generative models, when applied to specific domains, such as molecular generation and social networking modeling, may inadvertently create synthetic data that could mimic private structures and details. Thus, we encourage users to apply our model in compliance with relevant privacy regulations and to critically evaluate its outputs.

We are committed to full transparency in our research. All datasets used in our experiments are publicly available and documented in detail to ensure reproducibility and compliance with legal standards. The methods and assumptions made in our research are clearly documented, and we provide substantial reproduction details in Appendix L. Furthermore, there is no conflict of interest, financial or otherwise, that could influence the development or presentation of this work.

Based on these considerations, we do not anticipate any violations of the ICML Publication Ethics through the development or application of this model. However, we emphasize again that GraphBFN should not be used for malicious purposes, such as generating misleading data or creating harmful content.

References

- Athreya, A., Fishkind, D. E., Tang, M., Priebe, C. E., Park, Y., Vogelstein, J. T., Levin, K., Lyzinski, V., Qin, Y., and Sussman, D. L. Statistical inference on random dot product graphs: a survey. *Journal of Machine Learning Research*, 18(226):1–92, 2018.
- Austin, J., Johnson, D. D., Ho, J., Tarlow, D., and Van Den Berg, R. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.
- Beaini, D., Passaro, S., Létourneau, V., Hamilton, W., Corso, G., and Liò, P. Directional graph networks. In *International Conference on Machine Learning*, pp. 748–758. PMLR, 2021.
- Chen, T., Zhang, R., and Hinton, G. Analog bits: Generating discrete data using diffusion models with self-conditioning. *arXiv preprint arXiv:2208.04202*, 2022.
- Chen, X., He, J., Han, X., and Liu, L.-P. Efficient and degree-guided graph generation via discrete diffusion modeling, 2023. URL <https://arxiv.org/abs/2305.04111>.
- Chung, F. *Spectral Graph Theory*. Number no. 92 in CBMS Regional Conference Series. Conference Board of the Mathematical Sciences, 2006. ISBN 9780821889367. URL https://books.google.com/books?id=YUc38_MCuhAC.
- Chung, F. and Radcliffe, M. On the spectra of general random graphs. *Electr. J. Comb.*, 18, 10 2011. doi: 10.37236/702.
- De Cao, N. and Kipf, T. Molgan: An implicit generative model for small molecular graphs. *arXiv preprint arXiv:1805.11973*, 2018.
- Graves, A., Srivastava, R. K., Atkinson, T., and Gomez, F. Bayesian flow networks. *arXiv preprint arXiv:2308.07037*, 2023.
- Grone, R., Merris, R., and Sunder, V. The laplacian spectrum of a graph. *SIAM Journal on matrix analysis and applications*, 11(2):218–238, 1990.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020.
- Ingraham, J., Garg, V., Barzilay, R., and Jaakkola, T. Generative models for graph-based protein design. *Advances in neural information processing systems*, 32, 2019.
- Jiang, J. *AN INTRODUCTION TO SPECTRAL GRAPH THEORY*. The University of Chicago, 2012. URL <https://math.uchicago.edu/~may/REU2012/REUPapers/JiangJ.pdf>.
- Jin, W., Barzilay, R., and Jaakkola, T. Junction tree variational autoencoder for molecular graph generation. In *International conference on machine learning*, pp. 2323–2332. PMLR, 2018.
- Jin, W., Barzilay, R., and Jaakkola, T. Hierarchical generation of molecular graphs using structural motifs. In *International conference on machine learning*, pp. 4839–4848. PMLR, 2020.

- Jo, J., Lee, S., and Hwang, S. J. Score-based generative modeling of graphs via the system of stochastic differential equations. *arXiv preprint arXiv:2202.02514*, 2022.
- Jo, J., Kim, D., and Hwang, S. J. Graph generation with diffusion mixture. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=cZTFxktg23>.
- Karami, M. Higen: Hierarchical graph generative networks. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=KNvubydSB5>.
- Kong, L., Cui, J., Sun, H., Zhuang, Y., Prakash, B. A., and Zhang, C. Autoregressive diffusion model for graph generation, 2023. URL <https://arxiv.org/abs/2307.08849>.
- Li, X. L., Thickstun, J., Gulrajani, I., Liang, P., and Hashimoto, T. Diffusion-LM improves controllable text generation. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=3s9IrEsjLyk>.
- Liu, X., He, Y., Chen, B., and Zhou, M. Advancing graph generation through beta diffusion. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=x1An5a3U9I>.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization, 2019.
- Luo, T., Mo, Z., and Pan, S. J. Fast graph generation via spectral diffusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Luo, Y., Yan, K., and Ji, S. Graphdf: A discrete flow model for molecular graph generation. In *International Conference on Machine Learning*, pp. 7192–7203. PMLR, 2021.
- Martinkus, K., Loukas, A., Perraudin, N., and Wattenhofer, R. Spectre: Spectral conditioning helps to overcome the expressivity limits of one-shot graph generators. In *International Conference on Machine Learning*, pp. 15159–15179. PMLR, 2022.
- Niu, C., Song, Y., Song, J., Zhao, S., Grover, A., and Ermon, S. Permutation invariant graph generation via score-based generative modeling. In *International Conference on Artificial Intelligence and Statistics*, pp. 4474–4484. PMLR, 2020.
- Qin, Y., Vignac, C., and Frossard, P. Sparse training of discrete diffusion models for graph generation, 2023.
- Qu, Y., Qiu, K., Song, Y., Gong, J., Han, J., Zheng, M., Zhou, H., and Ma, W.-Y. Molcraft: Structure-based drug design in continuous parameter space. *arXiv preprint arXiv:2404.12141*, 2024.
- Shi, C., Xu, M., Zhu, Z., Zhang, W., Zhang, M., and Tang, J. Graphaf: a flow-based autoregressive model for molecular graph generation. *arXiv preprint arXiv:2001.09382*, 2020.
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, pp. 11918–11930, 2019.
- Song, Y., Gong, J., Zhou, H., Zheng, M., Liu, J., and Ma, W.-Y. Unified generative modeling of 3d molecules with bayesian flow networks. In *The Twelfth International Conference on Learning Representations*, 2023.
- Vignac, C., Krawczuk, I., Siraudin, A., Wang, B., Cevher, V., and Frossard, P. Digress: Discrete denoising diffusion for graph generation. *arXiv preprint arXiv:2209.14734*, 2022.
- Wyss-Gallifent, J. *Graph Theory*. The University of Maryland, 2021. URL https://www.math.umd.edu/~immortal/MATH401/book/ch_graph_theory.pdf.
- You, J., Liu, B., Ying, Z., Pande, V., and Leskovec, J. Graph convolutional policy network for goal-directed molecular graph generation. In *Advances in neural information processing systems*, pp. 6410–6421, 2018.
- Yu, J. J. Q. and Gu, J. Real-time traffic speed estimation with graph convolutional generative autoencoder. *IEEE Transactions on Intelligent Transportation Systems*, 20(10):3940–3951, 2019.
- Zang, C. and Wang, F. Moflow: an invertible flow model for generating molecular graphs. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 617–626, 2020.
- Zeng, X., Vahdat, A., Williams, F., Gojcic, Z., Litany, O., Fidler, S., and Kreis, K. LION: Latent point diffusion models for 3d shape generation. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=tHK5ntjp-5K>.
- Zhou, M., Chen, T., Wang, Z., and Zheng, H. Beta diffusion. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=zTS1m4nmlH>.

A. Detailed Introduction to Bayesian Flow Networks

In this section, we introduce the basic components of Bayesian Flow Networks (BFNs) (Graves et al., 2023). The graphical model of BFN could be expressed in the form of a latent variable model whose prior is non-Markov autoregressive, distinguishing it from the diffusion models that have Markov priors. That is, for learning the probability distribution p_ϕ over sample space \mathbf{x} , a series of noisy versions $\mathbf{y}_{\{1:n\}} = \langle \mathbf{y}_1, \dots, \mathbf{y}_n \rangle$ of \mathbf{x} are introduced as latent variables. Then the model seeks to optimize the *variational lower bound* of the likelihood:

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{\mathbf{y}_{\{1:n\}} \sim q} \left[\log \frac{p_\phi(\mathbf{x} | \mathbf{y}_{\{1:n\}}) p_\phi(\mathbf{y}_{\{1:n\}})}{q(\mathbf{y}_{\{1:n\}} | \mathbf{x})} \right] = -D_{KL}(q \| p_\phi(\mathbf{y}_{\{1:n\}})) + \mathbb{E}_{\mathbf{y}_{\{1:n\}} \sim q} \log [p_\phi(\mathbf{x} | \mathbf{y}_{\{1:n\}})] \quad (17)$$

The variational distribution q is defined by so-called sender distribution $p_S(\mathbf{y}_i | \mathbf{x}; \alpha_i)$. The sender distribution could be understood as adding noise to the data sample \mathbf{x} with a predefined accuracy α_i . The variational distribution is the fixed mean-field combination of sender distribution by $q(\mathbf{y}_{\{1:n\}} | \mathbf{x}) = \prod_{i=1}^n p_S(\mathbf{y}_i | \mathbf{x}; \alpha_i)$.

For the autoregressive prior $p_\phi(\mathbf{y}_1, \dots, \mathbf{y}_n)$, the joint distribution could be extended with Bayesian decomposition as:

$$p_\phi(\mathbf{y}_1, \dots, \mathbf{y}_n) = p_\phi(\mathbf{y}_1) \prod_{t=2}^n p_\phi(\mathbf{y}_t | \mathbf{y}_{\{1:t-1\}}) \quad (18)$$

Note in BFN, the \mathbf{y} is non-Markov contrast to diffusion models. In BFN, a new variable is introduced based on a fix-form functional transformation f , e.g. a weighted average function, as $\theta_t = f(\mathbf{y}_{\{1:t\}})$ with $t > 1$ and θ_0 as a constant. We elaborate more on the definition of t , which is implied by successive Bayesian updates as $\theta_t = h(\theta_{t-1}, \mathbf{y}_t, \alpha_t)$. Given the constant prior θ_0 and a series noisy versions $\mathbf{y}_{\{1:n\}}$ of \mathbf{x} , we could obtain $\theta_{\{1:n\}}$ through a sequence of Bayesian updates. The only randomness of the above process comes from \mathbf{y} , hence it could be summarized as a deterministic function f over \mathbf{y} . θ has the statistical interpretation as the parameter of the factorized distribution over the sample spaces by $p(\mathbf{x} | \theta) = \prod_{d=1}^D p(\mathbf{x}^{(d)} | \theta^{(d)})$. The conditional distribution $p_\phi(\mathbf{y}_t | \mathbf{y}_{\{1:t-1\}})$ is defined as $p_\phi(\mathbf{y}_t | \mathbf{y}_{\{1:t-1\}}) = p_\phi(\mathbf{y}_t | \theta_{t-1} = f(\mathbf{y}_{\{1:t-1\}}))$. Thus, Eq. 18 could be rewritten as: $p_\phi(\mathbf{y}_1, \dots, \mathbf{y}_n) = \prod_{t=1}^n p_\phi(\mathbf{y}_t | \theta_{t-1})$. We then clarify $p_\phi(\mathbf{y}_t | \theta_{t-1})$. To obtain this term, we will first input θ_{t-1} into the neural network Φ to capture the interdependency between different dimensions. Hence the distribution implied by θ_{t-1} , the neural network input, is also named as *input distribution*. The output of the neural networks denoted as $\hat{\theta}_{t-1} = \Phi(\theta_{t-1})$ still lies in the parameter space and the corresponding distribution is usually referred to as *output distribution*. Combining the formulation and accuracy α_t of *sender distribution* p_S , the distribution $p_\phi(\mathbf{y}_t | \theta_{t-1})$ is defined as

$$p_\phi(\mathbf{y}_t | \theta_{t-1}) = \mathbb{E}_{p_O(\hat{\mathbf{x}} | \Phi(\theta_{t-1}))} p_S(\mathbf{y}_t | \hat{\mathbf{x}}; \alpha_t) \quad (19)$$

Note the $p_\phi(\mathbf{y}_t | \theta_{t-1})$ is also formulated as *receiver distribution* $p_R(\mathbf{y}_t | \theta_{t-1}; \Phi, \alpha_t)$.

Integrate Eq. 19 to Eq. 17, with $p_\phi(\mathbf{x} | \mathbf{y}_{\{1:n\}}) = p_O(\mathbf{x} | \theta_n)$, we could get the objective as:

$$\mathcal{L}_{\text{ELBO}} = \sum_{t=1}^n D_{KL}(p_S(t) \| p_R(t)) + \log p_O(\mathbf{x} | \theta_n) \quad (20)$$

Here $p_S(t)$ and $p_R(t)$ is abbreviation for $p_S(\mathbf{y}_t | \mathbf{x}, \alpha_t)$ and $p_R(\mathbf{y}_t | \theta_{t-1}; \Phi, \alpha_t)$. To conduct simulation-free training on the ELBO, Graves et al. (2023) proposed to use the so-called Bayesian flow distribution $p_F(\theta_{\{0:n\}} | \mathbf{x})$ formulated as $p_F(\theta_t | \mathbf{x}) = \prod_{k=1}^t \mathbb{E}_{p_S(\mathbf{y}_k | \mathbf{x}; \alpha_k)} f(\mathbf{y}_{\{1:t\}})$, which allows for teacher-forcing-like training and generally has an analytical formulation.

B. Spectral Graph Theory

Following previous literature (Vignac et al., 2022; Martinkus et al., 2022), we consider the definition of graph spectrum as the eigenvalues of the Laplacian matrix. For a graph \mathbf{G} with n nodes, the graph Laplacian matrix is defined as $\mathbf{L} = \mathbf{D} - \mathbf{A}$, where \mathbf{A} is the $n \times n$ adjacency matrix with $\mathbf{A}_{i,j} = 1$ if the node v_i and node v_j are connected and $\mathbf{A}_{i,j} = 0$ otherwise. \mathbf{D} is the degree matrix, a diagonal matrix defined as $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$. Here d_i denotes the degree of the node v_i .

as $d_i = \sum_{j=1}^n \mathbf{A}_{i,j}$. Besides, the normalized Laplacian Matrix \mathbf{L}_{norm} is defined as $\mathbf{L}_{\text{norm}} = \mathbf{I} - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$. For the symmetric positive semi-definite matrix \mathbf{L} or \mathbf{L}_{norm} , the eigen-decomposition is $\mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$ is an orthogonal matrix for the eigenvectors and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ denotes the collection of eigenvalues. The eigenvalues and eigenvectors of the graph Laplacian imply some key properties of the graph topology structure, such as connectivity, clusterability, and distance between nodes (Grone et al., 1990).

Spectral Gap and Graph Clusterability In this paragraph, we elaborate on the relationship between graph clusterability and the spectral gap that we’ve used to quantify graphs’ topological information in Figure 1.

The spectral gap is computed by taking the difference between the first two smallest eigenvalues of graph matrices. Its definition varies slightly depending on whether the eigenvalues are taken from a graph’s adjacency matrix or the Laplacian matrix. Nevertheless, regardless of which matrix it takes, the correlation holds that a smaller spectral gap implies poor graph connectivity and good clusterability. In our paper, we stick to the definition according to the *Laplacian Matrix*.

The Laplacian matrix \mathbf{L} is symmetric and has non-negative and real eigenvalues. The multiplicity of 0 in the eigenvalues of \mathbf{L} is equivalent to the number of connected components in the graph (Jiang, 2012).

Let $\lambda_1, \dots, \lambda_n$ denote the eigenvalues of L in the increasing order. The spectral gap in this scenario is defined as $\lambda_2 - \lambda_1$ (Martinkus et al., 2022). Since the smallest eigenvalue is always 0 (as the graph must have at least one connected component), the spectral gap is nothing but the second smallest eigenvalue of the Laplacian.

The second smallest eigenvalue is known as the Fiedler Value or the algebraic connectivity of a graph, and it reflects the connectivity of a graph (Wyss-Gallifent, 2021). Specifically, given a fixed graph volume, the second smallest eigenvalue lowerbonds the diameter of a graph (Chung, 2006), which is a connectivity metric that measures the maximum distance between a pair of nodes in the graph. Therefore, if the second eigenvalue is small, the graph will contain nodes that are separated far apart from each other, implying loose connections and more clusters.

To gain more intuition for why a small spectral gap corresponds to a highly clustered graph, we can think about the extreme case where the spectral gap is 0. A zero spectral gap means that the second smallest eigenvalue of the Laplacian is also 0, implying the multiplicity of zero eigenvalues is at least 2. This means that the graph has at least 2 connected components (i.e. at least 2 clusters with no bridge between them).

C. Spectral Features in the Space of Probabilistic Adjacency Matrices

Clarification on Probabilistic Adjacency Matrix Space Instead of focusing on the properties customized on the graph domain, we generalize the graph featured by its adjacency matrix to the more general matrix domains. The eigenvalues or eigenvectors of the matrices could be effective tools to reflect the essential properties. Then there are the following facts that are proved in (Chung & Radcliffe, 2011):

Theorem C.1. (Properties and Bounds of PAM) For D^p defined as the $\mathbb{E}_{A \sim A^p} D$, i.e. the expectation of the Degree matrix of the graph whose adjacency matrix A sampled from A^p . Let δ be the minimum expected degree of G , and $L = I - D^{-1/2} A D^{-1/2}$ the (normalized) Laplacian matrix and similarly definition $L^p = I - (D^p)^{-1/2} A^p (D^p)^{-1/2}$. Choose $\epsilon > 0$:

- for $\Delta > \frac{4}{9} \ln(2n/\epsilon)$. Then with probability at least $1 - \epsilon$, for n sufficiently large, the eigenvalues of A and A^p satisfy

$$\forall 1 \leq j \leq n, |\lambda_j(A) - \lambda_j(A^p)| \leq \sqrt{4\Delta \ln(2n/\epsilon)}$$

- There exists a constant $k = k(\epsilon)$ such that if $\delta > k \ln n$, then with probability at least $1 - \epsilon$, the eigenvalues of L and L^p satisfy

$$\forall 1 \leq j \leq n, |\lambda_j(L) - \lambda_j(L^p)| \leq 3\sqrt{\frac{3 \ln(4n/\epsilon)}{\delta}}$$

The above theorem helps us connect the probabilistic matrix space to the graph space with theoretical justification. In particular, the Laplacian matrix could also get a similar definition in the probabilistic space whose eigenvalues could relate to the graph spectrum of sampled discrete graphs.

D. Beyond Simple Undirected Graph

Our experiment shows that GraphBFN could handle the generation of undirected graphs. Furthermore, it could be naturally extended to other types of graphs, and we elaborate on how it could be adapted to them in the following.

Directed graphs Compared to undirected graphs, directed graphs differ only in that they do not need the adjacency matrix to be symmetric. In our current implementation, we enforce the symmetry on the network input \mathbf{E}^{θ_t} by the operation $\mathbf{E}^{\theta_t} = (\mathbf{E}^{\theta_t} + (\mathbf{E}^{\theta_t})^T)/2$, and on the final sample by $\mathbf{E}_{final} = (\mathbf{E}_{final} + \mathbf{E}_{final}^T)/2$. To study the directed graphs, we can simply remove these operations.

Weighted graphs For weighted graphs where each edge in the graph is assigned a numerical value, the GraphBFN could still be applied by normalizing the numerical value to the $[0,1]$ space. In this way, we could obtain a reformatted adjacency matrix, where an entry is set as the normalized numerical value if the corresponding edge exists and 0 otherwise. Then, we could apply the same objective of GraphBFN to such samples. It should be noted that the weighted graph is a hybrid of continuous and discrete variables, we could also use a hybrid GraphBFN for modeling, *i.e.* continuous BFN for modeling continuous variables and discrete for other discrete types like node and edge types.

Attributed graphs When presenting our methodology, especially when we introduce the probabilistic adjacency matrix, we assume the edge type to be binary (existence and non-existence) for the simplicity of presentation. However, it could be noted that the GraphBFN framework could be easily extended to the attributed graphs where the node and edge have categorical attributes, *e.g.*, more than two node/edge types. Furthermore, BFN’s effectiveness on this general type of graph has been empirically verified in our experiments on molecule generation, where the graphical representations of molecules are attributed graphs. Lastly, we could still obtain the topological information under this setting by accumulating the probability of all categories other than the no-edge category to build the probabilistic adjacency matrix space.

E. Algorithms

Algorithm 1 Discrete Variable Bayesian Flow

```

1: Require:  $K \in \mathbb{R}$ 
    $K$  is the possible state number of a discrete variable
2: function BayesianFlow( $\mathbf{x} \in \mathbb{R}^{D \times K}, \beta \in \mathbb{R}$ )
3:    $\mathbf{y} \sim \mathcal{N}(\beta(K\mathbf{x} - \mathbf{1}), \beta K \mathbf{I})$ 
4:    $\boldsymbol{\theta} \leftarrow \text{softmax}(\mathbf{y})$ 
5:   Return  $\boldsymbol{\theta}$ 
6: end function

```

Algorithm 3 Training procedure for one step

```

1: Require: function  $\beta_v(t) = \int_0^t \alpha_v(t), \beta_e(t) = \int_0^t \alpha_e(t), c_v, c_e \in \mathbb{R}$ , an NeuralNetwork  $\phi(\cdot)$ 
2: Input:  $\mathbf{V} \in \mathbb{R}^{n \times c_v}, \mathbf{E} \in \mathbb{R}^{n \times n \times c_e}, t \in [0, 1]$ 
3:  $\mathbf{V}^{\theta_t} \leftarrow \text{BayesianFlow}(\mathbf{V}, t)$ 
4:  $\mathbf{E}^{\theta_t} \leftarrow \text{BayesianFlow}(\mathbf{E}, t)$ 
5:  $\mathbf{V}[\phi(\mathbf{G}^{\theta_t}, t)], \mathbf{E}[\phi(\mathbf{G}^{\theta_t}, t)] \leftarrow \phi(\mathbf{G}^{\theta_t} = (\mathbf{V}^{\theta_t}, \mathbf{E}^{\theta_t}), t)$ 
6:  $L = c_e \frac{1}{2} \alpha_e(t) \sum_{1 \leq i, j \leq n} \|\mathbf{E}_{i,j} - \mathbf{E}_{i,j}[\phi(\mathbf{G}^{\theta_t}, t)]\|^2 + c_v \frac{1}{2} \alpha_v(t) \sum_{1 \leq i \leq n} \|\mathbf{V}_i - \mathbf{V}_i[\phi(\mathbf{G}^{\theta_t}, t)]\|^2$ 
7: Minimize loss  $L$ 

```

F. Improved Techniques for GraphBFN

We introduced a more general accuracy scheduling with a simple yet effective formulation as in Eq. 14, we further suppliment a comprehensive visualization of the entropy change of the input distribution when a non-uniform prior is used.

As shown in Fig. 3, Our general accuracy scheduling exhibits more steady paces of information reconstruction. This

Algorithm 2 Sampling procedure with Adaptive Flowback

```

1: Require:  $T, n \in \mathbb{Z}^+, c_v, c_e \in \mathbb{R}, \phi(\cdot), \beta_v(\cdot), \beta_e(\cdot), \alpha_v(\cdot), \alpha_e(\cdot), \epsilon$ 
   #  $T$  denotes the sampling step,  $N$  is the number of nodes,  $c_v, c_e$  is the possible state of a node and edge variable,  $\phi(\cdot)$  is
   the optimized neural network.  $\beta$  and  $\alpha$  are the corresponding scheduler functions.  $\epsilon$  is the predefined constant
2: Initialize:  $\mathbf{V}^\theta, \mathbf{E}^\theta$  with  $\mathbf{V}^{\theta_0}, \mathbf{E}^{\theta_0}$ 
3: for  $i = 1$  to  $T$  do
4:    $t \leftarrow \frac{i-1}{T}$ 
5:    $\mathbf{V}[\phi(\mathbf{G}^{\theta_t}, t), \mathbf{E}[\phi(\mathbf{G}^{\theta_t}, t)]] \leftarrow \phi(\mathbf{G}^{\theta_t} = (\mathbf{V}^{\theta_t}, \mathbf{E}^{\theta_t}), t)$ 
6:    $\mathbf{V}, \mathbf{E} \sim \text{Categorical}(\mathbf{V}[\phi(\mathbf{G}^{\theta_t}, t)], \text{Categorical}(\mathbf{E}[\phi(\mathbf{G}^{\theta_t}, t)]))$ 
7:    $t_n = t + \frac{1}{T}, t_p = t - \frac{1}{T}$ 
8:   if  $i > 2$  and  $\|\phi(\mathbf{G}^{\theta_t}, t) - \phi(\mathbf{G}^{\theta_{t_p}}, t_p)\|^2 \geq \epsilon$  then
9:      $\mathbf{V}^{y_{t_n}}, \mathbf{E}^{y_{t_n}} \sim q(\mathbf{V}^{y_{t_n}} | \mathbf{V}, \alpha_v(t_n)), q(\mathbf{E}^{y_{t_n}} | \mathbf{E}, \alpha_e(t_n))$ 
10:     $\mathbf{V}^{\theta_{t_n}}, \mathbf{E}^{\theta_{t_n}} \leftarrow \text{Bayesian Update}(\mathbf{V}^{y_{t_n}}, \mathbf{E}^{y_{t_n}}, \mathbf{V}^{\theta_t}, \mathbf{E}^{\theta_t})$  as Eq. 7
11:   else
12:      $\mathbf{V}^{\theta_{t_n}} \leftarrow \text{BayesianFlow}(\mathbf{V}[\phi(\mathbf{G}^{\theta_t}, t)], \beta_v(t_n))$ 
13:      $\mathbf{E}^{\theta_{t_n}} \leftarrow \text{BayesianFlow}(\mathbf{E}[\phi(\mathbf{G}^{\theta_t}, t)], \beta_e(t_n))$ 
14:   end if
15: end for
16:  $\mathbf{V}[\phi(\mathbf{G}^{\theta_1}, 1)], \mathbf{E}[\phi(\mathbf{G}^{\theta_1}, 1)] \leftarrow \phi(\mathbf{G}^{\theta_1} = (\mathbf{V}^{\theta_1}, \mathbf{E}^{\theta_1}), 1)$ 
17:  $\mathbf{V}, \mathbf{E} \sim \text{Categorical}(\mathbf{V}[\phi(\mathbf{G}^{\theta_1}, 1)], \text{Categorical}(\mathbf{E}[\phi(\mathbf{G}^{\theta_1}, 1)]))$ 
18: Return:  $\mathbf{V}, \mathbf{E}$ 

```

scheduling could be compatible with low-entropy prior, *i.e.* assigning a nonzero accuracy for informative prior, and also have a faster input entropy decrease.

G. Analysis of the Generation Process

G.1. Transformation of Graph Topology

GraphBFN and the diffusion-based models share the same advantage of gradually recovering the graph topological information in the generation process. However, diffusion-based models experience more fluctuations in recovering graph topology. The instability arises because the denoising steps in these models are conditioned on the atomic samples G^t represented in binary adjacency matrices, and such samples could involve significant noises and randomness. In contrast, in GraphBFN, the intermediate steps are conditioned on the input distribution \mathbf{G}^{θ_t} represented in the probabilistic adjacency matrix, which changes smoothly across the Bayesian update at each time step.

Experiment Setup To verify this discrepancy, we train GraphBFN and the diffusion-based models, Digress and Grums, to overfit a 2-cluster community graph with 20 nodes. The topological information of the graph is captured by the spectral gap, as it reflects clusterability—the most important topological feature of our overfitting sample. Recall that the spectral gap is calculated by taking the difference between the 1st and 2nd eigenvalues of the Laplacian for diffusion-based models and the probabilistic Laplacian for GraphBFN. The spectral gaps are measured on the intermediate noisy samples E^t for DiGress and the input distribution \mathbf{G}^{θ_t} for GraphBFN, and, for GruM, we use the quantized E_t .

Results We plot the changes of the spectral gaps in Figure. 4. In this figure, we see that GruM transforms information more smoothly than DiGress but not as smoothly as GraphBFN. Such observation is further confirmed by quantifying smoothness using the following metrics three metrics. The results of these metrics are shown in Table. 6.

- **Mean Absolute Second Derivative(MASD):** MASD calculates the second derivative of the data, the 2nd derivative is then averaged out, and the equation is $S_{MASD} = \frac{\sum_{i=1}^{n-2} \left| \frac{d^2}{dx^2} (data_i) \right|}{n-2}$.
- **Mean Local Standard Deviation(MLSTD):** MLSTD computes standard deviation in a small window, where, in our case, the window is set to 5 and then averages the local standard deviation. Intuitively, MLSTD captures relatively long-range smoothness in the window.

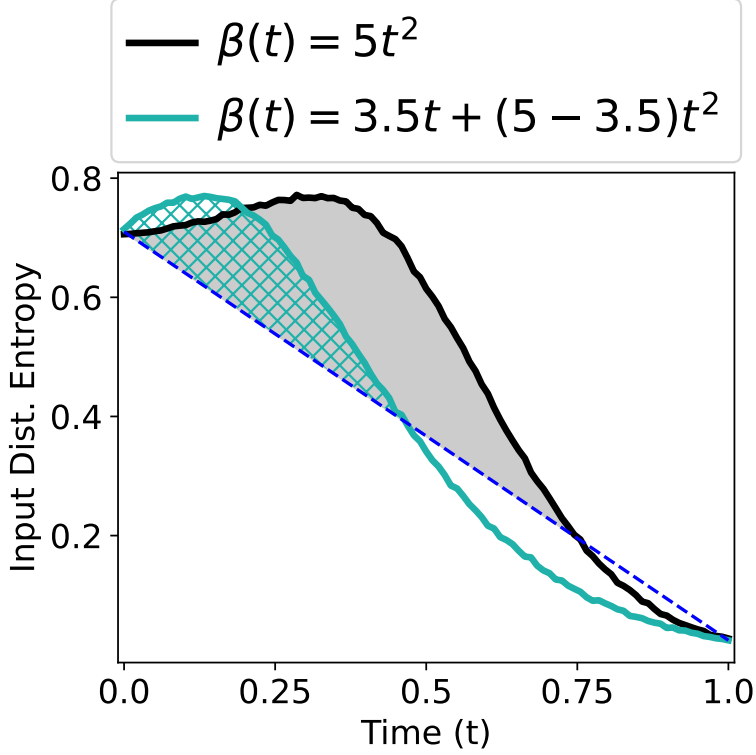


Figure 3: The entropy change of the input distribution \mathbf{G}^{θ_t} along the Bayesian flow equipped with a non-uniform prior.

Table 6: The smoothness of spectral gap transformation in the generation process of GraphBFN and the diffusion-based baseline models

Model	MASD(1e-2)↓	MLSTD(1e-2)↓	MVar(1e-2) ↓
DiGress	1.14	0.78	0.65
GruM	1.43	0.79	0.77
GraphBFN	0.39	0.25	0.26

- **Mean Variation(MVar):** MVar calculates the difference between neighboring data points, and we report the average. It emphasizes more on short-range smoothness.

G.2. Advantages of “Smooth” Transformation

We further conduct an empirical study that quantized the benefit of smooth graph topology transformation toward generation quality, thus revealing issues of diffusion-based graph generation models. We follow the same measure of smoothness described in the previous subsection.

Experiment Setup The toy overfitting experiments in the previous experiment does not reflect sample quality well, so in this experiment, we employ a more generalized setting of training on a complete dataset. In the first experiment, we train GraphBFN and DiGress (Vignac et al., 2022) on the complete **SBM** dataset (using the best configuration for both) and let each sample 200 graphs with 1,000 sample steps. We then collect intermediate noise samples and compute the smoothness metrics MLSTD and MVar on the spectral gap as in Sec. G.1. For the final generated samples, we compute the sample quality metrics as in Tab. 1.

Furthermore, we investigate how smoothness in transformation correlates with failure cases diffusion-based approach. We equally partition the 200 samples generated by DiGress into two categories based on the smoothness metric MLSTD, with one group containing the samples with the bottom 100 MLSTD and the other containing those with the top 100.

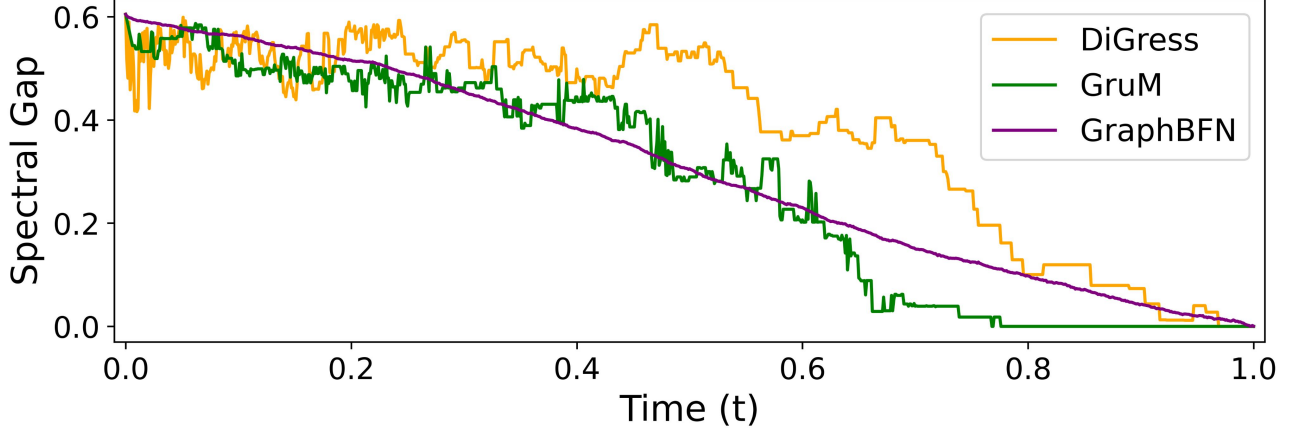


Figure 4: Topology recovery of GraphBFN.

Table 7: The study of the advantages of smoothness topology transformation at generation by comparing GraphBFN and the diffusion-based model DiGress. The smoothness metrics are in the scale of $1e-2$

Model	MLSTD↓	MVar ↓	Deg. ↓	Clus. ↓	Orbit ↓	Spec. ↓	V.U.N. ↑
DiGress	0.79	0.62	0.0015	0.0512	0.0396	0.0379	73.5
GraphBFN	0.23	0.25	0.0004	0.0509	0.0328	0.0050	90.0

Results According to the results of the above experiments, as shown in Tab. 7 and Tab. 8, it is evident that there is a significant performance gap between the sample groups with different smoothness metrics. By selecting the smoother transformation trajectories during generation, the performance of diffusion-based approaches can be significantly improved. This finding better justifies the issues of unstable dynamics in traditional diffusion-based methods.

G.3. Topology Recovery

We conducted a similar analysis as in Jo et al. (2024) on the recovery of graph topology in the generation process. The results of GraphBFN are shown in Figure. 5. Comparing to Figure.2 of Jo et al. (2024), we could find that GraphBFN shares with GruM the advantage of more speedy topology recovery, with GraphBFN recovering the graph topology even faster.

H. Ablation Studies

Output Distribution Conditioning We evaluate the role that the spectral features of graphs play in the generation process. We conduct an ablation study on SBM by training GraphBFN with and without conditioning on the features. When conditioning on the features, we consider extracting the spectral features from two sources, either from the PAM of the GraphBFN’s input distribution or the output distribution. Results in Tab. 9 show that incorporating information on the spectral features is essential for correctly modeling the internal rules of larger graphs like SBM. Besides, the output distribution could be a better source to extract the spectral features as shown by the further boost in performance. This

Table 8: The study of the advantages of smoothness topology transformation at generation through the failure cases of the diffusion-based model DiGress. The smoothness metrics are in the scale of $1e-2$.

Model	MLSTD↓	MVar ↓	Deg. ↓	Clus. ↓	Orbit ↓	Spec. ↓	V.U.N. ↑
DiGress (original)	0.79	0.62	0.0015	0.0512	0.0396	0.0379	73.5
DiGress (highest 100)	1.04	0.76	0.0018	0.0518	0.0461	0.0450	63
DiGress (lowest 100)	0.54	0.48	0.0007	0.0498	0.0336	0.0212	84

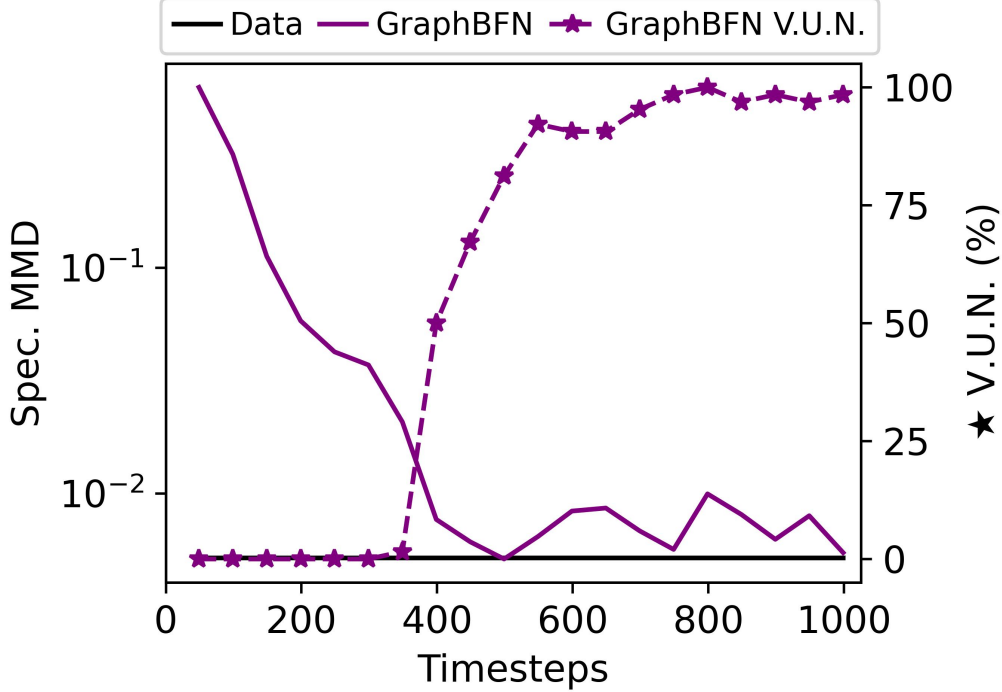


Figure 5: A comparison of the transformation of graph topology between GraphBFN and the diffusion-based models, DiGress and GruM. The curves for GruM, DiGress, and others can be found in Figure 2 of Jo et al. (2024).

Table 9: **Ablation study** on the technique of conditioning GraphBFN’s the training and generation processes on the extra graph features drawing from the PMA of GraphBFN’s *output distribution*. Experiments on SBM and Planar graph datasets show that this technique boosts performance in terms of both generation quality and diversity.

	Planar					SBM				
	Synthetic, $ V = 64$					Synthetic, $44 \leq V \leq 187$				
	Deg. ↓	Clus. ↓	Orbit ↓	Spec. ↓	V.U.N. ↑	Deg. ↓	Clus. ↓	Orbit ↓	Spec. ↓	V.U.N. ↑
Training set	0.0002	0.0310	0.0005	0.0052	100.0	0.0008	0.0332	0.0255	0.0063	100.0
None	0.0647	0.3525	1.695	0.0794	0.0	0.0149	0.0748	0.125	0.0116	0.0
Input Dist.	0.0006	0.0288	0.0013	0.0058	87.5	0.0009	0.058	0.054	0.0057	75.0
Output Dist.	0.0005	0.0294	0.0002	0.0046	96.7	0.0005	0.056	0.037	0.0053	87.5

phenomenon is within our expectation because the input distribution models the distribution of each edge in a graph independently and thus offers less clear spectral information.

Accuracy Scheduling We experimented with how the accuracy scheduling of GraphBFN impacted its generation results. We set up the experiment on SBM, the abstract graph at scale, and trained GraphBFN with both the original quadratic scheduling and our generalized scheduling, with $\beta(1)$ set to 3.0 and α set to 1.0 for the generalized. Results in Tab. 10 show that the generalized scheduling could improve the model’s performance.

I. Derivation of BFN over Discrete Variable

To obtain the exact form of the sender distribution, the key is to apply the central limit theorem and take the condition of $m\omega^2 = \alpha$ and is finite, which implies $m \rightarrow \infty$ and $\omega \rightarrow 0$. Also with the first-order Taylor extension could obtain the final result. For a detailed derivation of the BFN over discrete variables, we would like to suggest referring to the Eq.148 to Eq.157 in (Graves et al., 2023).

Table 10: **Ablation study** on different accuracy scheduling of GraphBFN. According to our experiment on SBM, the scheduling with the general accuracy rate function $\alpha(t) = a + 2(\beta(1) - a)t$ outperforms that with the quadratic rate in the original paper of Graves et al. (2023).

Acc. Sched.	SBM				
	Synthetic, $ V = 187$				
	Deg. ↓	Clus. ↓	Orbit ↓	Spec. ↓	V.U.N. ↓
Training Set	0.0008	0.0332	0.0255	0.0063	100.0
$\alpha(t) = 2\beta(1)t$	0.046	0.0574	0.069	0.0081	75.0
$\alpha(t) = a + 2(\beta(1) - a)t$	0.0005	0.056	0.037	0.0053	87.5

J. Proof of Proposition 4.2

We prove the proposition in the following:

Proof. We first prove the loss in Eq. 12 is permutation invariant. To this end, we first prove that

$$q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0}) = q(\pi \mathbf{G}^{\theta_t} | \pi \mathbf{G}, \pi \mathbf{G}^{\theta_0}) \quad (21)$$

We take the same decomposition as in Eq. 11 of the term $q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0})$:

$$\begin{aligned} q(\mathbf{E}^{\theta_t} | \mathbf{E}, \mathbf{E}^{\theta_0}) &= \prod_{1 \leq i, j \leq n} q(\mathbf{E}_{i,j}^{\theta_t} | \mathbf{E}_{i,j}, \mathbf{E}_{i,j}^{\theta_0}) \quad , \quad q(\mathbf{V}^{\theta_t} | \mathbf{V}, \mathbf{V}^{\theta_0}) = \prod_{1 \leq i \leq n} q(\mathbf{V}_i^{\theta_t} | \mathbf{V}_i, \mathbf{V}_i^{\theta_0}) \\ q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0}) &= q(\mathbf{E}^{\theta_t} | \mathbf{E}, \mathbf{E}^{\theta_0}) q(\mathbf{V}^{\theta_t} | \mathbf{V}, \mathbf{V}^{\theta_0}) \end{aligned} \quad (22)$$

We take the $q(\mathbf{V}^{\theta_t} | \mathbf{V}, \mathbf{V}^{\theta_0})$ as example:

$$\begin{aligned} q(\pi \mathbf{V}^{\theta_t} | \pi \mathbf{V}, \pi \mathbf{V}^{\theta_0}) &= \prod_{1 \leq i \leq n} q(\mathbf{V}_i^{\theta_t} | \mathbf{V}_i, \mathbf{V}_i^{\theta_0}) = \prod_{1 \leq i \leq n} q(\mathbf{V}_i^{\theta_t} | \mathbf{V}_i, \mathbf{V}_i^{\theta_0}) \\ &= q(\mathbf{V}^{\theta_t} | \mathbf{V}, \mathbf{V}^{\theta_0}) \end{aligned} \quad (23)$$

Similar derivation could be applied to \mathbf{E} . And hence we could finish the proof.

Next, we take the $L(\mathbf{V}, t)$ as an example to prove $L(\mathbf{V}, t) = L(\pi \mathbf{V}, t)$:

$$L(\pi \mathbf{V}, t) = c_v \frac{1}{2} \alpha_v(t) \mathbb{E}_{q(\pi \mathbf{G}^{\theta_t} | \pi \mathbf{G}, \pi \mathbf{G}^{\theta_0})} \sum_{1 \leq i \leq n} \|\pi \mathbf{V}_i - \mathbf{V}_i [\phi(\pi \mathbf{G}^{\theta_t}, t)]\|^2 \quad (24)$$

For simplicity, we focus on the term related to π , as ϕ is permutation equivariant, we have

$$\begin{aligned} &\mathbb{E}_{q(\pi \mathbf{G}^{\theta_t} | \pi \mathbf{G}, \pi \mathbf{G}^{\theta_0})} \sum_{1 \leq i \leq n} \|\pi \mathbf{V}_i - \mathbf{V}_i [\phi(\pi \mathbf{G}^{\theta_t}, t)]\|^2 = \mathbb{E}_{q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0})} \sum_{1 \leq i \leq n} \|\pi \mathbf{V}_i - \pi \mathbf{V}_i [\phi(\mathbf{G}^{\theta_t}, t)]\|^2 \\ &= \mathbb{E}_{q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0})} \pi \left(\sum_{1 \leq i \leq n} \|\mathbf{V}_i - \mathbf{V}_i [\phi(\mathbf{G}^{\theta_t}, t)]\|^2 \right) = \mathbb{E}_{q(\mathbf{G}^{\theta_t} | \mathbf{G}, \mathbf{G}^{\theta_0})} \sum_{1 \leq i \leq n} \|\mathbf{V}_i - \mathbf{V}_i [\phi(\mathbf{G}^{\theta_t}, t)]\|^2 \end{aligned} \quad (25)$$

Then we could also apply the same process to \mathbf{E} . Then we finish the proof over permutation invariance of the loss function.

Next, we discuss the invariant density modeling of the p_ϕ . This could be directly following the conclusion of Theorem 3.1 of Song et al. (2023). With a permutation invariant prior \mathbf{G}^{θ_0} and a permutation equivariant interdependency network ϕ , we only need to make sure the Bayesian Update function h is permutation equivariant. As we could find in Eq. 7,

$$h(\pi \theta_{t-1}, \pi y_t) = \pi \left(\frac{e^{y_t \theta_{t-1}}}{\sum_{k=1}^K e^{y_t^{(k)} \theta_{t-1}^{(k)}}} \right) = \pi h(\theta_{t-1}, y_t) \quad (26)$$

And then we could finish the proof. \square

K. Comparison with Diffusion-Based Approaches

K.1. Comparison with Discrete Diffusion Methods

The key intuition held by Bayesian Flow Networks(BFN) (Graves et al., 2023) is that the information in the generative process should be modeled in a smoothly increasing fashion for superior performance. Such intuition is shared by the BFN and diffusion in the continuous data and is also identified by several studies under the literature of diffusion models for image generations (Ho et al., 2020). GraphBFN is motivated by the compatibility of such "continuous information modeling" intuition and the special nature of graph data: that is, though generally represented in the discrete modality, the essential graph topology information could also be represented in the continuous manifold as implied by graph spectrum. This allows GraphBFN to better harness the inductive bias in modeling discrete data modality with continuous natures, compared to the previous discrete diffusion models (Vignac et al., 2022). Furthermore, compared to the methods that directly dequantized the data and applied continuous diffusion models, the GraphBFN needs no quantize/dequantize procedure and could be more compatible with the sparse adjacency matrix representation.

There are also some potential empirical advantages of GraphBFN. The input to the neural network of GraphBFN lies in the continuous manifold, which could enjoy the advantages of numerical stability. More importantly, the inputs to the continuous distribution hold a considerably small variance compared to the previous discrete diffusion models (Graves et al., 2023). Also, such properties could lead to a speed-up in convergence during training.

K.2. Comparison with Diffusion Methods modeled within a Bounded Space

Beta diffusion models (Zhou et al., 2023; Liu et al., 2025) present a compelling alternative to the widely used Gaussian diffusion models (Ho et al., 2020; Jo et al., 2022). Like the Gaussian counterparts, it admits analytical forms for the marginal and conditional distributions of the diffusion process. However, it offers the added advantage of modeling variables within a bounded domain between 0 and 1, which Gaussian diffusion does not naturally support. This boundedness is particularly beneficial in discrete settings such as graph generation, where the data is inherently constrained. Therefore, Graph Beta Diffusion (GBD) (Liu et al., 2025) is similar to GraphBFN in the sense that they both preserve the constrained nature of discrete graphs while modeling the evolution of continuous quantities.

Nevertheless, we would like to highlight the key similarities and distinctions between GraphBFN and GBD.

In the case of abstract graph generation, where there are only two edge types (has edge & no edge), it suffices to represent each entry of adjacent matrixes with a single variable within the range [0, 1]. In this simple setting, the variables modeled by GBD could be interpreted as the probability of edge existence. Therefore, both GBD and GraphBFN share the similar property of modeling the smooth interpolation of probabilistic matrices, although they differ in specific probability trajectories. As the results in Table 1 indicate, both GBD and GraphBFN could obtain a superior performance compared to the discrete-state-space models (e.g., DiGress (Vignac et al., 2022)) in abstract graph settings, which justifies the advantages of smooth interpolation, as we claimed.

However, for graphs with multiple edge types (e.g., molecular graphs), the modeling of GraphBFN and GBD have non-negligible differences. When the adjacency matrix contains k edge types, GraphBFN represents each matrix entry as a k -dimensional probability vector that lies on the simplex (i.e., the vector components sum to 1), making each dimension interpretable as the probability of a particular edge type. In contrast, GBD models each dimension independently as a value in [0, 1], without enforcing normalization. The trajectory over simplexes of GraphBFN makes it a more natural choice to model multiple-class discrete variables by making different dimensions contrastive. As the results in Table 2 indicate, GraphBFN achieves notably stronger performance in this setting, which justifies the potential of simplex-based approaches in modeling multi-class discrete graphs.

L. Experiment Details and Additional Results

L.1. Experiment Details

The datasets, dataset splits, and evaluation metrics of QM9 with explicit hydrogen and MOSES follow from Vignac et al. (2022), those of SBM, Planar, QM9 with implicit hydrogen, and ZINC250k follow from Martinkus et al. (2022), and those of Community-small, Ego-small, and Protein follow from Jo et al. (2022).

The network architectures and settings directly followed Vignac et al. (2022); Jo et al. (2024). The models are trained with

the AdamW optimizer (Loshchilov & Hutter, 2019) until convergence. The most sensitive hyperparameters of GraphBFN are β_E and β_V that control the speed of information transmission in the generation process. We search them in the discrete linear space $\{1.0, 2.0, \dots, 10.0\}$.

For all of our experiments, we use NVIDIA GeForce RTX 3090 with 24GB memory to train and evaluate our models. The batch sizes are adjusted to the maximum that can fit within the memory constraint.

All of our reported results are the average of three random runs. For a cleaner presentation and to follow the convention of previous work, we only add the error terms to the result of **QM9 with Explicit Hydrogen**.

For the results of **QM9 with Explicit Hydrogen** in Tab. 4, we only compare to one other method because, to our best knowledge, DiGress (Vignac et al., 2022) is the only previous graph generative model that has attempted to model QM9 molecules with *explicit* hydrogen.

Extra Features: We use extra features for the inputs of the GraphTransformer which directly follows (Vignac et al., 2022). Specifically, we have included cycle counts using pre-determined formulas to efficiently calculate cycles up to size 6. Node features count participation in up to 5-cycles, while graph features tally up to 6-cycles, utilizing node degree vectors and the Frobenius norm. Additionally, we integrate spectral features, which involve an $O(n^3)$ eigendecomposition—manageable for our graphs with up to 200 nodes. Key graph-level features derived from the graph Laplacian include the number of connected components and the first five nonzero eigenvalues. Node-level features estimate the largest connected component and include the first two eigenvectors associated with nonzero eigenvalues.

Each spectral feature or structural feature could be categorized into node-level features and graph-level features. Specifically, a node-level feature holds the shape of (batch-of-graphs, node-in-graph), such as an estimation of the biggest connected component (using the eigenvectors associated with eigenvalue 0) ; a graph-level feature holds the shape of (batch-of-graphs, 1). After computing the node-level features F_n and graph-level features F_g based on V^θ and E^θ , we concat V^θ and all F_n as new node representations V_{feat} , and concat all F_g as the finally graph-level feature y . Then the inputs for the graph transformers in GraphBFN are the triple of $(V_{\text{feat}}, E^\theta, y)$ similarly to DiGress (Vignac et al., 2022)

Scheduling Function: In short, our scheduling design in Eq. 14 is an improvement upon the original design in vanilla BFN (Graves et al., 2023) that is determined empirically based on the intuitive of linearly increased entropy. Our design attempts to change the non-informative prior to the original design to be informative.

In Eq. 14, $\beta(t)$ is the accumulative(integral of $\alpha(t)$ from 0 to t). The key intuition behind the design is that for informative prior, the initial step 0 already contains information of the ground truth sample. Then we set $\alpha(0)$ to a positive constant a . Besides, we also follow the heuristic of vanilla BFN (Graves et al., 2023) to use the square scheduler for accumulated accuracy β to ensure the entropy of the input distribution approximately changes linearly along the timesteps. Eq. 14 is a very simple form that satisfies both properties.

L.2. Computation Efficiency

Training time for QM9 takes approximately 3 hours, and for large molecules and abstract datasets, the longest training takes 4 days.

For the inference efficiency, we conducted an experiment that compares GraphBFN with two other multi-step graph generation models – DiGress (Vignac et al., 2022) and GruM (Jo et al., 2024). We record the time for each model to sample 10000 graphs of QM9 with explicit hydrogen with the same sampling steps of 500. The results are shown in the Tab. 11.

According to the results, GraphBFN has approximately the same inference speed as the other two diffusion-based models. This result is expected, as all three methods require only one network evaluation at each iteration. GruM enjoys the fastest sampling speed. We speculate that this is because GruM operates in the space of continuous data, so it does not require the relatively costly procedure of sampling from categorical distributions.

Table 11: Comparing GraphBFN with other multistep generation methods on the time to Sample 10,000 QM9 with *Explicit* Hydrogen Graphs.

Sampling Time (min.)	
DiGress	40
GruM	35
GraphBFN	44

Table 12: Additional results on smaller datasets. Most baseline performances are adapted from (Jo et al., 2022) and (Kong et al., 2023). Hyphens(-) denote unobtainable results due to either the reasons mentioned in (Jo et al., 2022) or the lack of official implementation.

		Ego-small			Community-small			Enzymes		
		Real, $4 \leq V \leq 18$			Synthetic, $12 \leq V \leq 20$			Real, $10 \leq V \leq 125$		
		Deg.↓	Clus.↓	Orbit↓	Deg.↓	Clus.↓	Orbit↓	Deg.↓	Clus.↓	Orbit↓
	Training set	0.002	0.021	0.005	0.02	0.07	0.01	0.0003	0.01079	0.0003
Autoreg.	DeepGMG	0.040	0.100	0.020	0.220	0.950	0.400	-	-	-
	GraphRNN	0.090	0.220	0.003	0.080	0.120	0.040	0.017	0.062	0.046
	GraphAF	0.03	0.11	0.001	0.18	0.20	0.02	1.669	1.283	0.266
	GraphDF	0.04	0.13	0.01	0.06	0.12	0.03	1.503	1.061	0.202
One-shot	GraphVAE	0.130	0.170	0.050	0.350	0.980	0.540	1.369	0.629	0.191
	GNF	0.030	0.100	0.001	0.200	0.200	0.110	-	-	-
	EDP-GNN	0.052	0.093	0.007	0.053	0.144	0.026	0.023	0.268	0.082
	GDSS	0.021	0.024	0.007	0.045	0.086	0.007	0.026	0.061	0.009
	SPECTRE	0.078	0.078	0.007	0.020	0.210	0.010	0.136	0.195	0.125
	DiGress	0.015	0.029	0.005	0.020	0.063	0.010	0.004	0.083	0.002
	EDGE	-	-	-	0.008	0.011	0.026	-	-	-
	GraphARM	0.019	0.017	0.010	0.034	0.082	0.004	0.029	0.054	0.015
	HiGen	-	-	-	-	-	-	0.012	0.038	7.2e-4
	GraphBFN	0.0034	0.0246	0.0057	0.0075	0.0095	0.0043	0.0064	0.0212	0.0051

L.3. Additional Results on Non-attributed Graphs

In Tab. 12, we provide results on non-attributed (i.e. graph with only one edge type and node type) graph datasets, along with comparisons to additional baselines (Chen et al., 2023; Kong et al., 2023; Karami, 2024). **Community-small** contains 100 randomly generated clustered graphs; **Ego-small** contains 200 small ego graphs; and **Enzymes** contains 587 chain-like graphs representing protein structures (Jo et al., 2022).

According to the table, GraphBFN has shown superior performance on Community-small and competitive performance on Ego-small and Enzymes. It should be noted that for **Community-small** and **Ego-small**, the metrics are very saturated due to the limited complexity of datasets; as we can see, the SOTA results are very close to the training set’s performance.

L.4. Code

More specific details about the experiment setting and our implementation of GraphBFN can be found in our code at <https://anonymous.4open.science/r/GraphBFN-5FEE/README.md>.

M. Flowback Sampling

The visualization of GraphBFN’s generation process with the Flowback technique is in Figure. 6.

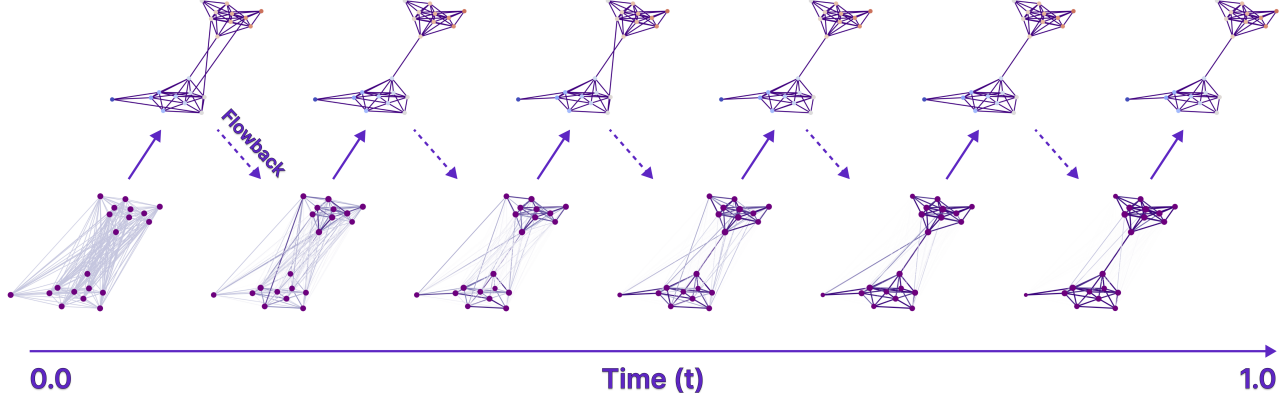


Figure 6: Visualization of GraphBFN’s generation process with the Flowback technique. The sample graphs in the first row are sampled from the output distribution $\phi(\mathbf{G}^{\theta_t}, t)$ at different time steps, while those on the second line are representations of the probabilistic adjacency matrix of the input distribution \mathbf{G}^{θ_t} . The solid lines denote the network prediction, while the solid line denotes the flowback Bayesian update.

N. Derivation from Eq. 7 to Eq. 10

Here we provide a formal derivation of the Bayesian Flow Distribution (Eq. 10) which is essentially the Eq. 172 to Eq. 185 in BFN (Graves et al., 2023). Firstly, we consider marginalizing the y from the sender distribution with the Bayesian update function in Eq. 7 which results in a distribution over θ , we term it as Bayesian update distribution:

$$p_U(\theta \mid \theta_{i-1}, \mathbf{x}; \alpha) = \mathcal{N}(\mathbf{y} \mid \alpha(K\mathbf{e}_x - \mathbf{1}), \alpha K\mathbf{I}) \delta\left(\theta - \frac{e^{\mathbf{y}\theta_{i-1}}}{\sum_{k=1}^K e^{\mathbf{y}_k (\theta_{i-1})_k}}\right). \quad (27)$$

Then we derive the additive accuracy property with two consecutive Bayesian updates of noisy samples $y_a \sim \mathcal{N}(\alpha_a(K\mathbf{e}_x - \mathbf{1}), \alpha_a K\mathbf{I})$ and $y_b \sim \mathcal{N}(\alpha_b(K\mathbf{e}_x - \mathbf{1}), \alpha_b K\mathbf{I})$ as:

$$\begin{aligned} h(y_b, h(y_a, \theta_{i-2})) &= \frac{\exp(y_b) \frac{\exp(y_a)\theta_{i-2}}{\sum_{k'=1}^K \exp((y_a)_{k'}) (\theta_{i-2})_{k'}}}{\sum_{k=1}^K \exp((y_b)_k) \frac{\exp((y_a)_{k'}) (\theta_{i-2})_{k'}}{\sum_{k'=1}^K \exp((y_a)_{k'}) (\theta_{i-2})_{k'}}} \\ &= \frac{\exp(y_b) \exp(y_a) \theta_{i-2}}{\sum_{k=1}^K \exp((y_b)_k) \exp((y_a)_k) (\theta_{i-2})_k} \\ &= \frac{\exp(y_a + y_b) \theta_{i-2}}{\sum_{k=1}^K \exp((y_a + y_b)_k) (\theta_{i-2})_k} \\ &= h(y_a + y_b, \theta_{i-2}). \end{aligned} \quad (28)$$

With the property of Gaussian distribution, we have that:

$$y_a + y_b \sim \mathcal{N}((\alpha_a + \alpha_b)(K\mathbf{e}_x - \mathbf{1}), (\alpha_a + \alpha_b) K\mathbf{I}). \quad (29)$$

Then the two-step update could be actually merged into a one-step update of added accuracies, i.e.:

$$p_U(\theta_{i-1} \mid \theta_{i-2}, \mathbf{x}; \alpha_a) p_U(\theta_i \mid \theta_{i-1}, \mathbf{x}; \alpha_b) = p_U(\theta_i \mid \theta_{i-2}, \mathbf{x}; \alpha_a + \alpha_b) \quad (30)$$

The the Bayesian flow distribution in Eq. 10:

$$p_F(\boldsymbol{\theta} \mid \mathbf{x}; t) = p_U(\boldsymbol{\theta}_1 \mid \boldsymbol{\theta}_0, \mathbf{x}; \alpha_1) \mathbb{E}_{p_U(\boldsymbol{\theta}_2 \mid \boldsymbol{\theta}_1, \mathbf{x}; \alpha_2)} \cdots \mathbb{E}_{p_U(\boldsymbol{\theta}_{n-1} \mid \boldsymbol{\theta}_{n-2}, \mathbf{x}; \alpha_{n-1})} p_U(\boldsymbol{\theta}_n \mid \boldsymbol{\theta}_{n-1}, \mathbf{x}; \alpha_n) = p_U\left(\boldsymbol{\theta}_n \mid \boldsymbol{\theta}_0, \mathbf{x}; \sum_{i=1}^n \alpha_i\right) \quad (31)$$

Note $\sum_{i=1}^n \alpha_i$ could be substituted as $\beta(t)$ as the integral of the $\alpha(t)$ from 0 to t in the continuous setting. Substituting Eq. 27 into Eq. 31, we could obtain that:

$$p_F(\boldsymbol{\theta} \mid \mathbf{x}; t) = \mathbb{E}_{\mathcal{N}(\mathbf{y} \mid \beta(t)(K\mathbf{e}_x - \mathbf{1}), \beta(t)KI)} \delta\left(\boldsymbol{\theta} - \frac{e^{\mathbf{y}\boldsymbol{\theta}_0}}{\sum_{k=1}^K e^{\mathbf{y}^k} (\boldsymbol{\theta}_0)_k}\right). \quad (32)$$

Changing the notation from y, θ to $\mathbf{E}_{i,j}^y$ and $\mathbf{E}_{i,j}^\theta$ we could finish the derivation.

O. Further Discussion over Adaptive Flowback

The generative process of BFN all start with the same initial points, e.g., vector $[\frac{1}{k}, \dots, \frac{1}{k}]$ for k -class discrete variable, and ends with various generated samples. Hence, the variance of the input distribution first grows and then decreases [2], which could be intuitively understood as first exploring the endpoint of the generative process and then finetuning towards that endpoint. Motivated by this fact, we suggest that when the generative endpoints are approximately determined, it could be more favorable to eliminate the extra variance in θ from sequential Bayesian Update and only focus on finetuning θ based on the Bayesian flow of the current step output prediction. The condition of $\|\phi(\mathbf{G}^{\theta_t}, t) - \phi(\mathbf{G}^{\theta_{t-1}}, t-1)\|^2 \geq \epsilon$ is to test whether the predicted outputs of two consecutive steps are close enough, i.e. whether the steps t and steps $t-1$ are generated towards the same targets. We use this as an indicator for whether the generation process is in a variance increase (exploration) or a variance decrease (finetuning) state.

P. Continuous-state Approaches or Discrete-state Approaches for Graph Generation

There have been several works making notable progress in continuous state-based methods, specifically spectrum-based approaches such as (Martinkus et al., 2022; Luo et al., 2023). More recently, several discrete state-based approaches that directly model the discrete variable of edges/nodes, e.g. (Vignac et al., 2022; Jo et al., 2024), have shown that combining the continuous state (e.g., spectrum) as an extra feature, the discrete state approaches could further achieve better results.

There is the particularity of graph data, whose form is discrete and important property enjoys continuity. And current empirical evidence (Vignac et al., 2022; Jo et al., 2024) suggests it could be favorable to take both into consideration. The motivation for choosing BFN lies in the fact that it could actually model the discrete variable in continuous categorical parameter space, which perfectly fits the discrete form and continuous property of graph data.

Additionally, though GraphBFN lies in the category of discrete state approaches, we would like to clarify that there is still a lot of research in the field to be conducted to finally answer the discrete or continuous question.

Q. Relationship to Spectral Graph Theory

Here we hope to provide a more detailed discussion of the relationship of proposed methods to spectral graph theory. The framework aims to introduce a new generative method that forms in the discrete state space and takes the intrinsic continuity of the topological information into consideration. The spectral graph theory provides an effective tool for us to implement and analyze the proposed approach.

From the implementation perspective, we use the spectral features of the input to enhance the network prediction following (Vignac et al., 2022). During both training and sampling, the spectral features of the parameters of the graphs are computed and directly appended to the network as additional inputs.

Additionally, we use spectral theory to help us analyze one key reason behind our method’s performance improvement compared to previous approaches: the smooth transformation of graph topology leads to enhanced sample quality. As demonstrated in the ablation study in Appendix F, we use clustered graphs as an example to show that, for both diffusion-based methods and GraphBFN, a smoother transition of the sample’s spectral features—such as the spectral gap in the case of clustered graphs—during the generation process significantly improves sample quality.

R. Visualization of Generated Samples

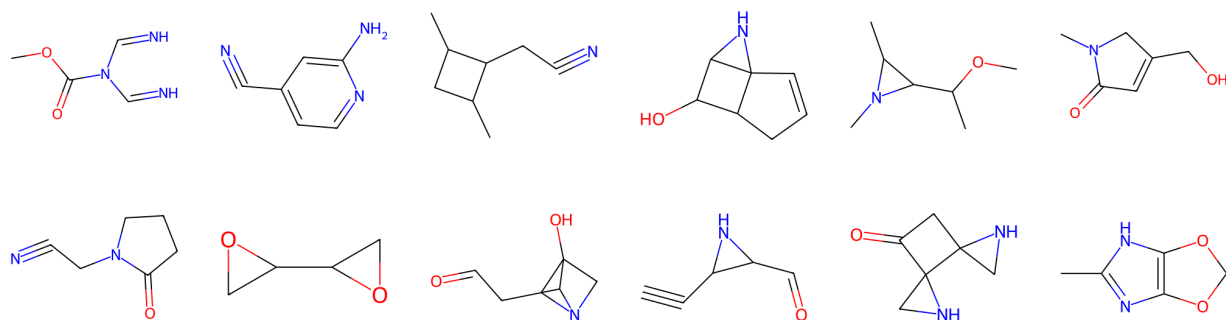


Figure 7: *Explicit* Hydrogens

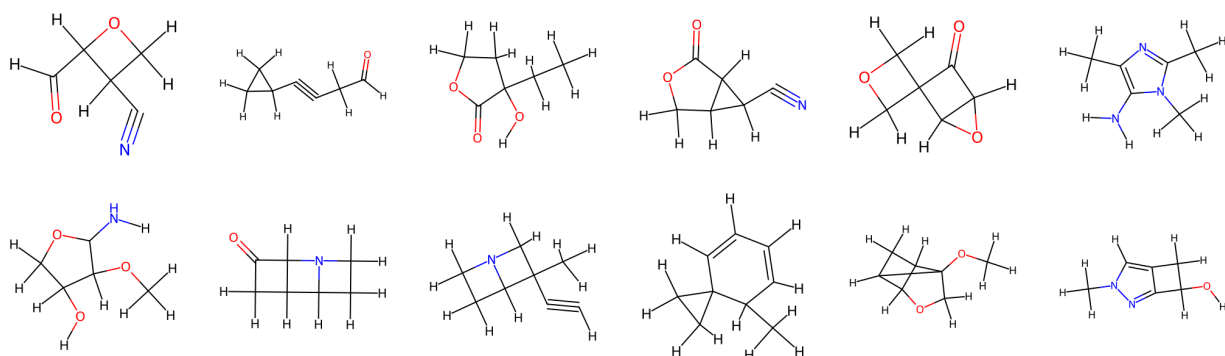


Figure 8: *Implicit* Hydrogens

Figure 9: Unfiltered samples generated by GraphBFN, trained on QM9

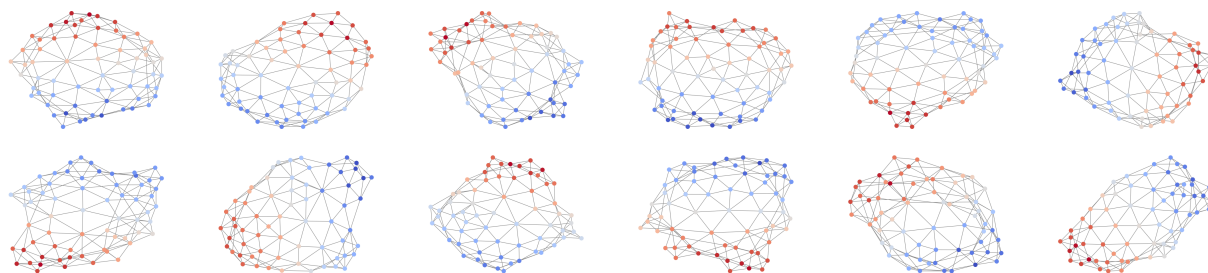


Figure 10: *Planar*

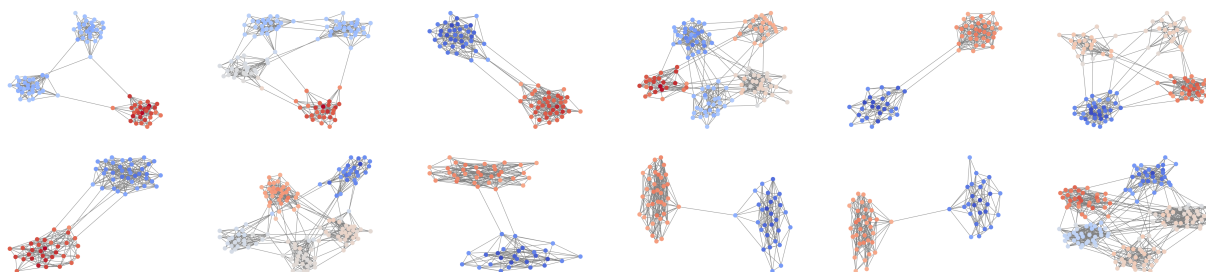


Figure 11: *SBM*

Figure 12: Unfiltered samples generated by GraphBFN, trained on abstract graph datasets

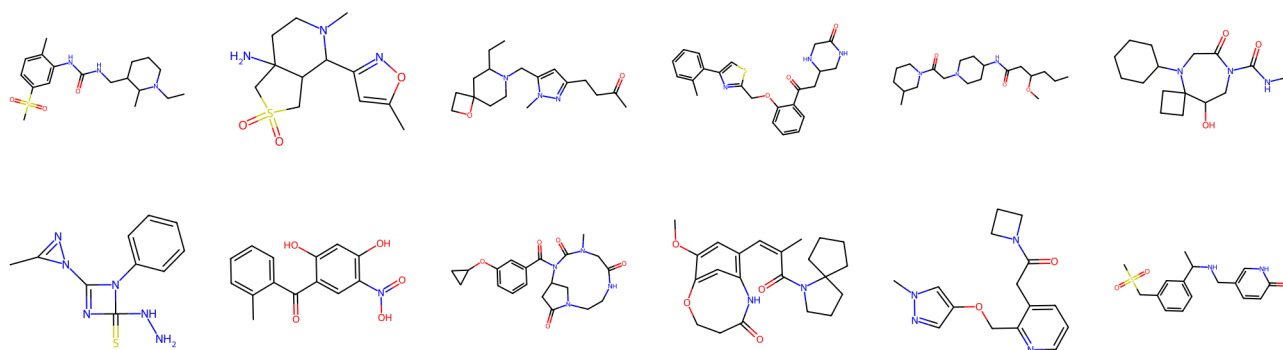


Figure 13: Unfiltered samples generated by GraphBFN, trained on ZINC250k with *implicit* hydrogens

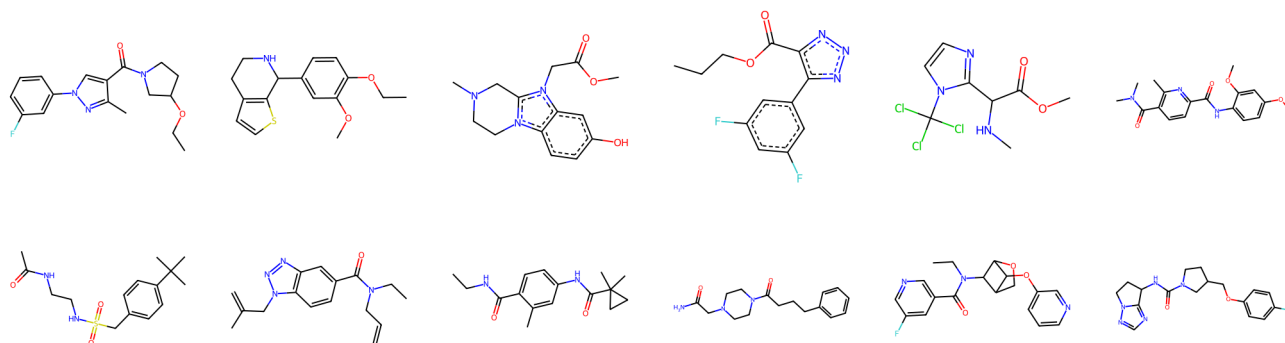


Figure 14: Unfiltered samples generated by GraphBFN, trained on MOSES graphs